

COMPUTER SECURITY HANDBOOK

COMPUTER SECURITY HANDBOOK

Sixth Edition

Volume 1

Edited by

SEYMOUR BOSWORTH

MICHEL E. KABAY

ERIC WHYNE

WILEY

Cover image: ©iStockphoto.com/Jimmy Anderson
Cover design: Wiley

Copyright © 2014 by John Wiley & Sons, Inc. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey.

Previous Edition: Computer Security Handbook, Fifth Edition. Copyright © 2009 by John Wiley & Sons, Inc.
All Rights Reserved. Published by John Wiley & Sons, Inc., Hoboken, New Jersey.

Published simultaneously in Canada.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 646-8600, or on the Web at www.copyright.com. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permissions>.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services or for technical support, please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Wiley publishes in a variety of print and electronic formats and by print-on-demand. Some material included with standard print versions of this book may not be included in e-books or in print-on-demand. If this book refers to media such as a CD or DVD that is not included in the version you purchased, you may download this material at <http://booksupport.wiley.com>. For more information about Wiley products, visit www.wiley.com.

Library of Congress Cataloging-in-Publication Data

Computer security handbook / [edited by] Seymour Bosworth, Michel E. Kabay,
Eric Whyne. – Sixth edition.

volumes cm

Includes index.

ISBN 978-1-118-13410-8 (vol. 1 : pbk.) – ISBN 978-1-118-13411-5 (vol. 2 : pbk.) –

ISBN 978-1-118-12706-3 (2 volume set : pbk.); ISBN 978-1-118-85174-6 (ebk);

ISBN 978-1-118-85179-1 (ebk) 1. Electronic data processing departments—Security measures.

I. Bosworth, Seymour. II. Kabay, Michel E. III. Whyne, Eric, 1981–

HF5548.37.C64 2014

658.4'78–dc23

2013041083

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

CONTENTS

PREFACE

ACKNOWLEDGMENTS

ABOUT THE EDITORS

ABOUT THE CONTRIBUTORS

A NOTE TO THE INSTRUCTOR

PART I FOUNDATIONS OF COMPUTER SECURITY

- 1. Brief History and Mission of Information System Security**
Seymour Bosworth and Robert V. Jacobson
- 2. History of Computer Crime**
M. E. Kabay
- 3. Toward a New Framework for Information Security**
Donn B. Parker, CISSP
- 4. Hardware Elements of Security**
Sy Bosworth and Stephen Cobb
- 5. Data Communications and Information Security**
Raymond Panko and Eric Fisher
- 6. Local Area Network Topologies, Protocols, and Design**
Gary C. Kessler
- 7. Encryption**
Stephen Cobb and Corinne LeFran ois
- 8. Using a Common Language for Computer Security Incident Information**
John D. Howard

vi CONTENTS

9. Mathematical Models of Computer Security
Matt Bishop
10. Understanding Studies and Surveys of Computer Crime
M. E. Kabay
11. Fundamentals of Intellectual Property Law
William A. Zucker and Scott J. Nathan

PART II THREATS AND VULNERABILITIES

12. The Psychology of Computer Criminals
Q. Campbell and David M. Kennedy
13. The Insider Threat
Gary L. Tagg, CISSP
14. Information Warfare
Seymour Bosworth
15. Penetrating Computer Systems and Networks
Chey Cobb, Stephen Cobb, M. E. Kabay, and Tim Crothers
16. Malicious Code
Robert Guess and Eric Salveggio
17. Mobile Code
Robert Gezelter
18. Denial-of-Service Attacks
Gary C. Kessler
19. Social-Engineering and Low-Tech Attacks
Karthik Raman, Susan Baumes, Kevin Beets, and Carl Ness
20. Spam, Phishing, and Trojans: Attacks Meant to Fool
Stephen Cobb
21. Web-Based Vulnerabilities
Anup K. Ghosh, Kurt Baumgarten, Jennifer Hadley, and Steven Lovaas
22. Physical Threats to the Information Infrastructure
Franklin Platt

PART III PREVENTION: TECHNICAL DEFENSES

23. Protecting the Physical Information Infrastructure
Franklin Platt

CONTENTS vii

- 24. Operating System Security**
William Stallings
- 25. Local Area Networks**
N. Todd Pritsky, Joseph R. Bumblis, and Gary C. Kessler
- 26. Gateway Security Devices**
Justin Opatrny
- 27. Intrusion Detection and Intrusion Prevention Devices**
Rebecca Gurley Bace
- 28. Identification and Authentication**
Ravi Sandhu, Jennifer Hadley, Steven Lovaas, and Nicholas Takacs
- 29. Biometric Authentication**
Eric Salveggio, Steven Lovaas, David R. Lease, and Robert Guess
- 30. E-Commerce and Web Server Safeguards**
Robert Gezelter
- 31. Web Monitoring and Content Filtering**
Steven Lovaas
- 32. Virtual Private Networks and Secure Remote Access**
Justin Opatrny and Carl Ness
- 33. 802.11 Wireless LAN Security**
Gary L. Tagg, CISSP and Jason Sinchak, CISSP
- 34. Securing VoIP**
Christopher Dantos and John Mason
- 35. Securing P2P, IM, SMS, and Collaboration Tools**
Carl Ness
- 36. Securing Stored Data**
David J. Johnson, Nicholas Takacs, Jennifer Hadley, and M. E. Kabay
- 37. PKI and Certificate Authorities**
Santosh Chokhani, Padgett Peterson, and Steven Lovaas
- 38. Writing Secure Code**
Lester E. Nichols, M. E. Kabay, and Timothy Braithwaite
- 39. Software Development and Quality Assurance**
Diane E. Levine, John Mason, and Jennifer Hadley
- 40. Managing Software Patches and Vulnerabilities**
Karen Scarfone, Peter Mell, and Murugiah Souppaya

viii CONTENTS

- 41. Antivirus Technology**
Chey Cobb and Allysa Myers
- 42. Protecting Digital Rights: Technical Approaches**
Robert Guess, Jennifer Hadley, Steven Lovaas, and Diane E. Levine

PART IV PREVENTION: HUMAN FACTORS

- 43. Ethical Decision Making and High Technology**
James Landon Linderman
- 44. Security Policy Guidelines**
M. E. Kabay and Bridgitt Robertson
- 45. Employment Practices and Policies**
M. E. Kabay and Bridgitt Robertson
- 46. Vulnerability Assessment**
Rebecca Gurley Bace and Jason Sinchak
- 47. Operations Security and Production Controls**
M. E. Kabay, Don Holden, and Myles Walsh
- 48. Email and Internet Use Policies**
M. E. Kabay and Nicholas Takacs
- 49. Implementing a Security-Awareness Program**
K. Rudolph
- 50. Using Social Psychology to Implement Security Policies**
M. E. Kabay, Bridgitt Robertson, Mani Akella, and D. T. Lang

- 51. Security Standards for Products**
Paul Brusil and Noel Zakin

PART V DETECTING SECURITY BREACHES

- 52. Application Controls**
Myles Walsh and Susan Baumes
- 53. Monitoring and Control Systems**
Caleb S. Coggins and Diane E. Levine
- 54. Security Audits**
Donald Glass, Richard O. Moore III, Chris Davis, John Mason, David Gursky, James Thomas, Wendy Carr, M. E. Kabay, and Diane Levine
- 55. Cyber Investigation**
Peter Stephenson

CONTENTS ix

PART VI RESPONSE AND REMEDIATION

- 56. Computer Security Incident Response Teams**
Michael Miora, M. E. Kabay, and Bernie Cowens
- 57. Data Backups and Archives**
M. E. Kabay and Don Holden
- 58. Business Continuity Planning**
Michael Miora
- 59. Disaster Recovery**
Michael Miora
- 60. Insurance Relief**
Robert A. Parisi, Jr., John F. Mullen, and Kevin Apollo
- 61. Working with Law Enforcement**
David A. Land

PART VII MANAGEMENT'S ROLE IN SECURITY

- 62. Quantitative Risk Assessment and Risk Management**
Robert V. Jacobson and Susan Baumes
- 63. Management Responsibilities and Liabilities**
Carl Hallberg, M. E. Kabay, Bridgett Robertson, and Arthur E. Hutt
- 64. U.S. Legal and Regulatory Security Issues**
Timothy Virtue
- 65. The Role of the CISO**
Karen F. Worstell
- 66. Developing Security Policies**
M. E. Kabay and Sean Kelley
- 67. Developing Classification Policies for Data**
Karthik Raman, Kevin Beets, and M. E. Kabay
- 68. Outsourcing and Security**
Kip Boyle, Michael Buglewicz, and Steven Lovaas

PART VIII PUBLIC POLICY AND OTHER CONSIDERATIONS

- 69. Privacy in Cyberspace: U.S. and European Perspectives**
Henry L. Judy, Scott L. David, Benjamin S. Hayes, Jeffrey B. Ritter, Marc Rotenberg, and M. E. Kabay

x CONTENTS

- 70. Anonymity and Identity in Cyberspace**
M. E. Kabay, Eric Salveggio, Robert Guess, and Russell D. Rosco
- 71. Healthcare Security and Privacy**
Paul Brusil
- 72. Legal and Policy Issues of Censorship and Content Filtering**
Lee Tien, Seth Finkelstein, and Steven Lovaas
- 73. Expert Witnesses and the *Daubert* Challenge**
Chey Cobb
- 74. Professional Certification and Training in Information Assurance**
M. E. Kabay, Christopher Christian, Kevin Henry, and Sondra Schneider
- 75. The Future of Information Assurance**
Jeremy A. Hansen

PREFACE

Computers are an integral part of our economic, social, professional, governmental, and military infrastructures. They have become necessities in virtually every area of modern life, but their vulnerability is of increasing concern. Computer-based systems are constantly under threats of inadvertent error and acts of nature, as well as those attributable to unethical, immoral, and criminal activities. It is the purpose of *The Computer Security Handbook* to provide guidance in recognizing these threats, eliminating them where possible and, if not, then reducing any losses attributable to them.

The Handbook will be most valuable to those directly responsible for computer, network, or information security, as well as those who must design, install, and maintain secure systems. It will be equally important to those managers whose operating functions can be affected by breaches in security and to those executives who are responsible for protecting the assets that have been entrusted to them.

With the advent of desktop, laptop, and handheld computers, and with the vast international networks that interconnect them, the nature and extent of threats to computer security have grown almost beyond measure. In order to encompass this unprecedented expansion, *The Computer Security Handbook* has grown apace.

When the first edition of the *Handbook* was published, its entire focus was on mainframe computers, the only type then in widespread use. The second edition recognized the advent of small computers, while the third edition placed increased emphasis on PCs and networks.

Edition	Publication Date	Chapters	Text Pages
First	1973	12	162
Second	1988	19	383
Third	1995	23	571
Fourth	2002	54	1,184
Fifth	2009	77	2,040
Sixth	2014	75	2,224

The fourth edition of *The Computer Security Handbook* gave almost equal attention to mainframes and microcomputers, requiring more than twice the number of chapters and pages as the third.

xii PREFACE

The fifth edition was as great a step forward as the fourth. With 77 chapters and the work of 86 authors, we increased coverage in both breadth and depth. In this sixth edition, we updated all chapters while continuing to cover all 10 domains of the Common Body of Knowledge, as defined by the International Information Systems Security Certification Consortium (ISC)²:

- 1.** Security Management Practices: Chapters 10, 12, 13, 14, 15, 19, 10, 31, 43, 44, 45, 46, 47, 48, 49, 50, 51, 54, 55, 62, 63, 64, 65, 66, 67, 68, 74, 75
- 2.** Security Architecture and Models: Chapters 1, 2, 3, 8, 9, 24, 26, 27, 51
- 3.** Access Control Systems and Methodology: Chapters 15, 19, 28, 29, 32
- 4.** Application Development Security: Chapters 13, 19, 21, 30, 38, 39, 52, 53
- 5.** Operations Security: Chapters 13, 14, 15, 19, 21, 24, 36, 40, 47, 53, 57
- 6.** Physical Security: Chapters 4, 13, 15, 19, 22, 23, 28, 29
- 7.** Cryptography: Chapters 7, 32, 37, 42
- 8.** Telecomm, Networks, and Internet Security: Chapters 4, 5, 6, 13, 14, 15, 16, 17, 18, 20, 21, 24, 25, 26, 27, 30, 31, 32, 33, 34, 35, 41, 48
- 9.** Business Continuity Planning: Chapters 22, 23, 56, 57, 58, 59, 60
- 10.** Law, Investigations, and Ethics: Chapters 11, 12, 13, 31, 42, 61

We have continued our practice from the fourth and fifth editions of inviting a security luminary to write the final chapter, “The Future of Information Assurance.” We are pleased to include this stellar contribution from Jeremy A. Hansen.

SEYMOUR BOSWORTH
Editor-in-Chief
February 2014

ACKNOWLEDGMENTS

Seymour Bosworth, Editor-in-Chief. I would like to give grateful recognition to Arthur Hutt and Douglas Hoyt, my coeditors of the first, second, and third editions of this *Handbook*. Although both Art and Doug are deceased, their commitment and their competence remain as constant reminders that nothing less than excellence is acceptable. Mich Kabay, my coeditor from the fourth and fifth editions, and Eric Whyne, our fellow editor from the fifth and now sixth editions, continue in that tradition. I would not have wanted to undertake this project without them.

Thanks are also due to our colleagues at John Wiley & Sons: Tim Burgard as former Acquisitions Editor, Helen Cho as Editorial Program Coordinator, Sheck Ho as Executive Editor, Kimberly Kappmeyer as Production Editor, Natasha Andrews as Senior Production Editor, and Darice Moore as Copyeditor. All have performed their duties in an exemplary manner and with unfailing kindness, courtesy, and professionalism.

M. E. Kabay, Technical Editor. I want to thank my beloved wife, Deborah Black, light of my life, for her support and understanding over the years that this project has taken away from our time together. I am also grateful to the authors who have selflessly contributed so much to updating the material presented in this text.

Eric Whyne, Administrative Editor. An undertaking as big as pulling together this handbook would not be possible without my wife Lindsay and the love and support she gives to me and to our son Colton. I'd also like to thank the friends and mentors that have helped me most in my career: Mich and Sy, Tom Aldrich, Tom Payne, Frank Vanecik, and my parents Len and Terri. Any successful undertakings I've had, including this book, have been from listening to the advice they've given and aspiring to internalize the virtues that they exemplify. The authors who have contributed to this book also deserve many thanks for sharing their experience and wisdom. It is something for which I, myself, and the readers are extremely grateful.

ABOUT THE EDITORS

Seymour Bosworth, M.S., CDP (email: sybosworth55@gmail.com) is president of S. Bosworth & Associates, Plainview, New York, a management consulting firm specializing in computing applications for banking, commerce, and industry. Since 1972, he has been a contributing editor for all six editions of the *Computer Security Handbook*, and for several editions has been Editor-in-Chief. He has written many articles and lectured extensively about computer security and other technical and managerial subjects. He has been responsible for the design and manufacture, systems analysis, programming, and operations, of both digital and analog computers. For his technical contributions, including an error-computing calibrator, a programming aid, and an analog-to-digital converter, he has been granted a number of patents, and is working on several others.

Bosworth is a former president and CEO of Computer Corporation of America, manufacturers of computers for scientific and engineering applications; president of Abbey Electronics Corporation, manufacturers of precision electronic instruments and digital devices; and president of Alpha Data Processing Corporation, a general-purpose computer service bureau. As a vice president at Bankers Trust Company, he had overall responsibility for computer operations, including security concerns.

For more than 20 years, Bosworth was an adjunct associate professor of management at the Information Technologies Institute of New York University, where he lectured on computer security and related disciplines. He has conducted many seminars and training sessions for the Battelle Institute, New York University, the Negotiation Institute, the American Management Association, and other prestigious organizations. For many years he served as arbitrator, chief arbitrator, and panelist for the American Arbitration Association. He holds a master's degree from the Graduate School of Business of Columbia University and a Certificate in Data Processing from the Data Processing Management Association.

M. E. Kabay, Ph.D., CISSP-ISSMP (email: mekabay@gmail.com) has been programming since 1966. In 1976, he received his Ph.D. from Dartmouth College in applied statistics and invertebrate zoology. After joining a compiler and relational database team in 1979, he worked for Hewlett-Packard (Canada) Ltd. from 1980 through 1983 as an HP3000 operating system performance specialist and then ran operations at a large service bureau in Montréal in the mid-1980s before founding his own operations management consultancy. From 1986 to 1996, he was an adjunct instructor in the John Abbott College professional programs in programming and in technical support. He was director of education for the National Computer Security Association from 1991 to the end of 1999 and was security leader for the INFOSEC Group of AtomicTangerine, Inc., from January 2000 to June 2001. In July 2001, he joined the

xvi ABOUT THE EDITORS

faculty at Norwich University as associate professor of computer information systems in the School of Business and Management. In January 2002, he took on additional duties as the director of the graduate program in information assurance in the School of Graduate and Continuing Studies at Norwich, where he was also chief technical officer for several years. He returned to full-time teaching in the School of Business and Management in 2009 and was promoted to professor of computer information systems in 2011. He serves as associate director of the Norwich University Center for Advanced Computing and Digital Forensics.

Kabay was inducted into the Information Systems Security Association Hall of Fame in 2004. He has published more than 1,500 articles in operations management and security in several trade journals since 1986. He wrote two columns a week for *Network World Security Strategies* between 2000 and 2011; archives are at www.mekabay.com/nwss. For the last three editions, Kabay has been Technical Editor of the *Computer Security Handbook*. He also has a Website with freely available teaching materials and papers at www.mekabay.com.

Eric Whyne (email: ericwhyne@gmail.com), administrative editor of the *Computer Security Handbook*, is a technical manager and engineer at Data Tactics Corporation where he develops solutions that benefit national security, currently managing and working on several DARPA-funded big data and data science projects. Prior to this position, he was employed by Exelis Corp. and managed the Joint Improvised Explosive Device Defeat Organization (JIEDDO) Counter-IED Operations Integration Center (COIC) Systems Integration Laboratory (SIL), which consisted of engineers and analysts tasked to develop, deploy, and maintain global software and systems designed to provide intelligence to help predict and prevent explosive devices in Iraq and Afghanistan. Previously, he worked as an engineer with Pennsylvania State University Applied Research Labs (PSU ARL) researching the development and use of immersive visualization systems, geospatial information systems, visual data mining, and deploying touch interfaces in operational environments.

Prior to his industry experience, Whyne spent nine years on active duty in the United States Marine Corps (USMC) in ground combat units, starting as enlisted and attaining the rank of captain. His accomplishments in the Marine Corps include two meritorious promotions and a Navy Commendation Medal with Valor distinction for actions in combat. During his time in the military, he worked in the fields of signals intelligence and communications and served as an advisor to the Iraqi Army. Since 2005, he has been the coordinating editor for the 5th and 6th editions of the *Computer Security Handbook*. He contributes to several open source projects, serves as an invited technical expert in the W3C HTML5 working group, and is a member of the AFCEA Technology Committee.

Whyne attended Norwich University and graduated magna cum laude with a B.S. in computer science and minor degrees in mathematics, information assurance, and engineering.

ABOUT THE CONTRIBUTORS

Wendy Adams Carr currently works for the U.S. Army Corps of Engineers as a member of the Computer Incident Response Team (CIRT). Prior to this she performed as an Information Assurance Security Engineer with Booz Allen & Hamilton, where she supported a Department of Defense client in developing and maintaining DITSCAP and DIACAP-based certification and accreditation of complex, large-scale Information Systems. She is retired from the U.S. Army. She is also an active member of Infragard.

Mani Akella, a director (technology), has been actively working with information-security architectures and identity protection for Consultantgurus and its clients. An industry professional for 20 years, he has worked with hardware, software, networking, and all the associated technologies that service information in all of its incarnations and aspects. Over the years, he has developed a particular affinity for international data law and understanding people and why they do what they do (or do not). He firmly believes that the best law and policy is that which understands and accounts for cross-cultural differences, and works with an understanding of culture and societal influences. To that end, he has been actively working with all his clients and business acquaintances to improve security policies and make them more people-friendly: His experience has been that the best policy is that which works with, instead of being antagonistic to, the end user.

Rebecca Gurley Bace is the president/CEO of Infidel, Inc., a strategic consulting practice headquartered in Scotts Valley, California. She is also a venture consultant for Palo Alto-based Trident Capital, where she is credited with building Trident's investment portfolio of security product and service firms. Her areas of expertise include intrusion detection and prevention, vulnerability analysis and mitigation, and the technical transfer of information-security research results to the commercial product realm. Prior to transitioning to the commercial world, she worked in the public sector, first at the National Security Agency, where she led the Intrusion Detection research program, then at the Computing Division of the Los Alamos National Laboratory, where she served as deputy security officer. Her publishing credits include two books, an NIST Special Publication on intrusion detection and prevention, and numerous articles on information-security technology topics.

Susan Baumes, MS, CISSP, is an information-security professional working in the financial services industry. In her current role, she works across the enterprise to develop information-security awareness and is responsible for application security. Her role also extends to policy development, compliance, and audit. She has 11 years'

xviii ABOUT THE CONTRIBUTORS

experience in application development, systems and network administration, database management, and information security. Previously, she worked in a number of different sectors, including government (federal and state), academia, and retail.

Kurt Baumgarten, CISA, is vice president of information security and a partner at Peritus Security Partners, LLC, a leader in providing compliance-driven information security solutions. He is also a lecturer, consultant, and the developer of the DDIPS intrusion prevention technology as well as a pioneer in using best practices frameworks for the improvement of information technology security programs and management systems. He has authored multiple articles about the business benefits of sound information technology and information assurance practices, and assists businesses and government agencies in defining strategic plans that enhance IT and IA as positive value chain modifiers. He holds both a master's of science in information assurance and an M.B.A. with a concentration in e-commerce, and serves as an adjunct professor of information assurance. He has more than 20 years of experience in IT infrastructure and information security and is an active member of ISSA, ISACA, ISSSP, and the MIT Enterprise Forum. He periodically acts as an interim Director within external organizations in order to facilitate strategic operational changes in IT and information security.

Kevin Beets has been a research scientist with McAfee for over nine years. His work has concentrated on vulnerability, exploit and malware analysis, and documentation for the Foundstone and McAfee Labs teams. Prior to working with McAfee, he architected private LANS as well as built, monitored, and supported CheckPoint and PIX firewalls and RealSecure IDS systems.

Matt Bishop is a professor in the Department of Computer Science at the University of California at Davis and a codirector of the Computer Security Laboratory. His main research area is the analysis of vulnerabilities in computer systems, especially their origin, detection, and remediation. He also studies network security, policy modeling, and electronic voting. His textbook, *Computer Security: Art and Science*, is used widely in advanced undergraduate and graduate courses. He received his Ph.D. in computer science from Purdue University, where he specialized in computer security, in 1984.

Kip Boyle is the chief information-security officer of PEMCO Insurance, a \$350 million property, casualty, and life insurance company serving the Pacific Northwest. Prior to joining PEMCO Insurance, he held such positions as chief security officer for a \$50 million national credit card transaction processor and technology service provider; authentication and encryption product manager for Cable & Wireless America; senior security architect for Digital Island, Inc.; and a senior consultant in the Information Security Group at Stanford Research Institute (SRI) Consulting. He has also held director-level positions in information systems and network security for the U.S. Air Force. He is a Certified Information System Security Professional and Certified Information Security Manager. He holds a bachelor's of science in computer information systems from the University of Tampa (where he was an Air Force ROTC Distinguished Graduate) and a master's of science in management from Troy State University.

Jennifer Bradley is a member of the first Master of Science in Information Assurance graduating class at Norwich University. She is the primary Systems and Security Consultant for Indiana Networking in Lafayette, Indiana, and has served as both a

ABOUT THE CONTRIBUTORS xix

network and systems administrator in higher education and private consulting. She has almost 15 years' experience as a programmer and instructor of Web technologies, with additional interests in data backup, virtualization, authentication/identification, monitoring, desktop and server deployment, and incident response. At present she serves as an independent consultant. She has previously worked as a tester for quality and performance projects for Google, Inc., and as a collegiate adjunct instructor in computer technologies. She received a bachelor's of science in Industrial and Computer Technology from Purdue University.

Timothy Braithwaite has more than 30 years of hands-on experience in all aspects of automated information processing and communications. He is currently the deputy director of strategic programs at the Center for Information Assurance of Titan Corporation. Before joining Titan, he managed most aspects of information technology, including data and communications centers, software development projects, strategic planning and budget organizations, system security programs, and quality improvement initiatives. His pioneering work in computer systems and communications security while with the Department of Defense resulted in his selection to be the first systems security officer for the Social Security Administration (SSA) in 1980. After developing security policy and establishing a nationwide network of regional security officers, he directed the risk assessment of all payment systems for the agency. In 1982, he assumed the duties of deputy director, systems planning and control of the SSA, where he performed substantive reviews of all major acquisitions for the associate commissioner for systems and, through a facilitation process, personally led the development of the first Strategic Systems Plan for the administration. In 1984, he became director of information and communication services for the Bureau of Alcohol, Tobacco, and Firearms at the Department of Treasury. In the private sector, he worked in senior technical and business development positions for SAGE Federal Systems, a software development company; Validity Corporation, a testing and independent validation and verification company; and J.G. Van Dyke & Associates, where he was director, Y2K testing services. He was recruited to join Titan Corporation in December 1999 to assist in establishing and growing the company's Information Assurance practice.

Dr. Paul Brusil founded Strategic Management Directions, a security and enterprise management consultancy in Beverly, Massachusetts. He has been working with various industry and government sectors, including healthcare, telecommunications, and middleware to improve the specification, implementation, and use of trustworthy, quality, security-related products and systems. He supported strategic planning that led to the National Information Assurance Partnership and other industry forums created to understand, promote, and use the Common Criteria to develop security and assurance requirements and to evaluate products. He has organized, convened, and chaired several national workshops, conferences, and international symposia pertinent to management and security. Through these and other efforts to stimulate awareness and cooperation among competing market forces, he spearheaded industry's development of the initial open, secure, convergent, standards-based network and enterprise management solutions. While at the MITRE Corp, he led research and development critical to the commercialization of the world's first LAN solutions. Earlier, at Harvard, he pioneered research leading to noninvasive diagnosis of cardiopulmonary dysfunction. He is a Senior Member of the IEEE, a member of the Editorial Advisory Board of the *Journal of Network and Systems Management* (JNSM), has been senior technical editor for JNSM, is the guest editor for all JNSM's Special Issues on Security and Management,

xx ABOUT THE CONTRIBUTORS

and is a lead instructor for the adjunct faculty supporting the master's of science in information assurance degree program at Norwich University. He has authored over 100 papers and book chapters. He graduated from Harvard University with a joint degree in Engineering and Medicine.

Michael Buglewicz is employed at National Security Technologies as a section manager whose team provides technology and communications solutions to various government agencies. He spent 17 years at Microsoft in various roles in services and business management. Prior to Microsoft, he was involved with building some of the first Internet ecommerce and banking solutions while at First Data Corporation; he also spent ten years in law enforcement. In addition to his contributions to *The Computer Security Handbook*, he was a contributing author to *The Encyclopedia of Information Assurance* and, most recently, contributed to the book *Cloud Migration* by Tobias Hollwarth.

Dr. Joseph R. Bumblis is currently a research specialist with the Institute of Electrical and Electronic Engineers (IEEE) Twin Cities (TC) Section's Phoenix Project, where he conducts research and engineering projects in the areas of sensors, signal processing, and embedded systems design. His expertise includes computer networks, embedded systems with FPGA and SoC codesign, IT systems security, and software engineering methodologies. As an Associate Professor of Computer Engineering at the University of Wisconsin-Stout (UW-Stout), he developed computer engineering curriculum and taught courses in digital design, solid-state devices, embedded systems design, and Verilog programming. Prior to joining UW-Stout, he served as an IT systems architect at BAE Systems and held several adjunct professor positions where he taught software engineering and computer networking courses.

Q. Campbell has worked in the information-security field for over six years. He specializes in information-security threat analysis and education.

Santosh Chokhani is the founder and president of CygnaCom Solutions, Inc., an Entrust company specializing in PKI. He has made numerous contributions to PKI technology and related standards, including trust models, security, and policy and revocation processing. He is the inventor of the PKI Certificate Policy and Certification Practices Statement Framework. His pioneering work in this area led to the Internet RFC that is used as the standard for CP and CPS by governments and industry throughout the world. Before starting CygnaCom, he worked for The MITRE Corporation from 1978 to 1994. At MITRE, he was senior technical manager and managed a variety of technology research, development, and engineering projects in the areas of PKI, computer security, expert systems, image processing, and computer graphics. Chokhani obtained his master's (1971) and Ph.D. (1975) in electrical engineering/computer science from Rutgers University, where he was a Louis Bevier Fellow from 1971 to 1973.

Christopher Christian is a first lieutenant and an aviator in the United States Army. He received a bachelor's degree in Computer Information Systems at Norwich University class of 2005. His primary focus of study was Information Assurance and Security. He worked as an intern for an engineering consulting company for three years. He developed cost/analysis worksheets and floor-plan layouts to maximize workspace efficiency for companies in various industries. He graduated flight school at Fort Rucker, Alabama, where he trained on the H-60 Blackhawk. He serves as a flight

ABOUT THE CONTRIBUTORS xxii

platoon leader in an air assault battalion. He is currently serving in Iraq in support of Operation Iraqi Freedom 08–09.

Chey Cobb, CISSP, began her career in information security while at the National Computer Security Association (now known as TruSecure/ICSA Labs). During her tenure as the NCSA award-winning Webmaster, she realized that Web servers often created security holes in networks and became an outspoken advocate of systems security. Later, while developing secure networks for the Air Force in Florida, her work captured the attention of the U.S. intelligence agencies. She moved to Virginia and began working for the government as the senior technical security advisor on highly classified projects. Ultimately, she went on to manage the security program at an overseas site. Now semiretired, she writes books and articles on computer security and is a frequent speaker at security conferences.

Stephen Cobb, CISSP, is an independent information-security consultant and an adjunct professor of information assurance at Norwich University, Vermont. A graduate of the University of Leeds, his areas of expertise include risk assessment, computer fraud, data privacy, business continuity management, and security awareness and education. A frequent speaker and seminar leader at industry conferences around the world, he is the author of numerous books on security and privacy as well as hundreds of articles. He cofounded several security companies whose products expanded the range of security solutions available to enterprises and government agencies. As a consultant, he has advised some of the world's largest companies on how to maximize the benefits of information technology by minimizing IT risks.

Caleb S. Coggins, MSIA, GSNA, CISSP, CISA, currently works for Sylint in the area of network forensics. Prior to Sylint, he operated in an internal audit and advisory services capacity for a healthcare IT company in Tennessee that focused on revenue and payment cycle management. Previous to that, he served in IT, security, and audit functions at Bridgestone Americas and its subsidiaries for over eight years. During his time working in the Americas and West Africa, he has enjoyed collaborating with management and teammates in identifying practical and effective solutions, while reducing risk to business operations. Prior to Bridgestone, he was the information manager for a private company as well as an information-security consultant to business clients. He holds a B.A. from Willamette University and an M.S. in information assurance from Norwich University.

Bernie Cowens, CISSP, CISA, is chief information-security officer at a Fortune 500 company in the financial services industry. He is an information risk, privacy, and security expert with more than 20 years' experience in industries including defense, high technology, healthcare, financial, and Big Four professional services. He has created, trained, and led a number of computer emergency, forensic investigation, and incident response teams over the years. He has real-world experience responding to attacks, disasters, and failures resulting from a variety of sources, including malicious attackers, criminals, and foreign governments. He has served as an advisor to and a member of national-level panels charged with analyzing cybersystem threats to critical infrastructures, assessing associated risks, and recommending both technical and nontechnical mitigation policies and procedures. He holds a master's degree in management information systems along with undergraduate degrees and certificates in systems management and information processing.

xxii ABOUT THE CONTRIBUTORS

Tim Crothers is an IT manager at 3M IT.

Rob Cryan is chief information-security officer for MAPFRE USA, consisting of seven property and casualty insurance companies in the United States. He has corporate responsibility for all aspects of information security and business continuity. He has worked in both the corporate and the consulting fields, managing and delivering security solutions. His current areas of focus include managing risk in cloud computing and delivering cost-effective risk management derived from larger, often cumbersome models to apply uniform controls across multiple compliance disciplines. He received his B.S. in business administration management from the University of Maine. He is currently pursuing his M.S. in information systems and technology management with a concentration in information assurance and security.

Christopher Dantos is a senior architectural specialist with Computer Science Corporation's Global Security Solutions Group. His areas of expertise include 802.11, VoIP, and Web application security. Prior to joining CSC, he spent 10 years as a security architect with Motorola Inc., including five years in the Motorola Labs Wireless Access Research Center of Excellence. He holds a master's of science degree in information assurance from Norwich University and a bachelor's of science degree in marine engineering from the Maine Maritime Academy.

Scott L. David is executive director of the Law, Technology, and Arts Group at the University of Washington School of Law.

Chris Davis, MBA, CISA, CISSP, CCNP, finds that his role as a cloud security and compliance product manager at VCE enables him to apply his past experiences to the latest in data center computing. He has trained and presented in information security, audit, forensic analysis, hardware security design, auditing, and systems engineering for government, corporate, and university requirements. He has written or contributed to nine books covering multiple security disciplines and teaches as an adjunct professor at Southern Methodist University covering graduate courses in Information Security and Risk Management (EMIS7380) and IT Controls (EMIS7382). His professional career has taken him through Texas Instruments followed by several startup and consulting roles. He holds a bachelor's degree in nuclear engineering technologies from Thomas Edison State College and a master's in business from the University of Texas McCombs School of Business at Austin. He served eight years in the U.S. Naval Submarine Fleet onboard the special projects Submarine NR-1 and the ballistic missile submarine USS Nebraska.

Seth Finkelstein is a professional programmer with degrees in Mathematics and in Physics from MIT. He cofounded the Censorware Project, an anti-censorware advocacy group. In 1998, his efforts evaluating the sites blocked by the library's Internet policy in Loudoun County, Virginia, helped the American Civil Liberties Union win a federal lawsuit challenging the policy. In 2001, he received a Pioneer of the Electronic Frontier Award from the Electronic Frontier Foundation for his groundbreaking work in analyzing content-blocking software. In 2003, he was primarily responsible for winning a temporary exemption in the Digital Millennium Copyright Act allowing for the analysis of censorware.

Eric Fisher is a systems architecture and network security engineer for Penn State University's Applied Research Laboratories and is currently a master's candidate

ABOUT THE CONTRIBUTORS xxiii

in the field of information security and forensics. Eric has extensive experience designing and implementing high performance and redundant compute systems for the Department of Defense and the IC at large. His areas of expertise are designing and implementing secure and scalable information systems, and the networks that connect them. Before joining Penn State, he worked for the Raytheon Corporation as a *nix/Windows/Network administrator where he managed large-scale high-availability clusters.

Robert Gezelter, CDP, has over 33 years of experience in computing, starting with programming scientific/technical problems. Shortly thereafter, his focus shifted to operating systems, networks, security, and related matters, where he has 32 years of experience in systems architecture, programming, and management. He has worked extensively in systems architecture, security, internals, and networks, ranging from high-level strategic issues to the low-level specification, design, and implementation of device protocols and embedded firmware. He is an alumnus of the IEEE Computer Society's Distinguished Visitor Program for North America, having been appointed to a three-year term in 2004. His appointment included many presentations at Computer Society chapters throughout North America. He has published numerous articles, appearing in *Hardcopy*, *Computer Purchasing Update*, *Network Computing*, *Open Systems Today*, *Digital Systems Journal*, and *Network World*. He is a frequent presenter at conference sessions on operating systems, languages, security, networks, and related topics at local, regional, national, and international conferences, speaking for DECUS, Encompass, IEEE, ISSA, ISACA, and others. He previously authored the mobile code and Internet-related chapters for the 4th edition of this *Handbook* (2002) as well as the "Internet Security" chapters of the 3rd edition (1995) and its supplement (1997). He is a graduate of New York University with B.A. (1981) and M.S. (1983) degrees in computer science. He founded his consulting practice in 1978, working with clients both locally and internationally. He maintains his offices in Flushing, New York. He may be contacted via his firm's Website at www.rlgsc.com.

Dr. Anup K. Ghosh is president and chief executive of Secure Command, LLC, a security software start-up developing next-generation Internet security products for corporate networks. He also holds a position as research professor at George Mason University. He was previously senior scientist and program Manager in the Advanced Technology Office of the Defense Advanced Research Projects Agency (DARPA), where he managed an extensive portfolio of information assurance and information operations programs. He previously served in executive management as Vice President of Research at Digital, Inc. He has served as principal investigator on contracts from DARPA, NSA, and NIST's Advanced Technology Program and has written more than 40 peer-reviewed conference and journal articles. He is also author of three books on computer network defense, serves on the editorial board of *IEEE Security and Privacy Magazine*, and has been guest editor for *IEEE Software* and *IEEE Journal on Selected Areas in Communications*. He is a Senior Member of the IEEE. For his contributions to the Department of Defense's information assurance, he was awarded the Frank B. Rowlett Trophy for Individual Contributions by the National Security Agency in November 2005, a federal government-wide award. He was also awarded the Office of the Secretary of Defense Medal for Exceptional Public Service for his contributions while at DARPA. In 2005, Worcester Polytechnic Institute awarded him its Hobart Newell Award for Outstanding Contributions to the Electrical and Computer Engineering Profession. He has previously been awarded the IEEE's Millennium Medal for Outstanding Contributions to E-Commerce Security. He completed his Ph.D. and

xxiv ABOUT THE CONTRIBUTORS

master's of science in electrical engineering at the University of Virginia and his bachelor's of science in electrical engineering at Worcester Polytechnic Institute.

Donald Glass, CISA, CISSP, has over 15 years of experience in the IT Auditing and Information Security fields. He is the current director of IT audit for Kerzner International. Author of several information security and IT audit articles, he is recognized as a leader in the IT audit field and information security.

Robert Guess, CISSP, NSA-IAM, NSA-IEM, is a senior security engineer at a Fortune 500 organization and an Associate Professor of Information Systems Technology. He possesses an M.S. in information assurance and has over a dozen industry certifications, including the Certified Information Systems Security Practitioner (CISSP), National Security Agency INFOSEC Assessment Methodologist (NSA-IAM), and National Security Agency INFOSEC Evaluation Methodologist (NSA-IEM). In addition to his academic experience, his professional experience includes work within the air, space, and defense sectors, serving as primary subject matter expert on a National Science Foundation cybersecurity grant, the development of workforce certification standards for information assurance professionals (DOD 8570.01-m), and periodic work as both an editor and author (e.g., the 5th and 6th editions of *The Computer Security Handbook*, etc.). His work in recent years has focused on penetration testing, incident response, the forensic analysis of digital evidence, virtualization, and cloud computing.

David Gursky, CISA, CISM, CISSP, is an information assurance manager and researcher at Raytheon Integrated Defense Systems working in Crystal City, Virginia. He is principal investigator for behavior-based intrusion detection systems, attribute-based access control, and resource-efficient authentication techniques. He held several senior positions as a Department of Defense contractor supporting information assurance programs and has over 30 years' experience in information technology and information security. He has conducted numerous security audits for PriceWaterhouse and Coopers. He has a bachelor's of science degree in business management from Southern New Hampshire University, a master's of science degree from Norwich University, and an M.B.A. from Northeastern University. In addition, he holds a CISA, CISM, and CISSP certifications. He lives in Northern Virginia and is an active member of (ISC)² and ISACA.

Jennifer Hadley is a member of the first master's of science in information assurance graduating class at Norwich University. She is the primary systems and security consultant for Indiana Networking in Lafayette, Indiana, and has served as both a network and systems administrator in higher education and private consulting. She has almost 10 years' experience as a programmer and instructor of Web technologies with additional interests in data backup, virtualization, authentication/identification, monitoring, desktop and server deployment, and incident response. At present, she serves as a technology consultant for Axcell Technologies, Inc. Previously she worked as a tester for quality and performance projects for Google, Inc., and as a collegiate adjunct instructor in computer technologies. She received a bachelor's of science degree in industrial and computer technology from Purdue University.

Carl Hallberg, CISSP, has been a UNIX systems administrator for years as well as an information-security consultant. He has also written training courses for subjects including firewalls, VPNs, and home network security. He has a bachelor's degree in

ABOUT THE CONTRIBUTORS xxv

psychology. Currently, he is a senior member of an incident response team for a major U.S. financial institution.

Jeremy A. Hansen is an information-security and cryptography educator and researcher.

David Harley, CITP, FBCS, CISSP, is CEO of Small Blue-Green World, COO of AVIEN, and ESET Senior Research Fellow, specializing in antimalware research and security/technical authoring and editing in a number of areas. His books include *Viruses Revealed* and *The AVIEN Malware Defense Guide*. Previously, he managed the UK National Health Service's Threat Assessment Centre, and before that he served as security analyst for the Imperial Cancer Research Fund (now Cancer Research UK). His academic background is in social sciences and computer science, and he is a fellow of the BCS Institute (formerly the British Computer Society).

Benjamin S. Hayes is an IT specialist at Travis Software.

Kevin Henry has been involved in computers since 1976, when he was an operator on the largest minicomputer system in Canada. He has since worked in many areas of information technology, including computer programming, systems analysis, and information technology audit. He was asked to become director of security based on the evidence of his audits and involvement in promoting secure IT operations. Following 20 years in the telecommunications field, he moved to a senior auditor position with the State of Oregon, where he was a member of the Governor's IT Security Subcommittee and performed audits on courts and court-related IT systems. He has extensive experience in risk management and business continuity and disaster recovery planning. He frequently presents papers at industry events and conferences and is on the preferred speakers list for nearly every major security conference. Since joining (ISC)² as their first full-time program manager in 2002, he has been responsible for research and development of new certifications, courseware, and development of educational programs and instructors. He has also been providing support services and consulting for organizations that require in-depth risk analysis and assistance with specific security-related challenges. This has led to numerous consulting engagements in the Middle East and Asia for large telecommunications firms, government departments, and commercial enterprises.

Don Holden, CISSP, ISSMP, is a principal consultant with Concordant specializing in information security. He has more than 20 years of management experience in information systems, security, encryption, business continuity, and disaster recovery planning in both industry and government. Previously he was a technology leader for RedSiren Technologies (formerly SRI Consulting). His achievements include leading a cyberinsurance joint venture project, developing privacy and encryption policies for major financial institutions, and recommending standards-based information technology security policies for a federal financial regulator. Holden is an adjunct professor for the Norwich University's master's of science in information assurance. He received an M.B.A. from Wharton and is a Certified Information System Security Professional and Information System Security Management Professional.

Dr. John D. Howard is a former Air Force engineer and test pilot who currently works in the Security and Networking Research Group at the Sandia National Laboratories, Livermore, California. His projects include development of the SecureLink software for

xxvi ABOUT THE CONTRIBUTORS

automatic encryption of network connections. He has extensive experience in systems development, including an aircraft–ground collision avoidance system for which he holds a patent. He is a graduate of the Air Force Academy, has master’s degrees in both aeronautical engineering and political science, and has a Ph.D. in engineering and public policy from Carnegie Mellon University.

Arthur E. Hutt, CCEP. The late Arthur E. Hutt was an information systems consultant with extensive experience in banking, industry, and government. He served as a contributing editor to the 1st, 2nd, and 3rd editions of *The Computer Security Handbook*. He was a principal of PAGE Assured Systems, Inc., a consulting group specializing in security and control of information systems and contingency/disaster recovery planning. He was a senior information systems executive for several major banks active in domestic and international banking. His innovative and pioneering development of online banking systems received international recognition. He was also noted for his contributions to computer security and to information systems planning for municipal government. He was on the faculty of the City University of New York and served as a consultant to CUNY on curriculum and on data processing management. He also served on the mayor’s technical advisory panel for the City of New York. He was active in development of national and international technical standards, via ANSI and ISO, for the banking industry.

John B. Ippolito, CISSP, PMP, is director of information assurance services for the Allied Technology Group, Inc. He has more than 35 years of IT experience including “hands-on” experience with legacy mainframes, PCs, and virtual systems. He has supported government-wide standards efforts and coauthored NIST IT security training guidelines. He has a broad range of IT experience and is currently supporting federal agencies in the development of information assurance and privacy programs and the secure introduction of new technologies such as cloud-based systems and “bring your own device” into the business environment.

Robert V. Jacobson, CPP, CISSP. The late Robert V. Jacobson was the president of International Security Technology, Inc., a New York City–based risk-management consulting firm. He founded IST in 1978 to develop and apply superior risk-management systems. Current and past government and industry clients are located in the United States, Europe, Africa, Asia, and the Middle East. He pioneered many of the basic computer security concepts now in general use. He served as the first information system security officer at Chemical Bank, now known as J P Morgan Chase. He was a frequent lecturer and had written numerous technical articles. He held B.S. and M.S. degrees from Yale University, and was a Certified Information Systems Security Professional. He was also a Certified Protection Professional of the American Society for Industrial Security. He was a member of the National Fire Protection Association and the Information Systems Security Association. In 1991, he received the Fitzgerald Memorial Award for Excellence in Security from the New York Chapter of the ISSA.

David J. Johnson is an information-security architect with a Fortune 500 financial services company where he provides enterprise guidance and solutions on topics such as cloud computing security, data protection technologies, and biometric authentication. His prior experience includes information-security analysis, advising business and IT management and staff on security risk, and development of policies and

ABOUT THE CONTRIBUTORS xxvii

procedures for a Fortune 1000 company. He also has experience designing, developing, and maintaining electronic commerce (EC/EDI) infrastructure and data transfers for a national financial services company. He holds a B.S. in business administration from Oregon State University, an M.S. in information assurance from Norwich University, and multiple information-security certifications from (ISC)², ISACA, and the Cloud Security Alliance.

Henry L. Judy is of counsel to Kirkpatrick & Lockhart, a national U.S. law firm. He advises clients on a wide range of corporate and financial law matters as well as federal and state securities, legislative and regulatory matters, with a particular emphasis on financial institutions, housing finance, and technology law. He is recognized for his work on the jurisdiction and dispute resolution issues of electronic commerce. He is a graduate of Georgetown University (J.D. and A.B.).

Sean Kelley is the vice president of technology services at Evolent Technologies in Herndon, Virginia. He is a retired Naval Officer with over 15 years of information assurance experience. Sean Kelley has received an M.A. in computer resource and information management from Webster University and an M.S. in information technology management (information assurance core competency) from the Naval Postgraduate School (NPS) in Monterey, California. He concentrated his studies on computer and network security by taking classes through the NPS Center for INFOSEC Studies and Research (NPS CISR), the world's foremost center for military research and education in information assurance (IA), defensive information warfare, and computer and network security. During his tenure at NPS, he received several certifications from the Committee of National Security Systems (CNSS), which operates under NSA. He is a Certified Information Systems Security Professional (CISSP), and is also a Project Management Professional (PMP). He has been responsible for information systems at several high-level offices in Washington, DC, and in the operational setting.

David M. Kennedy, CISSP, is TruSecure Corporation's chief of research. He directs the Research Group to provide expert services to TruSecure Corporation members, clients, and staff. He supervises the Information Security Reconnaissance (IS/R) team, which collects security-relevant information, both above- and underground in TruSecure Corporation's IS/R data collection. IS/R provides biweekly and special topic reports to IS/R subscribers. He is a retired U.S. Army Military Police officer. In his last tour of duty, he was responsible for enterprise security of five LANs with Internet access and over 3,000 personal computers and workstations. He holds a B.S. in forensic science.

Dr. Gary C. Kessler, CISSP, CCE, is an associate professor of Homeland Security at Embry-Riddle Aeronautical University in Daytona Beach, Florida, specializing in cybersecurity. He is a member of the North Florida Internet Crimes against Children (ICAC) Task Force and an adjunct faculty member at Edith Cowan University in Perth, Western Australia. From 2011 to 2012, he was the program director of the M.S. in Information Assurance program at Norwich University in Northfield, Vermont; from 2000 to 2010, he was an associate professor at Champlain College in Burlington, Vermont, where he designed and directed undergraduate and graduate programs related to information security and digital forensics. He is a Certified Information Systems Security Professional (CISSP), Certified Computer Examiner (CCE), and on the board of directors of the Consortium of Digital Forensic Specialists (CDFS). He holds a B.A.

xxviii ABOUT THE CONTRIBUTORS

in mathematics, an M.S. in computer science, and a Ph.D. in computing technology in education. He is the coauthor of two professional texts and over 70 articles and papers, a frequent speaker at industry events, and immediate past editor-in-chief of the *Journal of Digital Forensics, Security, and Law*. More information about him can be found at www.garykessler.net.

Dr. Minjeong Kim is associate professor of Journalism and Technical Communication at the Colorado State University. Her areas of expertise include communication law and policy, copyright law, and digital media. Her research has been published in scholarly journals, including *Communication Law and Policy*, the *Journal of the Copyright Society of the U.S.A.*, the *Journal of Computer-Mediated Communication*, and *Telecommunications Policy*. She completed her M.A. and Ph.D. from the School of Journalism and Mass Communication at the University of North Carolina at Chapel Hill. Before joining the CSU faculty in Fall 2008, she was an assistant professor at Hawaii Pacific University in Honolulu.

David A. Land served in the U.S. Army as a Counterintelligence Special Agent. With David Christie, he developed and hosted the first Department of Defense Computer Crimes Conference. Since then he has investigated espionage cases for both the Army and the Department of Energy. He serves as the information technology coordinator for Anniston City Schools in Alabama and as an adjunct professor for Norwich University, his alma mater.

David T. Lang joined the U.S. Civil Service on August 15, 2011. He has more than 30 years of experience in technical program management, counterespionage, anti-terrorism, security, training, risk management, and law enforcement in private industry and the military. Before assuming his current position as director of the DCIN-TS PMO, he was the chief of enterprise architecture and security for the DCIN-TS PMO. His industry positions included director of federal security, director of information assurance and forensics, director of digital forensics, and director of external training. His military assignments included nearly a decade as a Special Agent for the Air Force Office of Special Investigations (AFOSI) and over a decade in special weapons.

Dr. David R. Lease is the Chief Solution Architect at Computer Sciences Corporation. He has over 30 years of technical and management experience in the information technology, security, telecommunications, and consulting industries. His recent projects include a \$2 billion security architecture redesign for a federal law enforcement agency and the design and implementation of a secure financial management system for an organization operating in 85 countries. He is a writer and frequent speaker at conferences for organizations in the intelligence community, Department of Defense, civilian federal agencies, as well as commercial and academic organizations. He is also a peer reviewer of technical research for the IEEE Computer Society. Additionally, he is on the faculty of Norwich University and the University of Fairfax, where he teaches graduate-level information assurance courses and supervises doctoral-level research.

Corinne Lefrançois is an information assurance analyst at the National Security Agency. She graduated from Norwich University with a bachelor's of science in business administration and accounting in 2004 and is a current student in Norwich University's master's of science in information assurance program.

ABOUT THE CONTRIBUTORS xxix

Dr. Diane E. Levine, DBA, CISSP, CFE, FBCI, CPS, is an internationally sought after security expert with Strategic Systems Management Ltd, a leader in the computer security field for 25 years. She is one of the initial developers of the Certification for Information System Professionals (CISSP) who contributed to the creation, acceptance, and expansion of both the program and the written examinations. Besides teaching both the CISSP and CISA review courses for many years, as an adjunct professor, she created the popular graduate program in Computer Security at New York University, a program that became the model for computer security university programs throughout the world. A prominent consultant, author, and teacher, she has had a notable career developing and implementing risk management, business continuity, and enterprise security systems in corporate, nonprofit, and government institutions and organizations throughout the world. She is widely published in trade and academic press and contributed numerous chapters to the 3rd, 4th, and 5th editions of *The Computer Security Handbook*. A sought-after public speaker and member of technical panels, she served on the board of directors for numerous professional education and certification organizations (i.e., Information System Security Professionals[ISSP], Association of Certified Fraud Examiners [ACFE], Business Continuity Institute [BCI], Contingency Planning Exchange [CPE], and the Information Security, Auditing and Control Association [ISACA]). She has published regularly in print and online at *Information Week*, *Information Security*, *Executive Security*, *Internet Week*, *The TechWriter*, *EarthWeb*, *Planet IT*, *Security, Technology & Design*, *internet.com*, *Smart Computing*, *Labmice.techtarget.net*, *Auerbach Publications*, *Journal for the EDP Auditors*, and many other prominent publications. Both through teaching and publishing, much of her professional career has been devoted to developing and expanding the specialty of computer security and mentoring others. As of 2011, she has been presented with numerous awards for excellence in and her contribution to the computer security and business continuity fields. A former high-ranking corporate executive, she currently works with corporations, organizations, and governments throughout the world and is actively involved with law enforcement agencies in the apprehension and conviction of computer security criminals and the implementation of computer security legislation. Previously published obituaries by John Wiley & Sons were never corrected and remain in print, although they are totally untrue and misleading.

Dr. James Landon Linderman is an associate professor in the Computer Information Systems department at Bentley College, Waltham, Massachusetts, where he has taught for 30 years. He is a research fellow at Bentley's Center for Business Ethics, and past vice-chair of the Faculty Senate. A resident of Fitzwilliam, New Hampshire, he is a Permanent Deacon in the Roman Catholic Diocese of Worcester, Massachusetts, and a consultant in the area of computer-assisted academic scheduling and timetable construction.

Steven Lovaas, MSIA, CISSP, is the information technology security manager for Colorado State University, where he is pursuing a Ph.D. in public communication and technology. His areas of expertise include IT security policy and architecture; communication and teaching of complex technical concepts; and security issues in both K–12 and higher education. He has taught for the M.S. program in information assurance at Norwich University, and founded OMBRA Research Associates to provide technical and security consultation for social science research. As part of his volunteer commitment to educating the next generation of scientists and engineers, he serves as the president of the Colorado Science Olympiad.

xxx ABOUT THE CONTRIBUTORS

Allysa Myers is the Director of Research for West Coast Labs. Her primary responsibilities are researching and analyzing technology and security threat trends as well as reviewing and developing test methodologies. Prior to joining West Coast Labs, she spent 10 years working in the Avert group at McAfee Security, Inc. While there, she wrote for the Avert blog and *Sage* magazine, plus several external publications. She also provided training demonstrations to new researchers within McAfee, along with other groups such as the Department of Defense and McAfee Technical Support and Anti-Spyware teams. She has been a member of various security industry groups, such as the Wildlist and the Drone Armies mailing list.

John Mason is the vice president and audit manager at PacTrust Bank, a part-time director at SSAE 16 Professionals, and an adjunct professor for Norwich University's MSIA program. He has over 20 years' combined experience in internal audit, regulatory compliance, information security, SSAE 16s/SAS 70s, enterprise risk management, investigations/loss prevention, and process reengineering. He has held senior positions such as chief internal auditor and vice president of audit and compliance in a variety of companies. At two multi-billion-dollar institutions, he was the chief information security officer (CISO) and helped establish information risk-management programs, as well as designing risk-based audit programs several years before Sarbanes-Oxley. He has routinely authored, reviewed, and researched finance control policies and procedures; performed audits for governmental agencies; and managed a full spectrum of financial, operational, SOX compliance, and data processing audits. He possesses an M.B.A., a B.A. in economics, and numerous certificates, including a CISM, CISA, CGEIT, and CFE. Currently, he resides in southern California. Aside from trying to help his students maximize their value, John's favorite tenet is "Learn from the mistakes of others, because you'll never live long enough to make them all yourself."

Peter Mell is a senior computer scientist in the computer security division at the National Institute of Standards and Technology. He is the program manager for the National Vulnerability Database as well as the Security Content Automation Protocol validation program. These programs are widely adopted within the U.S. government and used for standardizing and automating vulnerability and configuration management, measurement, and policy compliance checking. He has written the NIST publications on patching, malware, intrusion detection, common vulnerability scoring system, and the common vulnerabilities and exposures standard. His research experience includes the areas of intrusion detection systems, vulnerability scoring, and vulnerability databases.

Michael Miora has designed and assessed secure, survivable, highly robust systems for industry and government over the past 30 years. He was one of the original professionals granted the CISSP in the 1990s and the ISSMP in 2004, and was accepted as a Fellow of the Business Continuity Institute (BCI) in 2005. He founded and currently serves as president of ContingenZ Corporation (www.contingenz.com). He was educated at the University of California Los Angeles and the University of California Berkeley, earning bachelor's and master's degrees in mathematics. He is an adjunct professor at Norwich University, and a frequent speaker and consultant in the area of information assurance.

Richard O. Moore III has over 20 years of information-security experience, starting with 15 years in the Marine Corps intelligence community. After leaving the Marines

ABOUT THE CONTRIBUTORS xxxii

as the Regimental Intelligence Chief, he worked as a security consultant for Investor's Bank and Trust, preparing them for Sarbanes Oxley and a security review of their environment. At the end of that contract, he went to work for KPMG within their security consulting practice, where he performed security reviews, audits, and penetration testing activities. He left KPMG as an associate director in the Information Security Group managing the global penetration testing team and moved to the Royal Bank of Scotland Citizens Bank, where he created the security auditing function as the senior audit manager. He has continued with the Bank and is now the senior information security manager responsible for the Technical Services group, which consists of the Forensic, Penetration Testing, Risk Assessment, Cyber Intelligence, and Information Security Consultancy services. He has served on numerous industry boards, such as Board of Directors Massachusetts InfraGard, Norwich Editorial Board, and the ISACA Academic SubCommittee, and is a contributing author to the 5th and 6th editions of *The Computer Security Handbook*.

Scott J. Nathan, Esq., is an attorney whose practice includes litigation concerning intellectual property and technology matters, computer fraud and abuse, and environmental and insurance coverage matters involving the exchange of millions of pages of documents. In addition, he advises clients about, among other things, Internet-related risks and risk avoidance, proprietary and open source software licensing, service-level agreements, and insurance coverage. He has written and spoken extensively about such issues as online privacy, cyberspace jurisdiction, and the legal issues surrounding the use of open source software. He is admitted to practice before the United States Supreme Court, the United States Court of Appeals for the First Circuit, the Federal District Court for the District of Massachusetts, and the Courts of the Commonwealth of Massachusetts. He is a member of the American Bar Association's Litigation and Computer Litigation Committees.

Carl Ness, CISSP, is a senior security analyst for the University of Iowa. He has more than 10 years' experience in the information technology and information-security fields, mainly in the academic and healthcare sector. He is a speaker, author, and educator on information assurance, including security in the academic environment, messaging security, disaster recovery and business continuity, safe home computing, and information technology operations. He previously served as a systems administrator, network administrator, information technology director, and information-security officer. He holds an M.S. degree and also provides consulting to several security software development organizations.

Lester E. Nichols III has more than 15 years' experience in information technology, including: technology computing, cybersecurity, information assurance, and enterprise and security architecture. His experiences span the medical, nonprofit, financial, and local and federal government sectors in a variety of roles, including: application developer, network engineering, security operations manager, and information-security architect. He is a published author, contributing to major security publications such as the fifth and sixth editions of *The Computer Security Handbook*, and is a frequent contributor to the (ISC)² blog and other security-related communities and lectures. He also teaches college-level computing courses and provides computing lectures at secondary schools. Prior to joining ApplyLogic, he served in various technical, consulting, and leadership roles with both federal and commercial clients, such as the Administrative Offices of U.S. Courts, KCG, Prolific Solutions, the University of California Irvine,

xxxii ABOUT THE CONTRIBUTORS

and the University of Phoenix, and as a federal employee with the U.S. Department of Homeland Security, Transportation Security Administration.

Justin Opatrny currently works for a Fortune 500 company specializing in security infrastructure. He also is an adjunct instructor teaching information assurance courses and an independent consultant providing technology and security services. He holds a master's degree in information assurance from Norwich University; he holds industry certifications, including CISSP, GCFA, GSNA, and CASP; and he is an active member of the Information Systems Security Association and InfraGard.

Dr. John Orlando speaks and consults on the use of social media in disaster response, marketing, sales, business development, and training. He developed and managed the Master's of Science in Business Continuity Management program at Norwich University in Vermont, created the Social Media in Disaster Response LinkedIn Group, and has written over 20 articles on business continuity and disaster response in publications such as *Continuity Insights*, *Disaster Recovery Journal*, *Global Assurance*, *CFO*, and *The Definitive Handbook on Business Continuity Management*.

Dr. Raymond Panko is a professor of information technology management in the Shidler College of Business at the University of Hawaii. His interest in security began during lunches with Donn Parker in the 1970s at SRI International and has grown ever since. His textbook on IT security, *Corporate Computer and Network Security*, is published by Prentice-Hall. His current research focuses are security for end-user applications (especially spreadsheets), how to deal with fraud, and human and organizational controls. His main teaching focus, however, remains networking. In his networking classes and textbook, he emphasizes security throughput, pointing out the security implications of network protocols and practices.

Robert A. Parisi, Jr. is the senior vice-president and national technology, network risk and telecommunications practice leader for the FINPRO unit of Marsh USA. He has spoken at various businesses, technology, legal, and insurance forums throughout the world and has written on issues affecting professional liability, privacy, technology, telecommunications, media, intellectual property, computer security, and insurance. In 2002, he was honored by *Business Insurance* magazine as one of the Rising Stars of Insurance. Immediately prior to joining Marsh, he was the senior vice-president and chief underwriting officer of eBusiness Risk Solutions (a unit of the property and casualty companies of American International Group, Inc.). He joined the AIG group of companies in 1998 as legal counsel for its Professional Liability group and held several executive and legal positions within AIG, including the position of chief underwriting officer for professional liability and technology. While at AIG, he oversaw the creation and drafting of underwriting guidelines and policies for all lines of professional liability. Prior to joining AIG, he had been in private practice, principally as legal counsel to various Lloyds of London syndicates handling a variety of professional liability lines. He graduated cum laude from Fordham College with a B.A. in economics and received his law degree from Fordham University School of Law. He is admitted to practice in New York and the U.S. District Courts for the Eastern and Southern Districts of New York.

Donn B. Parker, CISSP, Fellow of the Association for Computing Machinery and Distinguished Fellow of the Information Systems Security Association (ISSA) is a retired (1997) senior management consultant from SRI International in Menlo Park,

ABOUT THE CONTRIBUTORS xxxiii

California, who has specialized in information security and computer crime research for 40 of his 50 years in the computer field. He is the founder of the International Information Integrity Institute (I-4), which provides a forum for exchange of security information among 75 of the largest international enterprises. He has written numerous books, papers, articles, and reports in his specialty based on interviews of over 200 computer criminals and reviews of the security of many large enterprises. He received the ISSA 1992 Individual Outstanding Achievement Award, the 1994 National Computer Systems Security Award from the US NIST and NSA, the Aerospace Computer Security Associates 1994 Distinguished Lecturer Award, and the MIS Training Institute Infosecurity News 1996 Lifetime Achievement Award. He received the ISSA Hall of Fame award and was chosen as one of five pioneers and profiled in *Information Security Magazine* in 1999.

Padgett Peterson, P.E., CISSP, IAM/IEM, has been involved with computer security and encryption for over 40 years. Since 1979 he has been employed by different elements of a major aerospace contractor. He is also an adjunct professor in the Master's of Science in Information Assurance program at Norwich University.

Franklin Platt is the founder and president of Office Planning Services, a Wall Street consultancy for 20 years headquartered in Stark, New Hampshire, since 1990. He has worked extensively in security planning, management, and preparedness in both the private and public sectors. His academic background includes business administration and electrical engineering. He has received extensive government training in emergency management, including terrorism and weapons of mass destruction, much of which is not available to the public. He holds many security certifications and is currently vetted by the FBI and by several states. His areas of expertise include: security risk management; compliance with the latest Homeland Security procedures and other federal regulations that affect the private sector; risk identification and assessment; vulnerability analysis; cost-value studies; response planning; site security surveys and compliance auditing; briefing and training; second opinion; and due diligence.

N. Todd Pritsky is the director of e-learning courseware at Hill Associates, a telecommunications training company in Colchester, Vermont. He is a senior member of the technical staff and an instructor of online, lecture, and hands-on classes. His teaching and writing specialties include e-commerce, network security, TCP/IP, and the Internet, and he also leads courses on fast packet and network access technologies. He enjoys writing articles on network security and is a contributing author of *Telecommunications: A Beginner's Guide* (McGraw-Hill/Osborne). Previously, he managed a computer center and created multimedia training programs. He holds a B.A. in philosophy and Russian/Soviet studies from Colby College.

Karthik Raman, CISSP, is a security researcher on the Adobe Secure Software Engineering Team (ASSET) where he focuses on vulnerability analysis and technical collaboration with industry partners. Before joining Adobe, he was a research scientist at McAfee Labs where he worked on threat analysis, building automation systems, malware analysis, and developing advanced anti-malware technology. He holds an M.S. degree in computer science from the University of California Irvine and B.S. degrees in computer science and computer security from Norwich University. Karthik has spoken at Infosec Southwest, SOURCE Boston, SOURCE Seattle, LayerOne, and Kaspersky Security Analyst Summit and has delivered a Black Hat Webcast.

xxxiv ABOUT THE CONTRIBUTORS

Bridgitt Robertson has been teaching business and technology courses for over a decade years. Her multidisciplinary approach to security awareness analyzes threats in the global enterprise and investigates how an educated workforce can mitigate risks and enhance corporate competitiveness. Prior to teaching, she worked for global companies in the areas of project management, business analysis, and consulting. She is a member of InfraGard.

Russell D. Rosco is currently program chair/professor of the Computer Information Network Technology program at Ivy Tech Community College–NorthEast. He is also a graduate student at Purdue University, where he is working on his dissertation for an Interdisciplinary Ph.D. in information assurance; combining his interests of computers, business, and psychology. He has an M.S. in information assurance from Norwich University, a B.S. in psychology, and a B.A. in business administration from Washington State University. Mr. Rosco had several years of IT work experience prior to his decade of teaching IT courses, ranging from intro to PCs to Active Directory design.

Marc Rotenberg is executive director of the Electronic Privacy Information Center in Washington, DC. He teaches information privacy law at Georgetown University Law Center. He has published many articles in legal and scientific journals. He is the coauthor of several books, including *Information Privacy Law*, *Privacy and Human Rights*, *The Privacy Law Sourcebook*, and *Litigation under the Federal Open Government Laws*. He frequently testifies before the U.S. Congress and the European Parliament on emerging privacy issues. He is a Fellow of the American Bar Foundation and the recipient of several awards, including the World Technology Award in Law.

K. Rudolph is the president and chief inspiration officer for Native Intelligence, Inc. She is a Certified Information Systems Security Professional (CISSP), a Federal IT Security Professional (FITSP-M), and holds a degree from Johns Hopkins University. Her interests include the psychology of security awareness and influence as related to learning and behavior, storytelling, and security-awareness metrics. Her publications include *System Forensics, Investigation, and Response* (2010) and the second edition of *Computer Forensics Jumpstart* (2011). She was also the technical editor for Michael Solomon's book, *Security Strategies in Windows Platforms and Applications*. She is a named contributor and participant in the group that created *NIST Special Publication 800-16, Information Technology Security Training Requirements: A Role-and Performance-Based Model*. In 2006, she was honored by the Federal Information Systems Security Educators' Association (FISSEA) as the Security Educator of the Year. She is currently a volunteer with (ISC)²'s Safe and Secure Online program.

Ravi Sandhu is cofounder and chief scientist of SingleSignOn.Net in Reston, Virginia, and professor of information technology and engineering at George Mason University in Fairfax, Virginia. An ACM and an IEEE Fellow, he is the founding editor-in-chief of *ACM's Transactions on Information and System Security*, chairman of ACM's Special Interest Group on Security, Audit and Control, and security editor for *IEEE Internet Computing*. He has published over 140 technical papers on information security. He is a popular teacher and has lectured all over the world. He has provided high-level consulting services to numerous private and government organizations.

Eric Salveggio is an information technology security professional who enjoys teaching online courses in CMIS for Liberty University and auditing for Norwich University.

ABOUT THE CONTRIBUTORS xxxv

He works as a trained ISO 17799, NSTISSI 4011 and 4013 consultant for Dynetics Corporation of Huntsville, Alabama, in IT Security and Auditing, and as a Private Consultant in networking, network design, and security (wired and wireless) with 10 years experience. He previously worked as the IT director for the Birmingham, Alabama, campus of Virginia College, where he opened two start-up campuses—on ground and online—and created three accredited programs (two undergrad, one graduate level) at state and federal levels in Network and Cyber Security. While in this position, he was chosen as a nominee for the 2006 Information Security Executive Award, and enjoyed being the only educational facility recognized. He was personally awarded a plaque of recognition by the Stonesoft Corporation for the same. He is a published author and photographer, and enjoys working at times as a technical editor for Pearson Education and Thomson Publishing on cyberforensics, cybersecurity, and operating systems.

Karen Scarfone is the principal consultant for Scarfone Cybersecurity. She provides cybersecurity publication consulting services to federal agencies, specializing in security automation standards and network and system security guidelines. Karen was formerly a senior computer scientist for the National Institute of Standards and Technology (NIST) and has co-authored over 50 NIST Special Publications and Interagency Reports. She holds bachelor's and master's degrees in computer science, and she has 20 years of professional experience in the IT field. Her security domains include general security engineering and administration, wired and wireless network security, host security, incident response, intrusion detection, log management, vulnerability measurement, and security automation.

Sondra Schneider is CEO and founder of Security University (SU), an information security and information assurance training and certification company. She and SU have challenged security professionals for the past 10 years, delivering hands-on tactical security classes and certifications around the world. Starting in 2008, SU set up an exam server to meet the demand for tactical security certifications. In 2005, SU refreshed the preexisting AIS security training program to the new “SU Qualified Programs,” which meet and exceed security professionals requirements for hands-on tactical security “skills” training. SU delivers the Qualified/Information Security Professional and Qualified/Information Assurance Professional certifications, which are the first of their kind that measure a candidate’s tactical hands-on security skills. In 2004, she was awarded Entrepreneur of the Year for the First Annual Women of Innovation Award from the Connecticut Technology Council. In 2007, she was tech editor for the popular 2007 *CEH V5 Study Guide*, and a multiple chapter author for the 2007 *CHFI Study Guide*. She sits on three advisory boards for computer security (start-up) technology companies and is a frequent speaker at computer security and wireless industry events. She is a founding member of the NYC HTCIA and IETF chapters, works closely with (ISC)², ISSA, and ISACA chapters, and the security and wireless vendor community. She specializes in information security, intrusion detection, information assurance (PKI), wireless security, and security-awareness training.

Jason Sinchak, CISSP, is an information-security consultant for Emerging Defense, LLC and an independent researcher specializing in risk assessment and penetration testing. He possesses a broad range of experience, beginning his career in application development and systems management, to operational security monitoring, enterprise security program advisory, incident response, and penetration testing. His security

xxxvi ABOUT THE CONTRIBUTORS

experience primarily stems from industry positions and consulting, where has had the pleasure of serving multiple Fortune 500 clients and deeply focusing on performing technical security assessments in a multitude of industries and environments.

William Stallings has made a unique contribution to understanding the broad sweep of technical developments in computer networking and computer architecture. He has authored 17 titles and, counting revised editions, has published a total of over 40 books on various aspects of these subjects. In over 20 years in the field, he has been a technical contributor, technical manager, and an executive with several high-technology firms. Currently, he is an independent consultant whose clients have included computer and networking manufacturers and customers, software development firms, and leading-edge government research institutions. He created and maintains the Computer Science Student Resource Site at www.computersciencestudent.com. He is a member of the editorial board of *Cryptologia*, a scholarly journal devoted to all aspects of cryptology.

Dr. Peter Stephenson is a writer, researcher and lecturer on information assurance and risk, information warfare and counterterrorism, and digital investigation and forensics on large-scale computer networks. He has lectured extensively on digital investigation and security and has written or contributed to 14 books and several hundred articles in major national and international trade, technical, and scientific publications. He is the associate program director in the Master's of Science in Information Assurance program at the Norwich University School of Graduate Studies, where he teaches information assurance, cybercrime and cyberlaw, and digital investigation on both the graduate and undergraduate levels. He is senior research scientist at the Norwich University Applied Research Institutes, chair of the Department of Computing, and the chief information-security officer for the university. He has lectured or delivered consulting engagements for the past 23 years in 11 countries plus the United States and has been a technologist for over 40 years. He operated a successful consulting practice for over 20 years and has worked for such companies as Siemens, Tektronix, and QinetiQ (UK). He obtained his Ph.D. in computer science at Oxford Brookes University, Oxford, England, where his research was in the structured investigation of digital incidents in complex computing environments. He holds a master's of arts degree in diplomacy with a concentration in terrorism from Norwich University. He is on the editorial advisory boards of *International Journal of Digital Evidence* and the Norwich University *Journal of Information Assurance* among several others, and serves as technology editor for *SC Magazine* and the editor in chief for Norwich University Press. He is a Fellow of the Institute for Communications, Arbitration and Forensics in the United Kingdom and is a member of Michigan InfraGard and the International Federation of Information Processing Technical Committee 11, Working Group 11.9, Digital Forensics. He serves on the steering Committee of the Michigan Electronic Crime Task Force. His research is focused on information conflict.

Gary L. Tagg, CISSP, is a highly experienced information-security professional with over 20 years working in the financial and government sectors. The organizations he has worked with include Deutsche Bank, PA Consulting group, Clearstream, Pearl Assurance, and Lloyds TSB. He has performed a wide range of security roles including risk management, consulting, security architecture, policy and standards, project management, development, testing, and auditing. He is currently a risk consultant in Deutsche Bank's IT security Governance Group.

ABOUT THE CONTRIBUTORS xxxvii

Nicholas Takacs is an information-security professional and business systems director for a long-term care insurance company. He is also an adjunct professor of information assurance at Norwich University. He has expertise in the areas of security policy management, security awareness, business continuity planning, and execution. Prior to moving into the insurance industry, he spent several years in the public utility industry focusing on the areas of regulatory compliance, disaster recovery, and identity management.

James Thomas, CISSP, is a senior partner with Norwich Security Associates, a full-spectrum information assurance consultancy headquartered in Scotland. Thomas spends most of his professional time providing policy, process, and governance advice to large banking and financial organizations in the United Kingdom and Europe. He is a 2004 graduate of the Norwich University Master's of Science in Information Assurance program. Prior to focusing his efforts in the security space, he had a long career in information technology and broadcast engineering spanning the United Kingdom and the eastern United States.

Lee Tien, Esq., is a senior staff attorney with the Electronic Frontier Foundation in San Francisco, California. He specializes in free speech and surveillance law and has authored several law review articles. He received his undergraduate degree in psychology from Stanford University and his law degree from Boalt Hall School of Law, University of California Berkeley. He is also a former newspaper reporter.

Timothy Virtue, CISSP, CISA, CCE, CFE, CIPP/G, is an accomplished information assurance leader with a focus in strategic enterprise technology risk management, information security, data privacy, and regulatory compliance. He has extensive experience with publicly traded corporations, privately held businesses, government agencies, and nonprofit organizations of all sizes.

Myles Walsh is an adjunct professor at three colleges in New Jersey: Ramapo College, County College of Morris, and Passaic County Community. For the past 12 years, he has taught courses in Microsoft Office and Web page design. He also implements small Office applications and Websites. From 1966 until 1989, he worked his way up from programmer to director in several positions at CBS, CBS Records, and CBS News. His formal education includes an M.B.A. from the Baruch School of Business and a B.B.A. from St. John's University.

Karen F. Worstell, CISM, is cofounder and principal of W Risk Group, a consultancy serving clients across multiple sectors to define due diligence to a defensible standard of care for information protection. Her areas of expertise include incident detection and management, compliance, governance, secure data management, and risk management. She is coauthor of *Evaluating the Electronic Discovery Capabilities of Outside Law Firms: A Model Request for Information and Analysis* (BNA, 2006) and is a frequent speaker and contributor in risk management and information-security forums internationally. She participates in ISACA, IIA, and the ABA Science and Technology Section, Information Security Committee, and serves as president of the Puget Sound Chapter of the ISSA.

Noel Zakin is president of RANCO Consulting LLC. He has been an information technology/telecommunications industry executive for over 45 years. He has held

xxxviii ABOUT THE CONTRIBUTORS

managerial positions at the Gartner Group, AT&T, the American Institute of CPAs, and Unisys. These positions involved strategic planning, market research, competitive analysis, business analysis, and education and training. His consulting assignments have ranged from the Fortune 500 to small start-ups and have involved data security, strategic planning, conference management, market research, and management of corporate operations. He has been active with ACM, IFIP, and AFIPS and currently with ISSA. He holds an M.B.A. from the Wharton School.

William A. Zucker, Esq., is a partner at McCarter & English, LLP's Boston office. He serves as a senior consultant for the Cutter Consortium on legal issues relating to information technology, outsourcing, and risk management, and is a member of the American Arbitration Association's National Technology Panel and a member of the CPR Institute's working group on technology business alliances and conflict management. He has also served on the faculty of Norwich University, where he taught the intellectual property aspects of computer security. He is a trial lawyer whose practice focuses on negotiation/litigation of business transactions, outsourcing/ebusiness and technology/intellectual property. Among his publications are: "The Legal Framework for Protecting Intellectual Property in the Field of Computing and Computer Software," written for *The Computer Security Handbook*, 4th edition, coauthored with Scott Nathan; and "Intellectual Property and Open Source: Copyright, Copyleft and Other Issues for the Community User."

A NOTE TO THE INSTRUCTOR

This two-volume text will serve the interests of practitioners and teachers of information assurance. The fourth edition and fifth editions of the *Handbook* were well received in academia; at least one-quarter of all copies were bought by university and college bookstores. The design of this sixth edition continues in the same vein and includes many updates to the material.

University professors looking for texts appropriate for a two-semester sequence of undergraduate courses in information assurance will find the *Handbook* most suitable. In my own work at Norwich University in Vermont, Volume I is the text for our *IS340 Introduction to Information Assurance* and Volume II is the basis for our *IS342 Management of Information Assurance* courses.

The text will also be useful as a resource in graduate courses. In the School of Graduate and Continuing Studies at Norwich University, both volumes have been used as required and supplementary reading for the 18-month, 36-credit Master's of Science in Information Security and Assurance program (MISA).

I will continue to create, update, and post PowerPoint lecture slides based on the chapters of the *Handbook* on my Website for free access by anyone applying them to noncommercial use (e.g., for self-study, for courses in academic institutions, and for unpaid industry training); the materials will be available in the IS340 and IS342 sections:

www.mekabay.com/courses/academic/norwich/is340/index.htm

www.mekabay.com/courses/academic/norwich/is342/index.htm

M. E. KABAY
Technical Editor
October 2013

Computer Security Handbook, Sixth Edition, Volume I
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

INTRODUCTION TO PART I

FOUNDATIONS OF COMPUTER SECURITY

The foundations of computer security include answers to the superficially simple question “What is this all about?” Our first part establishes a technological and historical context for information assurance so that readers will have a broad understanding of why information assurance matters in the real world. Chapters focus on principles that will underlie the rest of the text: historical perspective on the development of our field; how to conceptualize the goals of information assurance in a well-ordered schema that can be applied universally to all information systems; computer hardware and network elements underlying technical security; history and modern developments in cryptography; and how to discuss breaches of information security using a common technical language so that information can be shared, accumulated, and analyzed.

Readers also learn or review the basics of commonly used mathematical models of information-security concepts and how to interpret survey data and, in particular, the pitfalls of self-selection in sampling about crimes. Finally, the first section of the text introduces elements of law (U.S. and international) applying to information assurance. This legal framework from a layman’s viewpoint provides a basis for understanding later chapters; in particular, when examining privacy laws and management’s fiduciary responsibilities.

Chapter titles and topics in Part I include:

- 1. Brief History and Mission of Information System Security.** An overview focusing primarily on developments in the second half of the twentieth century and the first decade of the twenty-first century
- 2. History of Computer Crime.** A review of key computer crimes and notorious computer criminals from the 1970s to the mid-2000s
- 3. Toward a New Framework for Information Security.** A systematic and thorough conceptual framework and terminology for discussing the nature and goals of securing all aspects of information, not simply the classic triad of confidentiality, integrity, and availability
- 4. Hardware Elements of Security.** A review of computer and network hardware underlying discussions of computer and network security

I · 2 FOUNDATIONS OF COMPUTER SECURITY

5. **Data Communications and Information Security.** Fundamental principles and terminology of data communications, and their implications for information assurance
6. **Local Area Network Topologies, Protocols, and Design.** Information assurance of the communications infrastructure
7. **Encryption.** Historical perspectives on cryptography and steganography from ancient times to today as fundamental tools in securing information
8. **Using a Common Language for Computer Security Incident Information.** An analytic framework for understanding, describing, and discussing security breaches by using a common language of well-defined terms
9. **Mathematical Models of Computer Security.** A review of the most commonly referenced mathematical models used to describe information-security functions
10. **Understanding Studies and Surveys of Computer Crime.** Scientific and statistical principles for understanding studies and surveys of computer crime
11. **Fundamentals of Intellectual Property Law.** An introductory review of cyberlaw: laws governing computer-related crime, including contracts, and intellectual property (trade secrets, copyright, patents, open-source models). Also, violations (piracy, circumvention of technological defenses), computer intrusions, and international frameworks for legal cooperation

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 1

BRIEF HISTORY AND MISSION OF INFORMATION SYSTEM SECURITY

Seymour Bosworth and Robert V. Jacobson

1.1 INTRODUCTION TO INFORMATION SYSTEM SECURITY	1·1	1.2.15 1980s: The Personal Computer	1·9
		1.2.16 Local Area Networks	1·10
		1.2.17 1990s: Interconnection	1·11
		1.2.18 1990s: Total Interconnection	1·12
1.2 EVOLUTION OF INFORMATION SYSTEMS	1·3	1.2.19 Telecommuting	1·12
1.2.1 1950s: Punched-Card Systems	1·4	1.2.20 Internet and the World Wide Web	1·12
1.2.2 Large-Scale Computers	1·4	1.2.21 Virtualization and the Cloud	1·13
1.2.3 Medium-Size Computers	1·5	1.2.22 Supervisory Control and Data Acquisition	1·13
1.2.4 1960s: Small-Scale Computers	1·6		
1.2.5 Transistors and Core Memory	1·7		
1.2.6 Time Sharing	1·7		
1.2.7 Real-Time, Online Systems	1·7	1.3 GOVERNMENT RECOGNITION OF INFORMATION ASSURANCE	1·13
1.2.8 A Family of Computers	1·7	1.3.1 IA Standards	1·13
1.2.9 1970s: Microprocessors	1·8	1.3.2 Computers at Risk	1·14
1.2.10 The First Personal Computers	1·8	1.3.3 InfraGard	1·18
1.2.11 The First Network	1·8		
1.2.12 Further Security Considerations	1·9	1.4 RECENT DEVELOPMENTS	1·19
1.2.13 The First “Worm”	1·9		
1.2.14 1980s: Productivity Enhancements	1·9	1.5 ONGOING MISSION FOR INFORMATION SYSTEM SECURITY	1·20
		1.6 NOTES	1·20

1.1 INTRODUCTION TO INFORMATION SYSTEM SECURITY. The growth of computers and of information technology has been explosive. Never before has an entirely new technology been propagated around the world with such speed and with so great a penetration of virtually every human activity. Computers have brought vast benefits to fields as diverse as human genome studies, space exploration, artificial intelligence, and a host of applications from the trivial to the most life-enhancing.

Unfortunately, there is also a dark side to computers: They are used to design and build weapons of mass destruction as well as military aircraft, nuclear submarines,

1 · 2 BRIEF HISTORY AND MISSION OF INFORMATION SYSTEM SECURITY

and reconnaissance space stations. The computer's role in formulating biologic and chemical weapons, and in simulating their deployment, is one of its least auspicious uses.

Of somewhat lesser concern, computers used in financial applications, such as facilitating the purchase and sales of everything from matchsticks to mansions, and transferring trillions of dollars each day in electronic funds, are irresistible to miscreants; many of them see these activities as open invitations to fraud and theft. Computer systems, and their interconnecting networks, are also prey to vandals, malicious egoists, terrorists, and an array of individuals, groups, companies, and governments intent on using them to further their own ends, with total disregard for the effects on innocent victims. Besides these intentional attacks on computer systems, there are innumerable ways in which inadvertent errors can damage or destroy a computer's ability to perform its intended functions.

Because of these security problems and because of a great many others described in this volume, the growth of information systems security has paralleled that of the computer field itself. Only by a detailed study of the potential problems, and implementation of the suggested solutions, can computers be expected to fulfill their promise, with few of the security lapses that plague less adequately protected systems. This chapter defines a few of the most important terms of information security and includes a very brief history of computers and information systems, as a prelude to the works that follow.

Security can be defined as the state of being free from danger and not exposed to damage from accidents or attack, or it can be defined as the process for achieving that desirable state. The objective of information system security¹ is to optimize the performance of an organization with respect to the risks to which it is exposed.

Risk is defined as the chance of injury, damage, or loss. Thus, risk has two elements: (1) chance—an element of uncertainty, and (2) potential loss or damage. Except for the possibility of restitution, information system security actions taken today work to reduce *future* risk losses. Because of the uncertainty about future risk losses, perfect security, which implies zero losses, would be infinitely expensive. For this reason, risk managers strive to optimize the allocation of resources by minimizing the total cost of information system security measures taken and the risk losses experienced. This optimization process is commonly referred to as *risk management*.

Risk management in this sense is a three-part process:

1. Identification of material risks
2. Selection and implementation of measures to mitigate the risks
3. Tracking and evaluating of risk losses experienced, in order to validate the first two parts of the process

The purpose of this *Handbook* is to describe information security system risks, the measures available to mitigate these risks, and techniques for managing security risks. (For a more detailed discussion of risk assessment and management, see Chapters 47 and 54.)

Risk management has been a part of business for centuries. Renaissance merchants often used several vessels simultaneously, each carrying a portion of the merchandise, so that the loss of a single ship would not result in loss of the entire lot. At almost the same time, the concept of insurance evolved, first to provide economic protection against the loss of cargo and later to provide protection against the loss of buildings

EVOLUTION OF INFORMATION SYSTEMS 1 · 3

by fire. Fire insurers and municipal authorities began to require adherence to standards intended to reduce the risk of catastrophes like the Great Fire of London in 1666. The Insurance Institute was established in London one year later. With the emergence of corporations as limited liability stock companies, corporate directors have been required to use prudence and due diligence in protecting shareholders' assets. Security risks are among the threats to corporate assets that directors have an obligation to address.

Double-entry bookkeeping, another Renaissance invention, proved to be an excellent tool for measuring and controlling corporate assets. One objective was to make insider fraud more difficult to conceal. The concept of separation of duties emerged, calling for the use of processing procedures that required more than one person to complete a transaction. As the books of account became increasingly important, accounting standards were developed, and they continue to evolve to this day. These standards served to make books of account comparable and to assure outsiders that an organization's books of account presented an accurate picture of its condition and assets. These developments led, in turn, to the requirement that an outside auditor perform an independent review of the books of account and operating procedures.

The transition to automated accounting systems introduced additional security requirements. Some early safeguards, such as the rule against erasures or changes in the books of account, no longer applied. Some computerized accounting systems lacked an audit trail, and others could have the audit trail subverted as easily as actual entries.

Finally, with the advent of the Information Age, intellectual property has become an increasingly important part of corporate and governmental assets. At the same time that intellectual property has grown in importance, threats to intellectual property have become more dangerous, because of information system (IS) technology itself. When sensitive information was stored on paper and other tangible documents, and rapid copying was limited to photography, protection was relatively straightforward. Nevertheless, document control systems, information classification procedures, and need-to-know access controls were not foolproof, and information compromises occurred with dismaying regularity. Evolution of IS technology has made information control several orders of magnitude more complex. The evolution and, more importantly, the implementation of control techniques have not kept pace.

The balance of this chapter describes how the evolution of information systems has caused a parallel evolution of information system security and at the same time has increased the importance of anticipating the impact of technical changes yet to come. This overview will clarify the factors leading to today's information system security risk environment and mitigation techniques and will serve as a warning to remain alert to the implication of technical innovations as they appear. The remaining chapters of this *Handbook* discuss information system security risks, threats, and vulnerabilities, their prevention and remediation, and many related topics in considerable detail.

1.2 EVOLUTION OF INFORMATION SYSTEMS. The first electromechanical punched-card system for data processing, developed by Herman Hollerith at the end of the nineteenth century, was used to tabulate and total census field reports for the U.S. Bureau of the Census in 1890. The first digital, stored-program computers developed in the 1940s were used for military purposes, primarily cryptanalysis and the calculation and printing of artillery firing tables. At the same time, punched-card systems were already being used for accounting applications and were an obvious choice for data input to the new electronic computing machines.

1 · 4 BRIEF HISTORY AND MISSION OF INFORMATION SYSTEM SECURITY

1.2.1 1950s: Punched-Card Systems. In the 1950s, punched-card equipment dominated the commercial computer market.² These electromechanical devices could perform the full range of accounting and reporting functions. Because they were programmed by an intricate system of plugboards with a great many plug-in cables, and because care had to be exercised in handling and storing punched cards, only experienced persons were permitted near the equipment. Although any of these individuals could have set up the equipment for fraudulent use, or even engaged in sabotage, apparently few, if any, actually did so.

The punched-card accounting systems typically used four processing steps. As a preliminary, operators would be given a “batch” of documents, typically with an adding machine tape showing one or more “control totals.” The operator keyed the data on each document into a punched card and then added an extra card, the batch control card, which stored the batch totals. Each card consisted of 80 columns, each containing, at most, one character. A complete record of an inventory item, for example, would be contained on a single card. The card was called a unit record, and the machines that processed the cards were called either unit record or punched-card machines. It was from the necessity to squeeze as much data as possible into an 80-character card that the later Year 2000 problem arose. Compressing the year into two characters was a universally used space-saving measure; its consequences 40 years later were not foreseen.

A group of punched cards, also called a “batch,” were commonly held in a metal tray. Sometimes a batch would be rekeyed by a second operator, using a “verify-mode” rather than actually punching new holes in the cards, in order to detect keypunch errors before processing the card deck. Each batch of cards would be processed separately, so the processes were referred to as “batch jobs.”

The first step would be to run the batch of cards through a simple program, which would calculate the control totals and compare them with the totals on the batch control card. If the batch totals did not reconcile, the batch was sent back to the keypunch area for rekeying. If the totals reconciled, the deck would be sort-merged with other batches of the same transaction type, for example, the current payroll. When this step was complete, the new batch consisted of a punched card for each employee in employee-number order. The payroll program accepted this input data card deck and processed the cards one by one. Each card was matched up with the corresponding employee’s card in the payroll master deck to calculate the current net pay and itemized deductions and to punch a new payroll master card, including year-to-date totals. The final step was to use the card decks to print payroll checks and management reports. These steps were identical with those used by early, small-scale electronic computers. The only difference was in the speed at which the actual calculations were made. A complete process was still known as a batch job.

With this process, the potential for abuse was great. The machine operator could control every step of the operation. Although the data was punched into cards and verified by others, there was always a keypunch machine nearby for use by the machine operator. Theoretically, that person could punch a new payroll card and a new batch total card to match the change before printing checks and again afterward. The low incidence of reported exploits was due to the controls that discouraged such abuse, and possibly to the pride that machine operators experienced in their jobs.

1.2.2 Large-Scale Computers. While these electromechanical punched card machines were sold in large numbers, research laboratories and universities were working to design large-scale computers that would have a revolutionary effect on

EVOLUTION OF INFORMATION SYSTEMS 1 · 5

the entire field. These computers, built around vacuum tubes, are known as the first generation. In March 1951, the first Universal Automatic Computer (UNIVAC) was accepted by the U.S. Census Bureau. Until then, every computer had been a one-off design, but UNIVAC was the first large-scale, mass-produced computer, with a total of 46 built. The word “universal” in its name indicated that UNIVAC was also the first computer designed for both scientific and business applications.³

UNIVAC contained 5,200 vacuum tubes, weighed 29,000 pounds, and consumed 125 kilowatts of electrical power. It dispensed with punched cards, receiving input from half-inch-wide metal tape recorded from keyboards, with output either to a similar tape or to a printer. Although not a model for future designs, its memory consisted of 1,000 72-bit words and was fabricated as a mercury delay line. Housed in a cabinet about six feet tall, two feet wide, and two feet deep was a mercury-filled coil running from top to bottom. A transducer at the top propagated slow-moving waves of energy down the coil to a receiving transducer at the bottom. There it was reconverted into electrical energy and passed on to the appropriate circuit, or recirculated if longer storage was required.

In 1956, IBM introduced the Random Access Method of Accounting and Control (RAMAC) magnetic disk system. It consisted of 50 magnetically coated metal disks, each 24 inches in diameter, mounted on a common spindle. Under servo control, two coupled read/write heads moved to span each side of the required disk and then inward to any one of 100 tracks. In one revolution of the disks, any or all of the information on those two tracks could be read out or recorded. The entire system was almost the size of a compact car and held what, for that time, was a tremendous amount of data—5 megabytes. The cost was \$10,000 per megabyte, or \$35,000 per year to lease. This compares with some of today’s magnetic hard drives that measure about $3\frac{1}{2}$ inches wide by 1 inch high, store as much as 1,000 gigabytes, and cost less than \$400, or about \$0.0004 per megabyte.

Those early, massive computers were housed in large, climate-controlled rooms. Within the room, a few knowledgeable experts, looking highly professional in their white laboratory coats, attended to the operation and maintenance of their million-dollar charges. The concept of a “user” as someone outside the computer room who could interact directly with the actual machine did not exist.

Service interruptions, software errors, and hardware errors were usually not critical. If any of these caused a program to fail or abort, beginning again was a relatively simple matter. Consequently, the primary security concerns were physical protection of the scarce and expensive hardware, and measures to increase their reliability. Another issue, then as now, was human fallibility. Because the earliest computers were programmed in extremely difficult machine language, consisting solely of ones (1s) and zeros (0s), the incidence of human error was high, and the time to correct errors was excessively long. Only later were assembler and compiler languages developed to increase the number of people able to program the machines and to reduce the incidence of errors and the time to correct them.

Information system security for large-scale computers was not a significant issue then for two reasons. First, only a few programming experts were able to utilize and manipulate computers. Second, there were very few computers in use, each of which was extremely valuable, important to its owners, and consequently closely guarded.

1.2.3 Medium-Size Computers. In the 1950s, smaller computer systems were developed with a very simple configuration; punched-card master files were replaced by punched paper tape and, later, by magnetic tape, and disk storage systems. The

1 · 6 BRIEF HISTORY AND MISSION OF INFORMATION SYSTEM SECURITY

electromechanical calculator with its patchboard was replaced by a central processor unit (CPU) that had a small main memory, sometimes as little as 8 kilobytes,⁴ and limited processing speed and power. One or two punched-card readers could read the data and instructions stored on that medium. Later, programs and data files were stored on magnetic tape. Output data were sent to cardpunches, for printing on unit record equipment, and later to magnetic tape. There was still no wired connection to the outside world, and there were no online users because no one, besides electronic data processing (EDP) people within the computer room, could interact directly with the system. These systems had very simple operating systems and did not use multiprocessing; they could run only one program at a time.

The IBM Model 650, as an example, introduced in 1954, measured about 5 feet by 3 feet by 6 feet and weighed almost 2,000 pounds. Its power supply was mounted in a similarly sized cabinet, weighing almost 3,000 pounds. It had 2,000 (10-digit) words of magnetic drum primary memory, with a total price of \$500,000 or a rental fee of \$3,200 per month. For an additional \$1,500 per month, a much faster core memory, of 60 words, could be added. Input and output both utilized read/write punch-card machines. The typical 1950s IS hardware was installed in a separate room, often with a viewing window so that visitors could admire the computer. In an early attempt at security, visitors actually within the computer room were often greeted by a printed sign saying:

Achtung! Alles Lookenspeepers!

Das computermachine ist nicht fur gefingerpoken und mittengrabben.
Ist easy schnappen der springenwerk, blowenfusen, und poppencorken mit spitzensparken.
Ist nicht fur gewerken bei das dumbkopfen.
Das rubbernecken sightseeren keepen hans in das pockets muss. Relaxen und watch das blinkenlichten.

Since there were still no online users, there were no user IDs and passwords. Programs processed batches of data, run at a regularly scheduled time—once a day, once a week, and so on, depending on the function. If the data for a program were not available at the scheduled run time, the operators might run some other job instead and wait for the missing data. As the printed output reports became available, they were delivered by hand to their end users. End users did not expect to get a continuous flow of data from the information processing system, and delays of even a day or more were not significant, except perhaps with paycheck production.

Information system security was hardly thought of as such. The focus was on batch controls for individual programs, physical access controls, and maintaining a proper environment for the reliable operation of the hardware.

1.2.4 1960s: Small-Scale Computers. During the 1960s, before the introduction of small-scale computers, dumb⁵ terminals provided users with a keyboard to send a character stream to the computer and a video screen that could display characters transmitted to it by the computer. Initially, these terminals were used to help computer operators control and monitor the job stream, while replacing banks of switches and indicator lights on the control console. However, it was soon recognized that these terminals could replace card readers and keypunch machines as well. Now users, identified by user IDs, and authenticated with passwords, could enter input data through a CRT terminal into an edit program, which would validate the input and then store it on a hard drive until it was needed for processing. Later, it was realized that users also could directly access data stored in online master files.

EVOLUTION OF INFORMATION SYSTEMS 1 · 7

1.2.5 Transistors and Core Memory. The IBM 1401, introduced in 1960 with a core memory of 4,096 characters, was the first all-transistor computer, marking the advent of the second generation. Housed in a cabinet measuring 5 feet by 3 feet, the 1401 required a similar cabinet to add an additional 12 kilobytes of main memory. Just one year later, the first integrated circuits were used in a computer, making possible all future advances in miniaturizing small-scale computers and in reducing the size of mainframes significantly.

1.2.6 Time Sharing. In 1961, the Compatible Time Sharing System (CTSS) was developed for the IBM 7090/7094. This operating system software, and its associated hardware, was the first to provide simultaneous remote access to a group of online users through multiprogramming.⁶ “Multiprogramming” means that more than one program can appear to execute at the same time. A master control program, usually called an operating system (OS), managed execution of the functional applications programs. For example, under the command of the operator, the OS would load and start application #1. After 50 milliseconds, the OS would interrupt the execution of application #1 and store its current state in memory. Then the OS would start application #2 and allow it to run for 50 milliseconds, and so on. Usually, within a second after users had entered keyboard data, the OS would give their applications a time slice to process the input. During each time slice, the computer might execute hundreds of instructions. These techniques made it appear as if the computer were entirely dedicated to each user’s program. This was true only so long as the number of simultaneous users was fairly small. After that, as the number grew, the response to each user slowed down.

1.2.7 Real-Time, Online Systems. Because of multiprogramming and the ability to store records online and accessible in random order, it became feasible to provide end users with direct access to data. For example, an airline reservation system stores a record of every seat on every flight for the next 12 months. A reservation clerk, working at a terminal, can answer a telephoned inquiry, search for an available seat on a particular flight, quote the fare, sell a ticket to the caller, and reserve the seat. Similarly, a bank officer can verify an account balance and effect money transfers. In both cases, each data record can be accessed and modified immediately, rather than having to wait for a batch to be run. Today, both the reservation clerk and the bank officer can be replaced by the customers themselves, who directly interface with the online computers.

While this advance led to a vast increase in available computing power, it also increased greatly the potential for breaches in computer security. With more complex operating systems, with many users online to sensitive programs, and with databases and other files available to them, protection had to be provided against inadvertent error and intentional abuse.

1.2.8 A Family of Computers. In 1964, IBM announced the S/360 family of computers, ranging from very small-scale to very large-scale models. All of the six models used integrated circuits, which marked the beginning of the third generation of computers. Where transistorized construction could permit up to 6,000 transistors per cubic foot, 30,000 integrated circuits could occupy the same volume. This lowered the costs substantially, and companies could buy into the family at a price within their means. Because all computers in the series used the same programming language and the same peripherals, companies could upgrade easily when necessary. The 360 family

1.8 BRIEF HISTORY AND MISSION OF INFORMATION SYSTEM SECURITY

quickly came to dominate the commercial and scientific markets. As these computers proliferated, so did the number of users, knowledgeable programmers, and technicians. Over the years, techniques and processes were developed to provide a high degree of security to these mainframe systems.

The year 1964 also saw the introduction of another computer with far-reaching influence: the Digital Equipment Corp. (DEC) PDP-8. The PDP-8 was the first mass-produced true minicomputer. Although its original application was in process control, the PDP-8 and its progeny quickly proved that commercial applications for minicomputers were virtually unlimited. Because these computers were not isolated in secure computer rooms but were distributed throughout many unguarded offices in widely dispersed locations, totally new risks arose, requiring innovative solutions.

1.2.9 1970s: Microprocessors. The foundations of all current personal computers (PCs) were laid in 1971 when Intel introduced the 4004 computer on a chip. Measuring 1/16 inch long by 1/8 inch high, the 4004 contained 2,250 transistors with a clock speed of 108 kiloHertz. The current generation of this earliest programmable microprocessor contains millions of transistors, with speeds over 1 gigaHertz, or more than 10,000 times faster. Introduction of microprocessor chips marked the fourth generation.

1.2.10 The First Personal Computers. Possibly the first personal computer was advertised in *Scientific American* in 1971. The KENBAK-1, priced at \$750, had three programming registers, five addressing modes, and 256 bytes of memory. Although not many were sold, the KENBACK-1 did increase public awareness of the possibility for home computers.

It was the MITS Altair 8800 that became the first personal computer to sell in substantial quantities. Like the KENBACK-1, the Altair 8800 had only 256 bytes of memory, but it was priced at \$375 without keyboard, display, or secondary memory. About one year later, the Apple II, designed by Steve Jobs and Steve Wozniak, was priced at \$1,298, including a CRT display and a keyboard.

Because these first personal computers were entirely stand-alone and usually under the control of a single individual, there were few security problems. However, in 1978, the VisiCalc spreadsheet program was developed. The advantages of standardized, inexpensive, widely used application programs were unquestionable, but packaged programs, as opposed to custom designs, opened the way for abuse because so many people understood their user interfaces as well as their inner workings.

1.2.11 The First Network. A national network, conceived in late 1969, was born as ARPANET⁷ (Advanced Research Projects Agency Network), a Department of Defense-sponsored effort to link a few of the country's important research universities, with two purposes: to develop experience in interconnecting computers and to increase productivity through resource sharing. This earliest connection of independent large-scale computer systems had just four nodes: the University of California at Los Angeles (UCLA), the University of California at Santa Barbara, Stanford Research Institute, and the University of Utah. Because of the inherent security in each leased-line interconnected node, and the physically protected mainframe computer rooms, there was no apparent concern for security issues. It was this simple network, with no thought of security designed in, from which evolved today's ubiquitous Internet and the World Wide Web (WWW), with their vast potential for security abuses.

EVOLUTION OF INFORMATION SYSTEMS 1 · 9

1.2.12 Further Security Considerations. With the proliferation of remote terminals on commercial computers, physical control over access to the computer room was no longer sufficient. In response to the new vulnerabilities, logical access control systems were developed. An access control system maintains an online table of authorized users. A typical user record would store the user's name, telephone number, employee number, and information about the data the user was authorized to access and the programs the user was authorized to execute. A user might be allowed to view, add, modify, and delete data records in different combinations for different programs.

At the same time, system managers recognized the value of being able to recover from a disaster that destroyed hardware and data. Data centers began to make regular tape copies of online files and software for offsite storage. Data center managers also began to develop and implement offsite disaster recovery plans, often involving the use of commercial disaster-recovery facilities. Even with such a system in place, new vulnerabilities were recognized throughout the following years, and these are the subjects of much of this *Handbook*.

1.2.13 The First “Worm”. A prophetic science-fiction novel, *The Shockwave Rider*, by John Brunner⁸ (1975), depicted a “worm” that grew continuously throughout a computer network. The worm eventually exceeded a billion bits in length and became impossible to kill without destroying the network. Although actual worms (e.g., the Morris Worm of 1988) later became real-and-present menaces to all networked computers, prudent computer security personnel install constantly updated antimalware programs that effectively kill viruses and worms without having to kill the network.

1.2.14 1980s: Productivity Enhancements. The decade of the 1980s might well be termed the era of productivity enhancement. The installation of millions of personal computers in commercial, industrial, and government applications enhanced efficiency and functionality of vast numbers of users. These advances, which could have been achieved in no other way, were made at costs that virtually any business could afford.

1.2.15 1980s: The Personal Computer. In 1981, IBM introduced a general-purpose small computer it called the “Personal Computer.” That model and similar systems became known generically as PCs. Until then, small computers were produced by relatively unknown sources, but IBM, with its worldwide reputation, brought PCs into the mainstream. The fact that IBM had demonstrated a belief in the viability of PCs made them serious contenders for corporate use.

There were many variations on the basic Model 5100 PC, and sales expanded far beyond IBM’s estimates. The basic configuration used the Intel 8088, operating at 4.77 megaHertz, with up to two floppy disk drives, each of 160 kilobytes capacity and with a disk-based operating system (DOS) in an open architecture. This open OS architecture, with its available “hooks,” made possible the growth of independent software producers, the most important of which was the Microsoft Corporation, formed by Bill Gates and Paul Allen.

IBM had arranged for Gates and Allen to create the DOS operating system. Under the agreement, IBM would not reimburse Gates and Allen for their development costs; rather, all profits from the sale of DOS would accrue to them. IBM did not have an exclusive right to the operating system, and Microsoft began selling it to many other customers as MS-DOS. IBM initially included with its computer the VisiCalc

1 · 10 BRIEF HISTORY AND MISSION OF INFORMATION SYSTEM SECURITY

spreadsheet program, but soon sales of Lotus 1-2-3 surpassed those of VisiCalc. The open architecture not only made it possible for many developers to produce software that would run on the PC, but also enabled anyone to put together purchased components into a computer that would compete with IBM's PC. The rapid growth of compatible application programs, coupled with the ready availability of compatible hardware, soon resulted in sales of more than 1 million units. Many subsequent generations of the original hardware and software are still producing sales measured in millions every year.

Apple took a very different approach with its Macintosh computer. Where IBM's system was wide open, Apple maintained tight control over any hardware or software designed to operate on the Macintosh so as to assure compatibility and ease of installation. The most important Apple innovations were the graphical user interface (GUI) and the mouse, both of which worked together to facilitate ease of use and both of which were derived from research and development at the Stanford Research Institute and the Xerox Palo Alto Research Institute in the 1960s and 1970s. Microsoft had attempted in 1985 to build these features into the Windows operating system, but early versions were generally rejected as slow, cumbersome, and unreliable. It was not until 1990 that Windows 3.0 overcame many of its problems and provided the foundation for later versions that were almost universally accepted.

1.2.16 Local Area Networks. During the 1980s, stand-alone desktop computers began to perform word processing, financial analysis, and graphic processing. Although this arrangement was much more convenient for end users than was a centralized facility, it was more difficult to share data with others.

As more powerful PCs were developed, it became practical to interconnect them so that their users could easily share data. These arrangements were commonly referred to as local area networks (LANs) because the hardware units were physically close, usually in the same building or office area. LANs have remained important to this day. Typically, a more powerful PC with a high storage capacity fixed⁹ disk was designated as the file server. Other PCs, referred to as workstations, were connected to the file server using network interface cards installed in the workstations with cables between these cards and the file server. Special network software installed on the file server and workstations made it possible for workstations to access defined portions of the file server fixed disk just as if these portions were installed on the workstations. Furthermore, these shared files could be backed up at the file server without depending on individual users. By 1997, it was estimated that worldwide there were more than 150 million PCs operating as LAN workstations. The most common network operating systems (NOS) were Novell NetWare and later Microsoft IE (Internet Explorer).

Most LANs were implemented using the Ethernet (IEEE 802.3) protocol.¹⁰ The server and workstations could be equipped with a modem (modulator/demodulator) connected to a dedicated telephone line. The modem enabled remote users, with a matching modem, to dial into the LAN and log on. This was a great convenience to LAN users who were traveling or working away from their offices, but such remote access created yet another new security issue. For the first time, computer systems were exposed in a major way to the outside world. From then on, it was possible to interact with a computer from virtually anywhere and from locations not under the same physical control as the computers themselves.

Typical NOS logical access control software provided for user IDs and passwords and selective authority to access file server data and program files. A workstation user logged on to the LAN by executing a log-in program resident on the file server.

EVOLUTION OF INFORMATION SYSTEMS 1 · 11

The program prompted the user to enter an ID and password. If the log-in program concluded that the ID and password were valid, it consulted an access-control table to determine which data and programs the user might access. Access modes were defined as read-only, execute-only, create, modify (write or append), lock, and delete, with respect to individual files and groups of files. The LAN administrator maintained the access control table using a utility program. The effectiveness of the controls depended on the care taken by the administrator, and so, in some circumstances, controls could be weak. It was essential to protect the ID and password of the LAN administrator since, if they were compromised, the entire access-control system became vulnerable. Alert information system security officers noted that control over *physical* access to LAN servers was critical in maintaining the logical access controls. Intruders who could physically access a LAN server could easily restart the server using their own version of the NOS, completely bypassing the installed logical access controls.

Superficially, a LAN appears to be the same as a 1970s mainframe with remote dumb terminals. The difference technically is that each LAN workstation user is executing programs on the workstation, not on the centralized file server, while mainframe computers use special software and hardware to run many programs concurrently, one program for each terminal. To the user at a workstation or remote terminal, the two situations appear to be the same, but from a security standpoint, there are significant differences. The mainframe program software stays on the mainframe and cannot, under normal conditions, be altered during execution. A LAN program on a workstation can be altered, for example, by a computer virus, while actually executing. As a rule, mainframe remote terminals cannot download and save files, whereas workstations usually have at least a removable disk drive. Furthermore, a malicious workstation user can easily install a rewritable CD device, which makes it much easier to copy and take away large amounts of data.

Another important difference is the character of the connection between the computer and the terminals. Each dumb terminal has a dedicated connection to its mainframe and receives only that data that is directed to it. A LAN operates more like a set of radio transmitters sharing a common frequency on which the file server and the workstations take turns “broadcasting” messages. Each message includes a “header” block that identifies the intended recipient, but every node (the file server and the workstations) on a LAN receives all messages. Under normal circumstances, each node ignores messages not addressed to it. However, it is technically feasible for a workstation to run a modified version of the NOS that allows it to capture all messages. In this way, a workstation could identify all log-in messages and record the user IDs and passwords of all other users on the LAN, giving it complete access to all of the LAN’s data and facilities.

Mainframe and LAN security also differ greatly in the operating environment. As noted, the typical mainframe is installed in a separate room and is managed by a staff of skilled technicians. The typical LAN file server, on the other hand, is installed in ordinary office space and is managed by a part-time, remotely located LAN administrator who may not be adequately trained. Consequently, the typical LAN has a higher exposure to tampering, sabotage, and theft. However, if the typical mainframe is disabled by an accident, fire, sabotage, or any other security incident, many business functions will be interrupted, whereas the loss of a LAN file server usually disrupts only a single function.

1.2.17 1990s: Interconnection. The Usenet evolved in the early 1980s as a free system for posting and retrieving news and commentary from participants—an

1 · 12 BRIEF HISTORY AND MISSION OF INFORMATION SYSTEM SECURITY

early form of disintermediation, since there were no controlling authorities to limit speech. Newsgroups developed on every conceivable topic, reaching tens of thousands of discussion areas within a few years. Computer enthusiasts and criminal hackers took to Usenet as an ideal channel for exchanging code, including details of hacks.

The commercial equivalents of the Usenet were the value-added networks (VANs) such as America On Line (AOL), CompuServe, and Prodigy. These services provided modems for telephone access, email, and facilities for defining discussion groups. Fees varied from hourly to monthly.

1.2.18 1990s: Total Interconnection. With the growing popularity of LANs, the technologies for interconnecting them emerged. These networks of physically interconnected local area networks were called wide area networks, or WANs. Any node on a LAN could access every node on any other interconnected LAN, and in some configurations, those nodes might also be given access to mainframe and minicomputer files and to processing capabilities.

1.2.19 Telecommuting. Once the WAN technology was in place, it became feasible to link LANs together by means of telecommunications circuits. It had been expensive to do this with the low-speed, online systems of the 1970s because all data had to be transmitted over the network. Now, since processing and most data used by a workstation were on its local LAN, a WAN network was much less expensive. Low-traffic LANs were linked using dial-up access for minimum costs, while major LANs were linked with high-speed dedicated circuits for better performance. Apart from dial-up access, all network traffic typically flowed over nonswitched private networks. Of the two methods, dial-up communications were considerably more vulnerable to security violations, and they remain so to this day.

1.2.20 Internet and the World Wide Web. The Internet, which began life in 1969 as the Advanced Research Projects Agency Network (ARPANET), slowly emerged onto the general computing scene during the 1980s. Initially, access to the Internet was restricted to U.S. Government agencies and their contractors. ARPANET users introduced the concept of email as a convenient way to communicate and exchange documents. Then, in 1989–1990, Sir Tim Berners-Lee conceived of the World Wide Web and the Web browser. This one concept produced a profound change in the Internet, greatly expanding its utility and creating an irresistible demand for access. During the 1990s, the U.S. Government relinquished its control, and the Internet became the gigantic, no-one-is-in-charge network of networks it is today. The explosive growth in participation in the global Internet is generally viewed as having started with the opening up of the .COM top-level domain to general use in 1993.

The Internet offers several important advantages: The cost is relatively low, connections are available locally in most industrialized countries, and by adopting the Internet protocol, TCP/IP, any computer becomes instantly compatible with all other Internet users.

The World Wide Web technology made it easy for anyone to access remote data. Almost overnight, the Internet became the key to global networking. Internet service providers (ISPs) operate Internet-compatible computers with both dial-up and dedicated access. A computer may access an ISP directly as a stand-alone ISP client or via a gateway from a LAN or WAN. A large ISP may offer dial-up access at many locations, sometimes called points of presence or POPs, interconnected by its own network. ISPs establish links with one another through the national access points (NAPs) initially set

GOVERNMENT RECOGNITION OF INFORMATION ASSURANCE 1 · 13

up by the National Science Foundation. With this “backbone” in place, any node with access can communicate with another node, connected to a different ISP, located half way around the globe, without making prior arrangements.

The unrestricted access provided by the Internet created new opportunities for organizations to communicate with clients. A company can implement a Web server with a full-time connection to an ISP and open the Web server, and the WWW pages it hosts, to the public. A potential customer can access a Website, download product information and software updates, ask questions, and even order products. Commercial Websites, as they evolved from static “brochure-ware” to online shopping centers, stock brokerages, and travel agencies, to name just a few of the uses, became known as e-businesses.

1.2.21 Virtualization and the Cloud. As far back as the late 1960s, software was available to create encapsulated versions of an operating system on mainframe computers. Users interacted with what appeared to be their own, private mainframe environment. By the late 1980s, vendors created simulations of operating environments that could run under different operating systems (e.g., one could run DOS programs on UNIX machines). The trend continued throughout the succeeding years so that it is commonplace now to run programs under hypervisors that simulate complete or functionally limited versions of required operating systems on shared hardware.¹¹

Today it is possible to provide users with instances of an operating environment on shared hardware, often at a distance, so that incremental increases in requirements can be satisfied at modest costs instead of having to purchase large-scale improvements in the hardware infrastructure. The situation is similar to what service bureaus offered in the decades when mainframe time-sharing was popular.

Another development in the last decade has been the availability of cloud computing, which refers to computer services, including storage (see Chapters 36 and 68), software as a service (SAAS), and infrastructure or platform as a service (IAAS and PAAS). See Chapter 68 for more details of managing and securing cloud computing.

1.2.22 Supervisory Control and Data Acquisition. The use of computers to control production of goods and services through supervisory control and data acquisition (SCADA) software and hardware has been growing throughout the four decades since this *Handbook* was first published in 1973. SCADA systems for critical infrastructure have been of great concern because contrary to initial design specifications, many of them have been connected to the general Internet, opening the systems they govern to subversion. For more about SCADA in information warfare, see Chapter 14 in this *Handbook*.

1.3 GOVERNMENT RECOGNITION OF INFORMATION ASSURANCE. Certain major events in the history of information assurance (IA) center on government initiatives. In particular, IA has been strongly influenced by the development of security standards starting in the 1980s, by the publication of the landmark publication *Computers and Risk* in 1991, and by the establishment of the InfraGard program in the late 1990s for protection of the U.S. critical infrastructure.

1.3.1 IA Standards. In the late 1970s, the U.S. Department of Defense “established a Computer Security Initiative to foster the wide-spread availability of trusted

1 · 14 BRIEF HISTORY AND MISSION OF INFORMATION SYSTEM SECURITY

computer systems.”¹² The author of the initial report that later became the *Trusted Computer Systems Evaluation Criteria* (TCSEC), DoD Standard 5200.28, wrote,

Trusted computer systems are operating systems capable of preventing users from accessing more information than that to which they are authorized. Such systems are in great demand as more processing is entrusted to computers, while less information should be shared by all the system’s users. With this demand comes a need to ascertain the integrity of computer systems on the market ...

The TCSEC was issued with a bright orange cover and became known as the “Orange Book.” Under the direction of National Computer Security Center (NCSC) director Patrick Gallagher and others, the National Security Agency (NSA) issued a “Rainbow Series” of books that profoundly affected the direction of IA in the USA and globally.¹³

The Rainbow Series led to similar efforts in other countries, culminating in the Common Criteria Evaluation and Validation Scheme (CCEVS), which has become the international standard for defining security levels for systems and software and for determining acceptable methods for testing and certifying system compliance with such standards.¹⁴

For details of the evolution of security standards, see Chapter 51 in this *Handbook*.

1.3.2 Computers at Risk.¹⁵ In 1988, the Defense Advanced Research Projects Agency (DARPA) asked the Computer Science and Technology Board (renamed the Computer Science and Telecommunications Board of the NRC in 1990) for a study of computer and communications security issues affecting U.S. Government and industry. The NRC’s System Security Study Committee published its results in a readable and informative book, *Computers at Risk: Safe Computing in the Information Age*.¹⁶

The Committee included experts with impeccable credentials, including executives from major computer vendors such as HP, DEC, and IBM; from high-technology companies such as Shearson, Lehman, Hutton Inc., and Rockwell International; universities such as Harvard and MIT; and think tanks like the RAND Corporation.

A public misconception is the supposed divergence in focus of the military and of commerce: The military is usually described as concerned with external threats and the problem of disclosure, whereas businesses are said to worry more about insider threats to data integrity. On the contrary, the military and commerce need to protect data in similar ways. The differences arise primarily from (1) the sophistication and resources available to governments that try to crack foreign military systems; (2) the relatively strong military emphasis on prevention compared with commercial need for proof that can be used in legal proceedings; and (3) the availability to the military of deep background checks on personnel, contrasted with the limits imposed on the invasion of privacy in the commercial sector.

Some of the more interesting points raised by the NRC Committee assert that:

- Because of the rapid and discontinuous pace of innovation in the computer field, “with respect to computer security, the past is not a good predictor of the future”;
- Embedded systems (those where the microprocessor is not accessible to reprogramming by the user; e.g., medical imaging systems) open us to greater risks from inadequate quality assurance (e.g., a software bug in a Therac 25 linear accelerator killed three patients by irradiating them with more than 100 times the intended radiation dosage);

GOVERNMENT RECOGNITION OF INFORMATION ASSURANCE 1 · 15

- Networking makes it possible to harm many more systems: “Interconnection gives an almost ecological flavor to security; it creates dependencies that can harm as well as benefit the community ...”

The Committee proposed major recommendations, summarized as follows:

1. Push for implementation of generally accepted system security principles including:
 - Quality assurance standards that address security considerations;
 - Access control for operations as well as data (e.g., any of the menu systems which preclude access to the operating system);
 - Unambiguous user identification (ID) and authentication (e.g., personal profiles and hand-held password generators)
 - Protection of executable code (e.g., flags to show that certain object modules are “production” or “installed” and thus apply strict access control that would prevent unauthorized modification—as found in configuration control systems);
 - Security logging (e.g., logging failed file-open attempts and logon password violations);
 - Assigning a security administrator to each enterprise;
 - Data encryption;
 - Operational support tools for verifying the state and effectiveness of security measures (e.g., audit tools);
 - Independent audits of system security by people not directly involved in programming or system management of the audited system;
 - Hazard analysis evaluating threats to safety from different malfunctions and breaches of security (e.g., consequences of tampering with patient data in hospitals).
2. Take specific short-term actions now:
 - Develop security policies for your organization before there’s a problem;
 - Form and train computer emergency response teams before a crisis to respond to security violations or attacks;
 - Use the Orange Book’s (TCSEC, from the National Computer Security Center’s Rainbow series) C2 and B1 criteria to define guidelines on security;
 - Improve software systems development by applying better quality-assurance methods;
 - Contribute to voluntary industry groups developing modern security standards and implement those standards in commercial software;
 - Make effective security the default in software and hardware (make the user explicitly disable security instead of having to enable it).
3. Learn and teach about security:
 - Build a repository of incident data;
 - Foster education in engineering secure systems, both by encouraging universities to provide postgraduate training in security and by urging industry to include security training as part of software engineering projects;

1 · 16 BRIEF HISTORY AND MISSION OF INFORMATION SYSTEM SECURITY

- Teach beginners about security and ethics in computer usage and programming (e.g., the NCSA is working on a research and development project to study beliefs, attitudes, and behavior about ethical issues in computing in grade and high schools, colleges, and universities).
- 4. Clarify export control criteria and set up a forum for arbitration (hardware and software vendors have been complaining for years that the arbitrary imposition of severe export restrictions hampers American competitiveness in overseas markets without materially helping national security).
- 5. Fund and pursue needed research in such areas as:
 - Security modularity: the effects on security of combining modules with known security properties;
 - Security policy models: more subtle requirements like integrity and availability are still not easily represented by control structures;
 - Cost estimation: there should be better ways of measuring the costs and benefits of security mechanisms in particular applications;
 - New technology: networking, in particular, leads to greater complexity (e.g., how to connect “mutually suspicious organizations”);
 - Quality assurance for security: how to measure effectiveness;
 - Modeling tools: standards for graphical representations of security relationships analogous to the diagrams used in functional decomposition and object-oriented methodologies for program design;
 - Automated procedures: audit and monitoring tools for the data center management team;
 - Nonrepudiation: combining the need for detailed records of user actions with the values of privacy;
 - Resource control: how to ensure that proprietary software and data are used legitimately (e.g., preventing more than the licensed number of users from accessing a system, preventing software theft);
 - Security perimeters: how to reconcile the desire for network interconnection with limitations due to security requirements (“If, for example, a network permits mail but not directory services ... less mail may be sent because no capability exists to look up the address of a recipient”).

Chapter 2 of the NRC report, Concepts of Information Security, is a 25-page primer on information systems security that could be handed to any manager who needs to be filled in on why you propose to spend so much money protecting the computer systems. The authors cover the fundamental aspects of information security (confidentiality, integrity, and availability); management controls (individual accountability, auditing, and separation of duties); risks (probabilities of attack or damage) and vulnerabilities (weak points); and privacy issues. In Appendix 2.2, the authors report an informal survey in April 1989 of 30 private companies in a variety of fields. The consensus among those polled included the following basic standards for information systems security (show these to your upper management if necessary):

- Unique IDs, block access after a maximum number of incorrect logon attempts, show last successful access at logon time, make passwords and IDs expire;

GOVERNMENT RECOGNITION OF INFORMATION ASSURANCE 1 · 17

- Disallow embedded passwords during logon, make passwords invisible during entry, force minimum length (6), store passwords encrypted, scan proposed passwords to eliminate easy words;
- Permit strict control over file access;
- Detect and interdict viruses, certify software as virus-free, provide data encryption, overwrite deleted files to prevent recovery, force tight binding of production data to production programs;
- Automated time-out for inactive sessions, unique identification of terminals and workstations during logon;
- Network security monitoring, modem locking, callback, automatic data encryption during transmission;
- Audit trails, including security violations;
- Generally applicable security standards that could be used by vendors and users to evaluate different equipment and software for specific environments.

Twenty years later, focus among information assurance experts has shifted beyond the technical to emphasize organizational controls. For example, the 2003 survey of members of the Information Systems Security Association included these information security function practices among the respondents:

- Access controls: 73%
- Written information security policy: 72%
- Compliance with existing laws and regulations: 66%
- Creation of organization and process to implement policy: 59%
- Awareness and training program: 57%
- Regular monitoring, reviewing, and auditing: 57%
- Business continuity planning: 57%
- Risk assessment and risk management: 56%

In 2007, Gary S. Miliefsky proposed the following seven priorities for corporate information security:

- 1. Policies**
- 2. Awareness and training**
- 3. Information security self-assessments**
- 4. Regulatory compliance self-assessments**
- 5. Corporate-wide encryption**
- 6. Manage all corporate assets**
- 7. Test Business Continuity Planning (BCP) and Disaster Recovery Planning (DRP)¹⁷**

The Computer Security Division of the Information Technology Laboratory at the National Institute of Standards and Technology issued a draft reference model that

1 · 18 BRIEF HISTORY AND MISSION OF INFORMATION SYSTEM SECURITY

included the following “programmatic, integration, and system security activities that are typically a part of an information security program”:

- Program Security Activities
 - Annual and Quarterly Review and Reporting of Information Security Program
 - Asset Inventory
 - Awareness and Specialized Security Training
 - Continuity of Operations
 - Incident Response
 - Periodic Testing and Evaluation
 - Plan of Action and Milestones
 - Policies and Procedures
 - Risk Management
- Integration Activities
 - Business Risk
 - Capital Planning and Investment Control (CPIC)
 - Configuration Management
 - Enterprise Architecture (EA)
 - Environmental Protection
 - Human Resources
 - Personnel Security
 - Physical Security
 - Privacy
 - Records Management
 - Strategic Plan
 - System Development Life Cycle (SDLC)
- System Security Activities
 - Categorize the Information System
 - Select Security Controls
 - Supplement Security Controls
 - Document Security Controls
 - Implement Security Controls
 - Assess Security Controls
 - Authorize the Information System
 - Monitor Security Controls

1.3.3 InfraGard.¹⁸ InfraGard is a nationwide program in the United States that brings together representatives from information technology departments in industry and academia for information sharing and analysis, especially to help protect critical infrastructure against cyberattacks and also to support the FBI in its cybercrime investigations and education projects.¹⁹

RECENT DEVELOPMENTS 1 · 19

The organization started in the Cleveland Field Office of the FBI in 1996 and expanded rapidly until there are now over 11,000 members in over 40 chapters. Joining InfraGard is easy and free for U.S. citizens residing in the United States. Using the Website (www.infragard.org), you can locate a nearby local chapter (“Find Chapters”) and contact your chapter officers. You can get application forms online and then send them in to the FBI liaison officer for that chapter to be vetted for admission. The FBI conducts a background check to ensure that all members are likely to be trustworthy to participate in confidential discussions of threats and vulnerabilities. Chapters usually conduct regular local meetings and organize list-servers for exchange of information among members. Many have newsletters as well.

1.4 RECENT DEVELOPMENTS. In recent years, a key development has been the dramatic increase in availability of inexpensive portable data storage devices. At the time of writing (2013), flash drives the size of a lipstick or even of an antacid pill are available with capacities in the dozens of gigabytes for a few dollars. Such devices are available in a wide range of concealable formats such as pens, music players, watches, and (no joke) sushi. Pocket-sized hard disks and solid-state drives with capacities in the hundreds of gigabytes to terabytes are available for less than US\$100. Digital cameras use storage cards that can be used for data transfers; mobile phones include cameras and recording capabilities. A 64 GB micro SD card for a phone costs about \$50 and can hold 6,000 songs from iTunes—or the entire customer database being stolen by a disaffected soon-to-be-fired employee. Controlling data leakage through unauthorized connection of such devices has become a significant problem for security managers. Systems for restricting connection of devices and controlling data transfers to such storage media (data-loss prevention or DLP) are spreading through government and corporate environments (see Chapter 13 in this *Handbook* for a detailed discussion of DLP).

Another issue that increasingly concerns security managers is the protection of personally identifiable information (PII) from customers or data subjects. Many organizations including government agencies, banks, and universities have suffered serious damage from loss of control over PII and the risks of identity theft resulting from exposure of such sensitive data. Legislators are responding to public concern by increasing legal requirements for protection of PII. The use of encryption on mobile data systems such as laptop computers, personal digital assistants (PDAs), mobile phones, and integrated systems that combine many functions (e.g., BlackBerries) has become a necessity. See Chapter 69 in this *Handbook* for extensive discussion of protection of PII.

A consequence of the growing interconnectivity of storage and communications devices is that corporate networks are no longer insulated from less-secure systems. Users who connect poorly protected laptops (or other devices) to public networks such as hotel-supplied ISPs or wireless access points in coffee shops may return to their home offices with malware-infected systems that contaminate the entire network. Security managers are increasingly turning to integrated systems for controlling connectivity via virtual private networks and supervisory software that monitors and restricts unauthorized connections, software installations, and downloads.

Another growing issue is the increasing speed and persistence of attacks on systems and networks by state-sponsored and criminal organizations engaged in industrial espionage and fraud. See Chapters 2, 14, 15, and 46, among many others in this *Handbook*, for further discussion of the changing threat profile for today’s information systems.

1 · 20 BRIEF HISTORY AND MISSION OF INFORMATION SYSTEM SECURITY

1.5 ONGOING MISSION FOR INFORMATION SYSTEM SECURITY.

There is no end in sight to the continuing proliferation of Internet nodes, to the variety of applications, to the number and value of online transactions, and, in fact, to the rapid integration of computers into virtually every facet of our existence. Nor will there be any restrictions as to time or place. With 24/7/365, always-on operation, and with global expansion even to relatively undeveloped lands, both the beneficial effects and the security violations can be expected to grow apace.

Convergence, which implies computers, televisions, cell phones, and other means of communication combined in one unit, together with continued growth of information technology, will lead to unexpected security risks. Distributed denial-of-service (DDoS) attacks, copyright infringement, child pornography, fraud, theft of identity, and industrial espionage are all ongoing security threats. So far, no perfect defensive measures have been developed.

The situation is currently (2013) changing from identifying vulnerabilities and preventing penetrations to identifying compromises and minimizing damage from long-lasting subversion of protection mechanisms. Situational awareness and rapid response are becoming an increasingly important element in long-term defenses of our information.

This *Handbook* provides a foundation for understanding and blunting both the existing vulnerabilities and those new threats that will inevitably arise in the future. Certainly, no one but the perpetrators could have foreseen the use of human-guided missiles to attack the World Trade Center. Besides its symbolic significance, the great concentration of resources within the WTC increased its attractiveness as a target. After 9/11, the importance of physical safety of personnel has become the dominant security issue, with disaster recovery of secondary, but still great, concern. This *Handbook* cannot foresee all possible future emergencies, but it does prescribe some preventative measures, and it does recommend procedures and resources for mitigation and remediation.

1.6 NOTES

1. Many technical specialists tend to use the term “security” to refer to logical access controls. A glance at the contents pages of this volume shows the much broader scope of information system security.
2. For further details, see, for example, www.cs.uiowa.edu/~jones/cards.
3. See <http://ei.cs.uiowa.edu/~history/UNIVAC.Weston.html> and inventors.about.com/library/weekly/aa062398.htm
4. It is notable that the IBM 1401 computer was so named because the initial model had 1,400 bytes of main memory. It was not long before memory size was raised to 8 kilobytes and then later to as much as 32 kilobytes. In 1980, the Series III minicomputer from Hewlett-Packard doubled its maximum memory from 1 megabyte to 2 megabytes at a cost of \$64,000 (about \$200,000 in 2008 dollars). This compares with today’s personal computers, typically equipped with no less than 512 megabytes and often a gigabyte or more.
5. The term “dumb” was used because the terminal had no internal storage or processing capability. It could only receive and display characters and accept and transmit keystrokes. Both the received characters and the transmitted ones were displayed on a cathode ray tube (CRT) much like a pre-color television screen. Consequently, these were also called “glass” terminals.

NOTES 1 · 21

6. “Multiprocessing,” “multiprogramming,” and “multitasking” are terms that are used almost interchangeably today. Originally, multitasking implied that several modules or subroutines of a single program could execute together. Multiprogramming was designed to execute several different programs, and their subroutines, concurrently. Multiprocessing most often meant that two or more computers worked together to speed program execution by providing more resources.
7. Also known as ARPAnet and Arpanet.
8. First published in 1975; reissued by Mass Market Paperbacks in May 1990.
9. “Fixed,” in contrast with the removable disk packs common in large data centers.
10. See standards.ieee.org/getieee802/802.3.html
11. Sean P. Conroy, “History of Virtualization,” *Everything VM*, 2010. www.everythingvm.com/content/history-virtualization
12. G. H. Nibaldi, Proposed Technical Evaluation Criteria for Trusted Computer Systems. Publication M79-225 (Bedford, MA: MITRE Corporation, 1999).
13. For access to all the Rainbow Series documents, see www.fas.org/irp/nsa/rainbow.htm
14. The CCEVS Website has extensive documentation; see www.niap-ccevs.org
15. This section is reprinted with slight modifications by permission of the author from the original manuscript for M. E. Kabay, *The NCSA Guide to Enterprise Security: Protecting Information Assets* (New York: McGraw-Hill, 1996), Chapter 1, pp. 2–5.
16. National Research Council, *Computers at Risk: Safe Computing in the Information Age* (Washington, DC: National Academy Press, 1991). Available as searchable openbook at www.nap.edu/openbook.php?isbn=0309043883
17. G. S. Miliefsky, “The 7 Best Practices for Network Security in 2007,” NetworkWorld Website, 2007, www.networkworld.com/columnists/2007/011707miliefsky.html?t51hb
18. M. E. Kabay (2005). “InfraGard is not a Deodorant,” NetworkWorld Website, 2005, www.networkworld.com/newsletters/sec/2005/0905sec2.html
19. www.infragard.org

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 2

HISTORY OF COMPUTER CRIME

M.E. Kabay

2.1 WHY STUDY HISTORICAL RECORDS?	2·2	2.10.1 1988 Flu-Shot Hoax	2·12
		2.10.2 Scrambler, 12-Tricks,	
		and PC Cyborg	2·12
2.2 OVERVIEW	2·2	2.10.3 1994: Datacomp	
		Hardware Trojan	2·13
2.3 1960S AND 1970S: SABOTAGE	2·3	2.10.4 Keylogger Trojans	2·13
2.3.1 Direct Damage to Computer Centers	2·3	2.10.5 Haephrati Trojan	2·14
2.3.2 1970–1972: Albert the Saboteur	2·4	2.10.6 Hardware Trojans and Information Warfare	2·14
		2.11 NOTORIOUS WORMS AND VIRUSES	2·15
2.4 IMPERSONATION	2·5	2.11.1 1970–1990: Early Malware Outbreaks	2·15
2.4.1 1970: Jerry Neal Schneider	2·5	2.11.2 December 1987: Christmas Tree Worm	2·16
2.4.2 1980–2003: Kevin Mitnick	2·5	2.11.3 November 2, 1988: Morris Worm	2·16
2.4.3 Credit Card Fraud	2·6	2.11.4 Malware in the 1990s	2·17
2.4.4 Identity Theft Rises	2·7	2.11.5 March 1999: Melissa	2·18
		2.11.6 May 2000: I Love You	2·20
		2.11.7 July 2010 Stuxnet	2·20
2.5 PHONE PHREAKING	2·8		
2.5.1 2600 Hz	2·8		
2.5.2 1982–1991: Kevin Poulsen	2·8		
2.6 DATA DIDDLING	2·9		
2.6.1 Equity Funding Fraud (1964–1973)	2·9	2.12 SPAM	2·20
2.6.2 1994: Vladimir Levin and the Citibank Heist	2·10	2.12.1 1994: Green Card Lottery Spam	2·20
		2.12.2 Spam Goes Global	2·21
2.7 SALAMI FRAUD	2·10	2.13 DENIAL OF SERVICE	2·21
		2.13.1 1996: Unemailer	2·21
2.8 LOGIC BOMBS	2·11	2.13.2 2000: MafiaBoy	2·22
2.9 EXTORTION	2·12	2.14 HACKER UNDERGROUND	2·22
		2.14.1 1981: Chaos Computer Club	2·23
2.10 TROJAN HORSES	2·12		

2 · 2 HISTORY OF COMPUTER CRIME

2.14.2	1982: The 414s	2·23	2.14.11	2004: Shadowcrew	2·27
2.14.3	1984: Cult of the Dead Cow	2·23	2.14.12	Late 2000s: Russian Business Network (RBN)	2·27
2.14.4	1984: <i>2600: The Hacker Quarterly</i>	2·24	2.14.13	Anonymous	2·28
2.14.5	1984: Legion of Doom	2·24	2.14.14	2013: Unlimited Operations	2·28
2.14.6	1985: <i>Phrack</i>	2·25			
2.14.7	1989: Masters of Deception	2·25	2.15 INDUSTRIAL ESPIONAGE	2·29	
2.14.8	1990: Operation Sundevil	2·26	2.16 CONCLUDING REMARKS	2·31	
2.14.9	1990: Steve Jackson Games	2·26	2.17 FURTHER READING	2·32	
2.14.10	1992: L0pht Heavy Industries	2·27	2.18 NOTES	2·33	

2.1 WHY STUDY HISTORICAL RECORDS? Every field of study and expertise develops a common body of knowledge that distinguishes professionals from amateurs. One element of that body of knowledge is a shared history of significant events that have shaped the development of the field. Newcomers to the field benefit from learning the names and significant events associated with their field so that they can understand references from more senior people in the profession, and so that they can put new events and patterns into perspective. This chapter provides a brief overview of some of the more famous (or notorious) cases of computer crime (including those targeting computers and those mediated through computers) of the last four decades.¹

2.2 OVERVIEW. This chapter illustrates several general trends from the 1960s through mid-2013:

- In the early decades of modern information technology (IT), computer crimes were largely committed by individual disgruntled and dishonest employees.
- Physical damage to computer systems was a prominent threat until the 1980s.
- Criminals often used authorized access to subvert security systems as they modified data for financial gain or destroyed data for revenge.
- Early attacks on telecommunications systems in the 1960s led to subversion of the long-distance phone systems for amusement and for theft of services.
- As telecommunications technology spread throughout the IT world, hobbyists with criminal tendencies learned to penetrate systems and networks.
- Programmers in the 1980s began writing malicious software, including self-replicating programs, to interfere with personal computers.
- As the Internet increased access to increasing numbers of systems worldwide, criminals used unauthorized access to poorly protected systems for vandalism, political action, and financial gain.
- As the 1990s progressed, financial crime using penetration and subversion of computer systems increased.
- The types of malware shifted during the 1990s, taking advantage of new vulnerabilities and dying out as operating systems were strengthened, only to succumb to new attack vectors.

1960s AND 1970s: SABOTAGE 2 · 3

- Illegitimate applications of email grew rapidly from the mid-1990s onward, generating torrents of unsolicited commercial and fraudulent email.
- Organized crime became increasingly involved in systematic penetration of financial systems and targeted fraud.
- Chinese government-supported civilian and military agents increasingly used computer-based industrial espionage to gain significant economic advantages over industry and commerce in North America and Europe.

2.3 1960s AND 1970s: SABOTAGE. Early computer crimes often involved physical damage to computer systems and subversion of the long-distance telephone networks.

2.3.1 Direct Damage to Computer Centers. In February 1969, the largest student riot in Canada was set off when police were called in to put an end to a student occupation of several floors of the Hall Building. The students had been protesting against a professor accused of racism, and when the police came in, a fire broke out and computer data and university property were destroyed. The damages totaled \$2 million, and 97 people were arrested.²

Thomas Whiteside cataloged a litany of early physical attacks on computer systems in the 1960s and 1970s³:

1968	Olympia, WA: An IBM 1401 in the state is shot twice by a pistol-toting intruder
1970	University of Wisconsin: Bomb kills one and injures three people and destroys \$16 million of computer data stored on site
1970	Fresno State College: Molotov cocktail causes \$1 million damage to computer system
1970	New York University: Radical students place fire-bombs on top of Atomic Energy Commission computer in attempt to free a jailed Black Panther
1972	Johannesburg, South Africa: Municipal computer is dented by four bullets fired through a window
1972	New York: Magnetic core in Honeywell computer attacked by someone with a sharp instrument, causing \$589,000 of damage
1973	Melbourne, Australia: Antiwar protesters shoot American firm's computer with double-barreled shotgun
1974	Charlotte, NC: Charlotte Liberty Mutual Life Insurance Company computer is shot by a frustrated operator
1974	Dayton, OH: Wright Patterson Air Force Base: Four attempts are made to sabotage computers, including by magnets, loosened wires, and gouges in equipment
1977	Rome, Italy: Four terrorists pour gasoline on university computer and burn it to cinders
1978	Lompoc, CA: Vandenburg Air Force Base: A peace activist destroys an unused IBM 3031 using a hammer, a crowbar, a bolt cutter, and a cordless power drill as a protest against the NAVSTAR satellite navigation system, claiming it gives the United States a first-strike capability

The incidents of physical abuse of computer systems did not stop as other forms of computer crime increased. For example, in 2001, *NewsScan* editors⁴ summarized a report from *Wired Magazine*:

A survey by British PC maker Novatech, intended to take a lighthearted look at techno-glitches, instead revealed the darker side of computing. One in every four computers has been physically assaulted by its owner, according to the 4,200 respondents.⁵

2 · 4 HISTORY OF COMPUTER CRIME

In April 2003, the National Information Protection Center and Department of Homeland Security reported:

Nothing brings a network to a halt more easily and quickly than physical damage. Yet as data transmission becomes the lifeblood of Corporate America, most big companies haven't performed due diligence to determine how damage-proof their data lifelines really are. Only 20 percent of midsize and large companies have seriously sussed out what happens to their data connections after they go beyond the company firewall, says Peter Salus of MatrixNetSystems, a network-optimization company based in Austin, TX.⁶

By the mid-2000s, concerns over the physical security of electronic voting systems had risen to public awareness. For example:

A cart of Diebold electronic voting machines was delivered today to the common room of this Berkeley, CA, boarding house, which will be a polling place on Tuesday's primary election. The machines are on a cart which is wrapped in plastic wrap (the same as the stuff we use in the kitchen). A few cable locks (bicycle locks, it seems) provide the appearance of physical security, but they aren't threaded through each machine. Moreover, someone fiddling with the cable locks, I am told, announced after less than a minute of fiddling that he had found the three-digit combination to be the same small integer repeated three times.⁷

2.3.2 1970-1972: Albert the Saboteur. One of the most instructive early cases of computer sabotage occurred at the National Farmers Union Service Corporation of Denver, where a Burroughs B3500 computer suffered 56 disk head crashes in the two years from 1970 to 1972. Downtime was as long as 24 hours per crash, with an average of 8 hours per incident. Burroughs experts were flown in from all over the United States at one time or another, and concluded that the crashes must be due to power fluctuations.

By the time all the equipment had been repaired and new wiring, motor generators, circuit breakers, and power-line monitors had been installed in the computer room, total expenditures for hardware and construction were over \$500,000 (in 1970 dollars). Total expenses related to down time and lost business opportunities because of delays in providing management with timely information are not included in this figure. In any case, after all this expense, the crashes continued sporadically as before.

By this time, the experts were beginning to wonder about their analysis. For one thing, all the crashes had occurred at night. Could it be sabotage? Surely not! Old Albert, the night-shift operator, had been so helpful over all these years; he had unfailingly called in the crashes at once, gone out for coffee and donuts for the repair crews, and been meticulous in noting the exact times and conditions of each crash. However, all the crashes had in fact occurred on his shift.

Management installed a closed-circuit television (CCTV) camera in the computer room—without informing Albert. For some days, nothing happened. Then one night another crash occurred. On the CCTV monitor, security guards saw good ol' Albert open up a disk cabinet and poke his car key into the read/write head solenoid, shorting it out and causing the 57th head crash.

The next morning, management confronted Albert with the film of his actions and asked for an explanation. Albert broke down in mingled shame and relief. He confessed to an overpowering urge to shut the computer down. Psychological investigation determined that Albert, who had been allowed to work night shifts for years without a change, had simply become lonely. He arrived just as everyone else was leaving; he left as everyone else was arriving. Hours and days would go by without the slightest human

IMPERSONATION 2 · 5

interaction. He never took courses, never participated in committees, never felt involved with others in his company. When the first head crashes occurred—spontaneously—he had been surprised and excited by the arrival of the repair crew. He had felt useful, bustling about, telling them what had happened. When the crashes had become less frequent, he had involuntarily, and almost unconsciously, re-created the friendly atmosphere of a crisis team. He had destroyed disk drives because he needed company.⁸

2.4 IMPERSONATION. Using the insignia and specialized language of officials as part of social engineering has a long history in crime; a dramatization of these techniques is in the popular movie *Catch Me If You Can*⁹ about Frank William Abagnale Jr., the teenage scammer and counterfeiter who pretended to be a pilot, a doctor, and a prosecutor before eventually becoming a major contributor to the U.S. Government's anticounterfeiting efforts and then founding a major security firm.¹⁰

Several criminals involved in computer-mediated or computer-oriented crime became notorious for using impersonation.

2.4.1 1970: Jerry Neal Schneider. A notorious computer-related crime started in 1970, when teenager Jerry Neal Schneider used Dumpster[®] diving to retrieve printouts from the Pacific Telephone and Telegraph (PT&T) company in Los Angeles. After years of collection, he had enough knowledge of procedures that he was able to impersonate company personnel on the phone. He collected yet more detailed information on procedures. Posing as a freelance magazine writer, he even got a tour of the computerized warehouse and information about ordering procedures. In June 1971, he ordered \$30,000 of equipment to be sent to a normal PT&T dropoff point—and promptly stole it and sold it. He eventually had a 6,000-square-foot warehouse and 10 employees. He stole over \$1 million of equipment—and sold some of it back to PT&T. He was finally denounced by one of his own disgruntled employees and became a computer security consultant after his prison term.¹¹

2.4.2 1980–2003: Kevin Mitnick. Born in 1963, Kevin Mitnick became involved in crime early, using a special punch for bus transfers to get free rides anywhere in the San Fernando Valley in California by the time he was a young teenager. His own autobiographical comments show him to have been involved in phone phreaking, malicious pranks, and breaking into computers at the Digital Equipment Corporation (DEC) using social engineering.¹²

In 1981, he and his friend Lewis De Payne used social engineering to gain unauthorized access to an operations center for Pacific Bell; “the juvenile court ordered a diagnostic psychological study of Mitnick and sentenced him to a year’s probation.”¹³ In 1987, he was arrested for breaking into the computers of the Santa Cruz Operation, makers of SCO UNIX, and sentenced to three years’ probation.

In the summer of 1988, Mitnick and his accomplice and friend Lenny DiCicco cracked the University of Southern California computers again and misappropriated hundreds of Mb of disk space (a lot at the time) to store VAX VMS source files stolen from Digital Equipment Corporation (DEC). Mitnick was arrested by the Federal Bureau of Investigation (FBI) for having stolen the VAX VMS source code. During his trial, he was described as suffering from an impulse-control disorder. In July 1989, he was sentenced to a year in jail and six months’ rehabilitation. He later tried to become a private investigator and security specialist. He was generally treated with hostility by the established information security community.

2 · 6 HISTORY OF COMPUTER CRIME

In November 1992, Mitnick went underground again when the FBI got a warrant for his arrest on charges of stealing computer time from a phone company. He was located two years later when he made the mistake of leaving insulting messages on the computer and voicemail systems of a physicist and Internet security expert Tsutomu Shimomura. Shimomura was so irritated that he helped law enforcement authorities track the fugitive to North Carolina, where Mitnick was arrested in February 1995 and imprisoned pending trial.

Mitnick was convicted in federal court for the Central District of California on August 9, 1999, and sentenced to 46 months imprisonment for “four counts of wire fraud, two counts of computer fraud, and one count of illegally intercepting a wire communication.”¹⁴ Mitnick was previously sentenced by Judge Pfaelzer to an additional 22 months in prison, this for possessing cloned cellular phones when he was arrested in North Carolina in 1995, and for violating terms of his supervised release imposed after being convicted of an unrelated computer fraud in 1989. He admitted to violating the terms of supervised release by hacking into PacBell voicemail and other systems, and to associating with known computer hackers, in this case codefendant Louis De Payne. Following his release from prison in September 2000, Mitnick was to be on three years’ parole, during which his access to computers was restricted¹⁵ and his profits from writing or speaking about his criminal career were to be turned over to reimburse his victims.

Mitnick earned a living on the talk circuit and eventually founded his own security consulting firm. In the years since his release from prison, he has collaborated in writing several books on social engineering.¹⁶

Perhaps his most significant position in the history of computer crime is that he became an icon in the criminal underground. “FREE KEVIN” was a popular component of Web vandalism for many years, and Eric Corley, the longtime editor of the criminal-hacking publication *2600: The Hacker Quarterly*, even made a movie, *Freedom Downtime*, about what the criminal underground describes as the grossly unfair treatment of Mitnick by the federal government and the news media.¹⁷

2.4.3 Credit Card Fraud. Credit at local businesses dates back into the undocumented past.¹⁸ In the United States, credit cards appeared in the mid-1920s when gasoline companies began issuing cards that were recognized at stations across the country.¹⁹ In 1950, Frank X. McNamara started the Diners Club, the first credit card company serving multiple types of businesses; the company began the practice of charging a percentage fee for each transaction and also charged its clients a membership fee.²⁰ The VISA card evolved from the 1951 BankAmericard from the Bank of America, and a consortium of California banks established MasterCard shortly thereafter. American Express started its card program in 1958.

Card use rose and, unsurprisingly, credit card fraud was rampant. Mail theft also became widespread as unscrupulous individuals discovered that envelopes containing credit cards were just like envelopes full of cash. And there was little to stop card companies from sending out cards that customers had never asked for, were not expecting, and could not have known had been stolen, until the issuing company began demanding payment for the charges that had been run up. These crimes and other problems stemming from the relentless card-pushing by banks led directly to the passage of the Fair Credit Billing Act of 1974²¹ as well as many other laws²² designed to protect the consumer.²³

IMPERSONATION 2 · 7

By the mid-1990s, credit card fraud was a rapidly growing problem for consumers and for law enforcement. A 1997 FBI report stated:

Around the world, bank card fraud losses to Visa and Master-Card alone have increased from \$110 million in 1980 to an estimated \$1.63 billion in 1995 The United States has suffered the bulk of these losses—approximately \$875 million for 1995 alone. This is not surprising because 71 percent of all worldwide revolving credit cards in circulation were issued in this country Law enforcement authorities continually confront new and complex schemes involving credit card frauds committed against financial institutions and bank card companies. Perpetrators run the gamut from individuals with easy access to credit card information—such as credit agency officials, airline baggage handlers, and mail carriers, both public and private, to organized groups, usually from similar ethnic backgrounds, involved in large-scale card theft, manipulation, and counterfeiting activities. Although current bank card fraud operations are numerous and varied, several schemes account for the majority of the industry's losses by taking advantage of dated technology, customer negligence, and laws peculiar to the industry.²⁴

2.4.4 Identity Theft Rises. By the late 1990s and in the decade following the year 2000, credit card fraud was subsumed into the broader category of *identity theft*. Instead of limiting their depredations to running up bills on stolen or forged credit card accounts, thieves, often in organized rings, created entire bogus parallel identities, initiating unpaid bank loans, buying cars with other people's credit, and wreaking havoc with innocent victims' credit ratings, financial situations, and even their daily life. Victims of extreme cases lost their ability to obtain mortgages, buy new homes, and accept new jobs. Worse, the burden of proof of innocence fell on the victims, in a bitter reversal of the assumption of innocence underlying British common law and its offshoot in the commonwealth and the United States.

In August 2008, the U.S. Department of Justice announced²⁵ the single largest and most complex case of identity theft ever charged in this country. It involved eleven people from five different countries, including two from the United States and two from the People's Republic of China, who had stolen more than 40,000,000 credit card records from a major U.S. retailer. They drove by, or loitered at, buildings in which wireless networks were housed, and installed sniffers that recorded passwords, card numbers, and account data. Unless adequate preventative measures are installed quickly, more such horrendous events will be sure to occur. For more on wireless network security, see Chapter 33 in this *Handbook*.

The 2011 report from the U.S. Bureau of Justice Statistics, *Identity Theft Reported by Households, 2005–2010* provides additional details. The abstract includes these highlights:

- In 2010, 7.0 percent of households in the United States, or about 8.6 million households, had at least one member age 12 or older who experienced one or more types of identity theft victimization.
- Among households in which at least one member experienced one or more types of identity theft, 64.1 percent experienced the misuse or attempted misuse of an existing credit card account in 2010.
- From 2005 to 2010, the percentage of all households with one or more type of identity theft that suffered no direct financial loss increased from 18.5 percent to 23.7 percent.²⁶

2 · 8 HISTORY OF COMPUTER CRIME

A report from Javelin Strategy & Research covering identity fraud in 2012 included the following observation in the overview: “Identity fraud incidence increased in 2012 for the second consecutive year, affecting 5.26% of U.S. adults. This increase was driven by dramatic jumps in the two most severe fraud types, new account fraud (NAF) and account takeover fraud (ATF).”²⁷

2.5 PHONE PHREAKING. Even in the earliest days of telephony, teenage boys played with the new technology to cause havoc. In the late 1870s, the new AT&T system in America had to stop using teenagers as switchboard operators:

The boys were openly rude to customers. They talked back to subscribers, saucing off, uttering facetious remarks, and generally giving lip. The rascals took Saint Patrick’s Day off without permission. And worst of all they played clever tricks with the switchboard plugs: disconnecting calls, crossing lines so that customers found themselves talking to strangers, and so forth.

This combination of power, technical mastery, and effective anonymity seemed to act like catnip on teenage boys.²⁸

2.5.1 2600 Hz. In the late 1950s, AT&T began switching its telephone networks to direct-dial long distance, using specific frequency tones to communicate among its switches. Around 1957, a blind seven-year-old child named Josef Engressia with perfect pitch and an emotional fixation on telephones learned to whistle the 2600-Hz pitch that interrupted long-distance telephone calls and allowed him to place a free long-distance call to anywhere in the world.²⁹ This emotionally disturbed person eventually renamed himself “Joybubbles” and is often described as the founder of phone phreaking—the manipulation of the phone system for unauthorized access to services.

John Draper was in the U.S. Air Force in 1964 when he began helping his colleagues place free phone calls. At the suggestion of Joybubbles, he used the whistles in Cap’n Crunch cereal boxes to generate the 2600-Hz tone and then, calling himself Captain Crunch, went on to create electronic tone synthesizers called *blue boxes*.³⁰ In the 1970s, Apple founders Steve Wozniak and Steve Jobs built blue boxes and, using the devices, perpetrated such pranks as calling the Vatican while pretending to be Henry Kissinger.³¹

A significant contributor to the growth of phreaking in the 1970s was the publication in 1971 of an article about phreaking in *Esquire Magazine*, which attracted the attention of many young technophiles.³²

2.5.2 1982–1991: Kevin Poulsen. As the phone system shifted to greater reliance on computers, the border between phreaking and hacking began to blur. One of the important names from the 1980s period of fascination with everything phone-related was Kevin Poulsen.

Kevin Poulsen’s autobiographical sketch is shown next.

Kevin Poulsen first gained notoriety in 1982, when the Los Angeles County District Attorney’s Office raided him for gaining unauthorized access to a dozen computers on the ARPANET, the forerunner of the modern Internet. Seventeen years old at the time, he was not charged, and went on to work as a programmer and computer security supervisor for SRI International in Menlo Park, California, then as a network administrator at Sun Microsystems.

In 1987, Pacific Bell security agents discovered that Poulsen and his friends had been penetrating telephone company computers and buildings. After learning that Poulsen had also worked for a defense contractor where he’d held a SECRET level security clearance, the FBI began building an espionage case against the hacker.

DATA DIDDLING 2 · 9

Confronted with the prospect of being held without bail, Poulsen became a fugitive. While on the run, he obtained information on the FBI's electronic surveillance methods, and supported himself by hacking into Pacific Bell computers to cheat at radio-station phone-in contests, winning a vacation to Hawaii and a Porsche 944-S2 Cabriolet in the process.

After surviving two appearances on NBC's *Unsolved Mysteries*, Poulsen was finally captured on April 10th, 1991, in a Van Nuys grocery store, by a Pacific Bell security agent acting on an informant's tip. On December 4th, 1992, Poulsen became the first hacker to be indicted under U.S. espionage laws when the Justice Department charged him with stealing classified information. (18 U.S.C. 793).

Poulsen was held without bail while he vigorously fought the espionage charge. The charge was dismissed on March 18th, 1996.

Poulsen served five years, two months, on a 71-month sentence for the crimes he committed as a fugitive, and the phone hacking that began his case. He was freed June 4th, 1996, and began a three-year period of supervised release, barred from owning a computer for the first year, and banned from the Internet for the next year and a half.

Since his release, Poulsen has appeared on MSNBC, and on ABC's *Nightline*, and he was the subject of Jon Littman's flawed book, "The Watchman—the Twisted Life and Crimes of Serial Hacker Kevin Poulsen." His case has earned mention in several computer security and infowar tracts—most of which still report that he broke into military computers and stole classified documents.³³

After his release from prison, Kevin Poulsen turned to journalism. He became an editor for *SecurityFocus* and then was hired as a senior editor at *Wired News*. He is a serious investigative reporter (e.g., he broke the story of sexual predators in MySpace)³⁴ and a frequent contributor to the "Threat Level" blog.³⁵

2.6 DATA DIDDLING. One of the most common forms of computer crime since the start of electronic data processing is *data diddling*—illegal or unauthorized data alteration. These changes can occur before and during data input, or before output. Data-diddling cases have included bank records, payrolls, inventory data, credit records, school transcripts, telephone switch configurations, and virtually all other applications of data processing.

2.6.1 Equity Funding Fraud (1964–1973). One of the classic early data-diddling frauds was the Equity Funding case, which began with computer problems at the Equity Funding Corporation of America, a publicly traded and highly successful firm with a bright idea. The idea was that investors would buy insurance policies from the company and also invest in mutual funds at the same time, with profits to be redistributed to clients and to stockholders. Through the late 1960s, Equity's shares rose dizzyingly in price, and there were news magazine stories about this wunderkind of the Los Angeles business community.

The computer problems occurred just before the close of the financial year in 1964. An annual report was about to be printed, yet the final figures simply could not be extracted from the mainframe. In despair, the head of data processing told the president the bad news; the report would have to be delayed. Nonsense, said the president expansively (in the movie, anyway); simply make up the bottom line to show about \$10 million in profits and calculate the other figures so it would come out that way. With trepidation, the DP chief obliged. He seemed to rationalize it with the thought that it was just a temporary expedient, and could be put to rights later in the real financial books.

2 · 10 HISTORY OF COMPUTER CRIME

The expected profit did not materialize, and some months later, it occurred to the executives at Equity that they could keep the stock price high by manufacturing false insurance policies that would make the company look good to investors. They therefore began inserting false information about nonexistent policyholders into the computerized records used to calculate the financial health of Equity.

In time, Equity's corporate staff got even greedier. Not content with jacking up the price of their stock, they decided to sell the policies to other insurance companies via the redistribution system known as reinsurance. Reinsurance companies pay money for policies they buy and spread the risk by selling parts of the liability to other insurance companies. At the end of the first year, the issuing insurance companies have to pay the reinsurers part of the premiums paid in by the policyholders. So in the first year, selling imaginary policies to the reinsurers brought in large amounts of real cash. However, when the premiums came due, the Equity crew "killed" imaginary policyholders with heart attacks, car accidents, and, in one memorable case, cancer of the uterus—in a male imaginary policyholder.

By late 1972, the head of DP calculated that by the end of the decade, at this rate, Equity Funding would have insured the entire population of the world. Its assets would surpass the gross national product of the planet. The president merely insisted that this showed how well the company was doing.

The scheme fell apart when an angry operator who had to work overtime told the authorities about shenanigans at Equity. Rumors spread throughout Wall Street and the insurance industry. Within days, the Securities and Exchange Commission had informed the California Insurance Department that they had received information about the ultimate form of data diddling: Tapes were being erased. The officers of the company were arrested, tried, and condemned to prison terms.³⁶

2.6.2 1994: Vladimir Levin and the Citibank Heist. In February 1998, Vladimir Levin was sentenced to three years in prison by a court in New York City. Levin masterminded a major conspiracy in 1994 in which the gang illegally transferred \$12 million in assets from Citibank to a number of international bank accounts. The crime was spotted after the first \$400,000 was stolen in July 1994, and Citibank cooperated with the FBI and Interpol to track down the criminals. Levin was ordered to pay back \$240,000, the amount he actually managed to withdraw before he was arrested.³⁷ The incident led to Citibank's hiring of Stephen R. Katz as the banking industry's first chief information security officer (CISO).

2.7 SALAMI FRAUD. In the salami technique, criminals steal money or resources a bit at a time. Two different etymologies are circulating about the origins of this term. One school of security specialists claim that it refers to slicing the data thin—like a salami. Others argue that it means building up a significant object or amount from tiny scraps—like a salami.

There were documented cases of salami frauds in the 1970s and 1980s, but one of the more striking incidents came to light in January 1993, when four executives of a Value Rent-a-Car franchise in Florida were charged with defrauding at least 47,000 customers using a salami technique. The federal grand jury in Fort Lauderdale claimed that the defendants modified a computer billing program to add five extra gallons to the actual gas tank capacity of their vehicles. From 1988 through 1991, every customer who returned a car without topping it off ended up paying inflated rates for an inflated

LOGIC BOMBS 2 · 11

total of gasoline. The thefts ranged from \$2 to \$15 per customer—rather thick slices of salami but nonetheless difficult for the victims to detect.

Unfortunately, salami attacks are *designed* to be difficult to detect. The only hope is that random audits, especially of financial data, will pick up a pattern of discrepancies and lead to discovery. As any accountant will warn, even a tiny error must be tracked down, since it may indicate a much larger problem. For example, Cliff Stoll's famous adventures tracking down spies in the Internet began with an unexplained \$0.75 discrepancy between two different resource accounting systems on UNIX computers at the Keck Observatory of the Lawrence Berkeley Laboratories. Stoll's determination to understand how the problem could have occurred revealed an unknown user; investigation led to the discovery that resource-accounting records were being modified to remove evidence of system use. The rest of the story is told in Clifford Stoll's book *The Cuckoo's Egg*.

2.8 LOGIC BOMBS. A logic bomb is a program that has deliberately been written or modified to produce results when certain conditions are met that are unexpected and unauthorized by legitimate users or owners of the software. Logic bombs may be within standalone programs, or they may be part of worms (programs that hide their existence and spread copies of themselves within a computer systems and through networks) or viruses (programs or code segments which hide within other programs and spread copies of themselves).

Time bombs are a subclass of logic bombs that “explode” at a certain time.

According to a National Security Council employee, the United States Government authorized insertion of a time bomb in software to control the Trans-Siberian natural gas pipeline that they knew would be stolen from U.S. sources by the Soviet government. “The result was the most monumental non-nuclear explosion and fire ever seen from space,” said Thomas C. Reed.³⁸

The infamous Jerusalem virus (also known as the Friday the 13th virus) of 1988 was a time bomb. It duplicated itself every Friday and on the thirteenth of the month, causing system slowdown; on every Friday the 13th after May 13, 1988, it also corrupted all available disks on the infected systems.

Other examples of notorious time bombs include:

- A common PC virus from the 1980s, *Cascade*, made all the characters fall to the last row of the display during the last three months of every year.
- The Michelangelo virus of 1992 was designed to damage hard disk directories on the sixth of March every year.
- In 1992, computer programmer Michael Lauffenburger was fined \$5,000 for leaving a logic bomb at General Dynamics. His intention was to return after his program had erased critical data and be paid to fix the problem.³⁹

The most famous time bomb of recent years was the Y2K (year 2000) problem. In brief, old programs used two-digit year codes that were based on the assumption that they applied to the twentieth century. As the twenty-first century approached, analysts warned of catastrophic consequences if the programs were not corrected to use four-digit years or otherwise adapt to the change of century.⁴⁰ In the event, the corrective measures worked and there were no disasters. Later analysis showed a positive correlation between investments in Y2K remediation and later profitability.⁴¹

2 · 12 HISTORY OF COMPUTER CRIME

2.9 EXTORTION. Computer data can be held for ransom. For example, according to Whiteside, in 1971, two reels of magnetic tape belonging to a branch of the Bank of America were stolen at Los Angeles International Airport. The thieves demanded money for their return. The owners ignored the threat of destruction because they had adequate backup copies.

Other early cases of extortion involving computers:

- In 1973, a West German computer operator stole 22 tapes and received \$200,000 for their return. The victim did not have adequate backups.
- In 1977, a programmer in the Rotterdam offices of Imperial Chemical Industries, Ltd. (ICI) stole all his employer's tapes, including backups. Luckily, ICI informed Interpol of the extortion attempt. As a result of the company's forthrightness, the thief and an accomplice were arrested in London by officers from Scotland Yard.

In the 1990s, one of the most notorious cases of extortion was the 1999 theft of 300,000 records of customer credit cards from the CD Universe Web site by "Maxus," a 19-year-old Russian. He sent an extortion note that read: "Pay me \$100,000 and I'll fix your bugs and forget about your shop forever ... or I'll sell your cards [customer credit data] and tell about this incident in news." Refused by CD Universe owners, he promptly released 25,000 credit card numbers via a Web site that became so popular with criminals that Maxus had to limit access to one stolen number per visit.

2.10 TROJAN HORSES. Trojans are programs that pretend to be useful but that also contain harmful code or are just plain harmful.

2.10.1 1988 Flu-Shot Hoax. One of the nastiest tricks played on the shell-shocked world of early microcomputer users was the FLU-SHOT-4 incident of March 1988. With the publicity given to damage caused by destructive, self-replicating virus programs distributed through electronic bulletin board systems (BBSs), it seemed natural that public-spirited programmers would rise to the challenge and provide protective screening.

Flu-Shot-3 was a useful program for detecting viruses. Flu-Shot-4 appeared on BBSs and looked just like version 3; however, it actually destroyed critical areas of hard disks and any floppies present when the program was run. The instructions that caused the damage were not present in the program file until it was running; this self-modifying code technique makes it especially difficult to identify Trojans by simple inspection of the assembler-level code.

2.10.2 Scrambler, 12-Tricks, and PC Cyborg. Other early and notorious PC Trojans from the late 1980s that are still remembered in the industry included:

- The Scrambler (also known as the KEYBGR Trojan), which pretended to be a keyboard driver (KEYBGR.COM), but actually made a smiley face move randomly around the screen

TROJAN HORSES 2 · 13

- The 12-Tricks Trojan, which masqueraded as CORETEST.COM, a program for testing the speed of a hard disk, but actually caused 12 different kinds of damage (e.g., garbling printer output, slowing screen displays, and formatting the hard disk).
- The PC Cyborg Trojan (or “AIDS Trojan”), which claimed to be an AIDS information program but actually encrypted all directory entries, filled up the entire C disk, and simulated COMMAND.COM, but produced an error message in response to nearly all commands.

2.10.3 1994: Datacomp Hardware Trojan. On November 8, 1994, a correspondent reported to the *RISKS Forum Digest* that he had been victimized by a curious kind of Trojan:

I recently purchased an Apple Macintosh computer at a “computer superstore,” as separate components—the Apple CPU, and Apple monitor, and a third-party keyboard billed as coming from a company called Sicon.

This past weekend, while trying to get some text-editing work done, I had to leave the computer alone for a while. Upon returning, I found to my horror that the text “welcome datacomp” had been inserted into the text I was editing. I was certain that I hadn’t typed it, and my wife verified that she hadn’t, either. A quick survey showed that the “clipboard” (the repository for information being manipulated via cut/paste operations) wasn’t the source of the offending text.

As usual, the initial reaction was to suspect a virus. Disinfectant, a leading anti-viral application for Macintoshes, gave the system a clean bill of health; furthermore, its descriptions of the known viruses (as of Disinfectant version 3.5, the latest release) did not mention any symptoms similar to my experiences.

I restarted the system in a fully minimal configuration, launched an editor, and waited. Sure enough, after a (rather long) wait, the text “welcome datacomp” once again appeared, all at once, on its own.

Further investigation revealed that someone had put unauthorized code in the ROM chip used in several brands of keyboard. The only solution was to replace the keyboard. Readers will understand the possible consequences of a keyboard that inserts unauthorized text into, say, source code. Winn Schwartau, the renowned computer security expert, has coined the word “chipping” to refer to such unauthorized modification of firmware.

2.10.4 Keylogger Trojans. By the mid-2000s, software and hardware Trojans designed to capture logs of keystrokes and sometimes to transmit those logs via covert Internet connections had become a well-known tool of industrial espionage. The United States Department of Homeland Security issued a warning in December 2005 that included this overview:

According to industry security experts, the biggest security vulnerability facing computer users and networks is email with concealed Trojan Horse software—destructive programs that masquerade as benign applications and embedded links to ostensibly innocent websites that download malicious code. While firewall architecture blocks direct attacks, email provides a vulnerable route into an organization’s internal network through which attackers can destroy or steal information.

Attackers try to circumvent technical blocks to the installation of malicious code by using social engineering—getting computer users to unwittingly take actions that allow the code to be installed and organization data to be compromised.

2 · 14 HISTORY OF COMPUTER CRIME

The techniques attackers use to install Trojan Horse programs through email are widely available, and include forging sender identification, using deceptive subject lines, and embedding malicious code in email attachments.

Developments in thumb-sized portable storage devices and the emergence of sophisticated keystroke logging software and devices make it easy for attackers to discover and steal massive amounts of information surreptitiously.⁴²

2.10.5 Haephrati Trojan. A case that made the news in the mid-2000s began when Israeli author Amon Jackont was upset to find parts of the manuscript on which he was working posted on the Internet. Then someone tried to steal money from his bank account. Suspicion fell on his stepdaughter's ex-husband, Michael Haephrati. Police discovered a keystroke logger on Jackont's computer. It turned out that Haephrati had also sold spy software to clients; the Trojan was concealed in what appeared to be confidential email. Once installed on the victims' computers, the software sent surveillance data to a server in London, England.

Haephrati was detained by U.K. police and investigations began in Germany and Israel. Twelve people were detained in Israel; eight others were under house arrest. Suspects included private investigators and top executives from industrial firms. Victims included Hewlett-Packard, Ace hardware stores, and a cable-communications company.

Michael and Ruth Haephrati were extradited from Britain for trial in Israel on January 31, 2006. They were accused of installing the Trojan horse program that activated a key logger with remote-reporting capabilities.⁴³

In March 2006, the couple were indicted in Tel Aviv for corporate espionage.⁴⁴ They pleaded guilty to the charges⁴⁵ and were sentenced to four and two years of jail, respectively, as well as punished with fines.⁴⁶

The story did not end there, however. Two years later, “Four members of the Israeli Modi'in Ezrahi private investigation firm were sentenced on Monday after they were found guilty of using Trojan malware to steal commercially sensitive information from their clients’ competitors.”⁴⁷ The report continues:

Asaf Zlotovsky, a manager at the Modi'in Ezrahi detective firm, was jailed for 19 months. Two other employees, Haim Zissman and Ron Barhoum, were sent to prison for 18 and nine months respectively. The firm’s former chief exec, Yitzhak Rett, the victim of an apparent accident when he fell down a stairwell during a break in police questioning back in 2005, escaped a jail sentence under a plea bargaining agreement. Rett was fined 250,000 Israeli Shekels (£36,500) and ordered to serve ten months’ probation over his involvement in the scam.

However, an article in April 2008 reported that Michael Haephrati “claimed that there was no jail time, and that he was completely free. As a matter of fact he was going to continue to offer his Trojan Horse service but this time he would only work with ‘law enforcement agencies.’”⁴⁸

2.10.6 Hardware Trojans and Information Warfare. In the late 2000s, a flurry of news stories discussed the dangers of growing reliance on Chinese-manufactured computing components.

U.S. Defense Department sources say privately that the level of Chinese cyberattacks obliges them to avoid Chinese-origin hardware and software in all classified systems and as many unclassified systems as fiscally possible. The high threat of Chinese cyberpenetrations into U.S. defense networks will be magnified as the Pentagon increasingly loses domestic sources of “trusted and classified” microchips.⁴⁹

NOTORIOUS WORMS AND VIRUSES 2 · 15

The discovery of counterfeit Cisco routers worsened concerns about the reliability of Chinese-manufactured network equipment.⁵⁰ The FBI, Immigration and Customs Enforcement (ICE), Customs and Border Protection (CBP), and the Royal Canadian Mounted Police (RCMP) worked together to track a massive pattern of counterfeit network hardware including Cisco routers; these investigations and seizures raised questions about the reliability and trustworthiness of such equipment, much of which was manufactured in the People's Republic of China. Although Cisco scientists examined some of the counterfeit equipment and found no back doors, concern was serious enough that government agencies created test chips to challenge quality assurance processes at military contractors:

In April [2008], the Defense Advanced Research Projects Agency, part of the Defense Department, began distributing chips with hidden Trojan horse circuitry to military contractors participating in an agency program, Trusted Integrated Circuits. The goal is to test forensic techniques for finding hidden electronic trap doors, which can be maddeningly elusive. The agency is not yet ready to announce the results of the test, said Jan Walker, a spokeswoman for the agency.⁵¹

A 2011 report on hearings before the U.S. House Oversight and Government Reform Committee about the issue of Trojan backdoors in imported software and hardware included this assertion with references for each topic:

... [E]mbedded malware lurking in consumer tech is not a new development. Since it's been happening for years and is hardly a national security secret, it's unclear why Schaffer hesitated so long before answering. There have been many incidents of malware-infected products being shipped to consumers, from hardware, to software, and even tainted peripheral devices. Malware has been sent pre-loaded in products like USBs, microchips, cameras, battery chargers, digital photo frames, webcams, printers, cell phones, motherboards or system boards, and hard drives.⁵²

2.11 NOTORIOUS WORMS AND VIRUSES. The next sections briefly describe some of the outstanding incidents that are often mentioned in discussions of the history of malware.⁵³

2.11.1 1970-1990: Early Malware Outbreaks. The ARPANET was the precursor of the Internet.⁵⁴ According to several reports:

Sometime in the early 1970s, the Creeper virus was detected on ARPANET, a US military computer network which was the forerunner of the modern Internet. Written for the then-popular Tenex operating system, this program was able to gain access independently through a modem and copy itself to the remote system. Infected systems displayed the message, "I'M THE CREEPER: CATCH ME IF YOU CAN."

Shortly thereafter, the Reaper program was anonymously created to delete Creeper. Reaper was a virus: it spread to networked machines and if it located a Creeper virus, Reaper would delete it. Even the participants are unable to say whether Reaper was a response to Creeper, or if it was created by the same person or persons who created Creeper in order to correct their mistake.⁵⁵

By 1981, the Apple II computer was a popular system among hobbyists; the Elk Cloner virus spread via infected floppy disks and is regarded as "the first large-scale computer virus outbreak in history."⁵⁶

2 · 16 HISTORY OF COMPUTER CRIME

In 1986, the Brain boot-sector virus was the first IBM-PCs malware to spread around the world. It was created by two brothers from Lahore, Pakistan, and included this text:

Welcome to the Dungeon (c) 1986 Brain & Amjads (pvt) Ltd VIRUS_SHOE RECORD V9.0
Dedicated to the dynamic memories of millions of viruses who are no longer with us today -
Thanks GOODNESS!! BEWARE OF THE er...VIRUS: this program is catching program
follows after these messages....\$#@%@!!

The Lehigh Virus appeared at Lehigh University in Pennsylvania in 1987 and damaged the files of several professors and students. This early program-infector targeted only *command.com* and was therefore extremely limited in its spread.

In 1988, the Jerusalem virus, a file infector that reproduced by inserting its code into EXE and COM files, caused a global PC epidemic.

Another noteworthy infection of 1998 came from the self-encrypting Cascade virus of 1988, which confused many naive users who interpreted the falling symbols on their screen as part of an unexpected screen saver. This virus was one of the earliest examples of the attempts to counter signature-based antivirus products.

2.11.2 December 1987: Christmas Tree Worm. In December 1987, users of IBM mainframe computers connected to the European Academic Research Network (EARN), BITNET, and the IBM company VNET were flooded with email bearing a character-based representation of a Christmas tree. A student at Technische Universität Clausthal⁵⁷ in Germany launched “a worm, written in an IBM-specific language called REXX.”⁵⁸ The worm used the victim’s list of correspondents to send copies of itself to everyone on the list.⁵⁹

2.11.3 November 2, 1988: Morris Worm. On November 2, 1988, the Internet was rocked by the explosive appearance of unauthorized code on systems all over the world. At 17:00 EST on November 2, 1988, Robert T. Morris, a student at Cornell University in Ithaca, New York, released a worm into the Internet. By midnight, it had attacked VAX computers running 4 BSD UNIX and SUN Microsystems Sun 3 computers throughout the United States. One of the most interesting aspects of the worm’s progress through the Internet was the almost complete independence of its path from normal geographical constraints. It sometimes leaped from coast to coast faster than it reached physically neighboring computer systems. The worm graphically demonstrated that cyberspace has its own geography.

The worm often superinfected its hosts, leading to slowdowns in overall processing speed. The first Internet warning (“We are under attack”) was posted at 02:38 on November 3 to the TCP-IP list by a scientist at University of California at Berkeley. At 03:34, Andy Sudduth, a friend of Morris’s at Harvard, posted a warning message (“There may be a virus loose on the internet”) anonymously and included a few comments on how to stop the worm. Unfortunately, Spafford writes, the Internet was so severely impeded by the worm that this message was not widely distributed for over 24 hours.

By 6:00 on the morning of November 3, messages were creeping through the Internet with details of how the worm worked. The news spread via news groups such as the TCP-IP list, Usenix 4bsd-ucb-fixes, and the Usenet news.announce.important group. Spafford and his friends and colleagues on the Internet collaborated feverishly on providing patches against the worm.

NOTORIOUS WORMS AND VIRUSES 2 · 17

Meanwhile, as word spread of the attack, some systems administrators began cutting their networks out of the Internet. The Defense Communications Agency isolated its Milnet and Arpanet networks from each other around 11:30 on November 3. At noon, machines in the science and technology center at the Stanford Research Institute were shut down.

By late on November 4, a comprehensive set of patches was posted on the Internet to defend systems against the worm. That evening, a *New York Times* reporter told Spafford that the author of the worm had been found.

By November 8, the Internet seemed to be back to normal. A group of concerned computer scientists met at the National Computer Security Center to study the incident and think about preventing recurrences of such attacks. Spafford put the incident into perspective with the comment that the affected systems were no more than 5 percent of the hosts on the Internet. It would be foolish to dismiss Morris's electronic vandalism as a prank or to claim that the worm alerted managers to weak security on their systems. Nonetheless, it is true that the incident contributed to the establishment of the Computer Emergency Response Team at the Software Engineering Institute of Carnegie-Mellon University. For these blessings, however, we owe no gratitude to Robert T. Morris.

In 1990, Morris was found guilty under the Computer Fraud and Abuse Act of 1986. The maximum penalties included five years in prison, a \$250,000 fine, and restitution costs. Morris was ordered to perform 400 hours of community service, sentenced to three years probation, and required to pay \$10,000 in fines. He was expelled from Cornell University.

His lawyers appealed the conviction to the Supreme Court of the United States. Their arguments included lack of evil intent (he did not mean to cause harm, honest—even though his worm took extraordinary precautions to conceal itself) and they deplored the scandalous behavior of Cornell University authorities, who had the temerity to search their own electronic mail message system to locate evidence that incriminated Morris. The lawyers also argued that sending a mail message might become a crime if Morris's conviction were upheld.

The Supreme Court upheld the decision by declining to hear the appeal.⁶⁰

Robert T. Morris eventually became an associate professor in the Electrical Engineering and Computer Science Department of the Massachusetts Institute of Technology and a member of the Computer Science and Artificial Intelligence Laboratory.⁶¹

2.11.4 Malware in the 1990s. The most significant malware development of the 1990s was the release in July 1995 of the world's first widely distributed macro-language virus. The *macro.concept* virus made its appearance in MS-Word for Windows documents. It demonstrated how to use the macro programming language, common to many Microsoft products, to generate self-reproducing macros that spread from document to document. Within a few months, clearly destructive versions of this demonstration virus appeared.

Macro viruses were a dangerous new development. As explained in a recent history of viruses and antivirus:

- Putting self-reproducing code in easily and frequently exchanged files, such as documents, greatly increased the infectiousness of the viruses.
- Virus writers shifted their attention to a much easier programming language than assembly.

2 · 18 HISTORY OF COMPUTER CRIME

- Email exchanges of infected documents were a far more effective mechanism for virus infection than exchanges of infected programs or disks.
- “[M]acro viruses were neither platform-specific, nor OS-specific. They were application-based.”⁶²

In the latter half of the 1990s, macro viruses replaced boot sector viruses and file infector viruses as a major type of malicious self-reproducing malware; during that period, additional types of script-based, network worms also increased.

Exhibit 2.1 shows the rise and fall of prevalence of macro viruses over the decade from discovery to extinction using data from the WildList archives. The WildList shows malware identified on user systems by at least two virus researchers.⁶³

Roger Thompson summarizes the developments in malware in the 1990s in this way:

By around 2000, macro viruses ceased to be a problem because the new version of MS-Office 2000 included features that blocked macro viruses. The next step in the evolution of malware was the mass mailers like the ILOVEYOU worm and then the network worms. These were easy to write and easy to obfuscate by varying the text contents, thus defeating signature scanners. These worms spread very quickly until the release of Windows XP Service Pack 2, which forced the Windows Firewall to be on by default. After that extinction-level event, criminals moved onward to creating mass mailers and bots which could spread malware and spam or cause distributed denial-of-service through communication via the trusted Web sites accessed through browsers that created a tunnel through the firewall.⁶⁴

2.11.5 March 1999: Melissa. On Friday, March 26, 1999, the CERT/CC received initial reports of a fast-spreading new MS-Word macro virus. “Melissa” was written to infect such documents; once loaded, it uses the victim’s MAPI-standard email address book to send copies of itself to the first 50 people on the list. The virus attaches an infected document to an email message with subject line “Subject: Important Message From <name>” where <name> is that of the inadvertent sender. The email message reads: “Here is that document you asked for … don’t show anyone else;-)” and includes a MS-Word file as an infected attachment. The original infected document, “list.doc,” was a compilation of URLs for pornographic Websites. However, as the virus spread, it was capable of sending any other infected document created by the victim.

Because of this high replication rate, the virus spread faster than any previous virus in history. On many corporate systems, the rapid rate of internal replication saturated email servers with outbound automated junk email. Initial estimates were in the range of 100,000 downed systems. Antivirus companies rallied immediately, and updates for all the standard products were available within hours of the first notices from CERT/CC.

The search for the originator of the Melissa email computer virus/worm began immediately after the outbreak. Initial findings traced the virus to Access Orlando, a Florida Internet Service Provider (ISP), whose servers were shut down by order of the FBI for forensic examination; the systems were then confiscated. That occurrence was then traced back to Source of Kaos, a free-speech Website where the virus may have lain dormant for months in a closed but not deleted virus-distributor’s pages. Investigators discovered a serial number in the vector document, written with MS-Word; the undocumented serial number helped law enforcement when investigators circulated it on the Net to help track down the perpetrator.

The next steps turned to the value-added network AOL, where the virus was released to the public. The giant ISP’s information helped to identify a possible suspect and

NOTORIOUS WORMS AND VIRUSES 2 · 19**EXHIBIT 2.1** Rise and Fall in Macro Viruses in the WildList, 1996–2008

Year	Macro Viruses	Total Entries	Percentage Macro Virus
1996 ^a	1	183	0.6%
1997 ^b	27	239	11%
1998 ^c	77	258	30%
1999 ^d	46	129	36%
2000 ^e	108	175	62%
2001 ^f	145	228	64%
2002 ^g	103	198	52%
2003 ^h	68	205	33%
2004 ⁱ	51	261	20%
2005 ^j	22	399	6%
2006 ^k	19	804	2%
2007 ^l	5	797	0.6%
2008 ^m	0	590	0.0%

^aWildList Organization International, "PC Viruses in the Wild—January 10, 1996," www.wildlist.org/WildList199601.htm.

^bWildList Organization International, "PC Viruses in the Wild—February, 1997," www.wildlist.org/WildList199702.htm.

^cWildList Organization International, "PC Viruses in the Wild—January, 1998," www.wildlist.org/WildList199801.htm.

^dWildList Organization International, "PC Viruses in the Wild—January 1999," www.wildlist.org/WildList199901.htm.

^eWildList Organization International, "PC Viruses in the Wild—January, 2000," www.wildlist.org/WildList200001.htm.

^fWildList Organization International, "PC Viruses in the Wild—January, 2001," www.wildlist.org/WildList200101.htm.

^gWildList Organization International, "PC Viruses in the Wild—January, 2002," www.wildlist.org/WildList200201.htm.

^hWildList Organization International, "PC Viruses in the Wild—January, 2003," www.wildlist.org/WildList200301.htm.

ⁱWildList Organization International, "PC Viruses in the Wild—January, 2004," www.wildlist.org/WildList200401.htm.

^jWildList Organization International, "PC Viruses in the Wild—January, 2005," www.wildlist.org/WildList200501.htm.

^kWildList Organization International, "PC Viruses in the Wild—January, 2006," www.wildlist.org/WildList200601.htm.

^lWildList Organization International, "PC Viruses in the Wild—January, 2007," www.wildlist.org/WildList200701.htm.

^mWildList Organization International, "PC Viruses in the Wild—January, 2008," www.wildlist.org/WildList200801.htm.

by April 2, the FBI arrested David L. Smith (age 30) of Aberdeen, New Jersey. Smith apparently panicked when he heard the FBI was on the trail of the Melissa spawner and he threw away his computer—stupidly, into the trash at his own apartment building.

Smith was charged with second-degree offenses of interruption of public communication, conspiracy to commit the offense and attempt to commit the offense, third-degree theft of computer service, and third-degree damage or wrongful access to computer systems. If convicted, Smith faced a maximum penalty of \$480,000 in fines and 40 years in prison. On December 10, 1999, Smith pleaded guilty to all federal charges and agreed to every particular of the indictment, including the estimates by the International Computer Security Association of at least \$80 million of consequential damages due to the Melissa infections.⁶⁵

2 · 20 HISTORY OF COMPUTER CRIME

2.11.6 May 2000: I Love You. Starting around May 4, 2000, email users opened messages from familiar correspondents with the subject line “I love you”; many then opened the attachment, LOVE-LETTER-FOR-YOU.txt.vbs, which infected the user’s email address book and initiated mass mailing of itself to all the contacts. The “Love Bug” was the fastest-spreading worm to that time, infecting computers all over the world, starting in Asia, then Europe.⁶⁶

On May 11, Filipino computer science student Onel de Guzman of AMA Computer College in Manila admitted to authorities that he may “accidentally have launched the destructive Love Bug virus out of youthful exuberance.” He did not admit that he had created the malware himself; however, the name GRAMMERSoft appeared in the computer code of the virus, and that was the name of a computer group to which the 23-year-old de Guzman belonged.⁶⁷

In September 2000, de Guzman participated in a live chat hosted by CNN.com; he vigorously defended virus-writing and blamed the creators of vulnerable systems for releasing poorly designed software. He refused to take responsibility for writing the worm.⁶⁸

Philippine authorities tried to prosecute de Guzman but had to drop their attempts in August 2000 for lack of sufficient evidence. Due to the lack of computer crime laws at the time, it was impossible for other countries such as the United States to extradite the suspect: International principles of dual criminality require equivalent laws in both jurisdictions before extradition can proceed.

By October 2000, de Guzman had refused to take responsibility for writing the worm and publicly stated, “I admit I create viruses, but I don’t know if it’s one of mine.... If the source code was given to me, I could look at it and see. Maybe it is somebody else’s, or maybe it was stolen from me.”⁶⁹

The “I Love You” case was a wake-up call for the international community to think about standardizing computer crime laws around the globe.⁷⁰

2.11.7 July 2010 Stuxnet. In July 2010, reports surfaced of a zero-day threat to SCADA systems using Siemens AG’s Simatic WinCC and PCS 7 software. Analysts found that the Stuxnet worm was designed for industrial espionage; however, the same techniques could have been used for sabotage. Experts expressed concern that the worm was signed using valid digital certificates from Taiwanese companies and that the complex code implied considerable knowledge of the SCADA software.⁷¹ Further analysis of the malware code suggested that the software was developed by the United States and Israel and used at least as early as November 2007.⁷²

2.12 SPAM. Chapter 20 in this *Handbook* includes a detailed history of unsolicited commercial email and the reason it is called *spam*. This section looks solely at a seminal abuse of the USENET in 1994 and trends in spam over the next decade.

2.12.1 1994: Green Card Lottery Spam. On April 2, 1994, Laurence A. Canter and Marthas S. Siegel posted an advertisement for legal services connected to the U.S. Government’s Green Card Lottery to over 6,000 USENET groups. Instead of cross-posting their commercial message, they used a script to post a copy of the message separately to every group. The former method would have shown the message to USENET users once; Canter and Siegel’s abuse of the USENET made their ad show up in every affected group to which users subscribed.⁷³

DENIAL OF SERVICE 2 · 21

Reaction worldwide was massive. Automated cancelbots trolled the USENET deleting the unwanted messages; the attorneys' ISP was so overloaded with email complaints that its servers crashed. Canter and Siegel were reviled in postings and newspaper articles.⁷⁴ Their unsavory backgrounds were posted in discussion groups, including details of disciplinary hearings before the Florida Bar and accusations of dishonesty and unprofessional behavior.⁷⁵

Unfazed, the couple published a book about how to abuse the Internet using spam and defended their actions in interviews as an expression of freedom of speech; they dismissed critics as "wild-eyed zealots" or as commercial interests intent on controlling the Internet for their own gain.⁷⁶

Canter was eventually disbarred in Tennessee, in part for his spamming.⁷⁷ He remained unrepentant; in 2002, he spammed 50,000 K–12 teachers with an advertisement for a book whose title he liked so he could harvest payments for referrals from Amazon.⁷⁸

2.12.2 Spam Goes Global. Over the next decade, the incidence of spam grew explosively. By 2007, spam watchers and anti-spam companies reported that around 88 percent of all email traffic on the Internet was spam. Spammers caused so much irritation that companies developed software and hardware solutions for filtering email by content. Spammers responded by increasing the number of images in their spam, making content filtering more difficult. At one point, the amount of spam grew 17 percent between one day and the next as spammers began pumping PDF files into spam pipelines.⁷⁹

Botnets spawned through infected zombie machines established rogue SMTP nodes using innocent (and ignorant) PC users' computers and persistent high-speed Internet connections.⁸⁰ Spam currently provides a major vector for fraud by deceit, including in particular 4-1-9 advance fee fraud and phishing attacks.⁸¹ Advance-fee fraud usually consists of enticements to participate in the theft of ill-gotten gains such as bank deposits belonging to dead people or stolen from poor countries; the dupes who agree to participate in such illegality are promised millions of dollars—only to be told that they suddenly have to send cash for unexpected bribes or fees. If they do so, they are asked for more ... and more ... and more. Phishing involves sending email messages that are supposed to look like official, usually alarming, warnings from banks and other institutions; victims click on links that look like one thing but actually go to the criminals' Websites. There the victims cheerfully type in their user identification, passwords, bank account numbers, and all manner of other confidential information useful for identity theft.⁸² Advance-fee fraud and phishing are discussed in Chapter 20 in this *Handbook*.

2.13 DENIAL OF SERVICE. Denial of service results from exhaustion or destruction of necessary resources and is thoroughly discussed in Chapter 18. However, a couple of denial-of-service attackers stand out among all the others in the last two decades: the Unemailer and Mafiaboy.

2.13.1 1996: Unemailer. In August 1996, someone using the pseudonym "johnny [x]chaotic" claimed the blame for a massive mail-bombing run based on fraudulently subscribing dozens of victims to hundreds of mailing lists. The denial of service was the result in part of the naïveté of list managers who accepted subscriptions for any email address from any other email address. In a rambling and incoherent letter posted on the Net, (s)he made rude remarks about famous and not-so-famous

2 · 22 HISTORY OF COMPUTER CRIME

people, whose capacity to receive meaningful email was then obliterated by up to thousands of unwanted messages a day.⁸³ “The first attack, in August, targeted more than 40 individuals, including Bill Clinton and Newt Gingrich and brought a torrent of complaints from the people who found their names sent as subscribers to some 3,000 E-mail lists.”⁸⁴

Someone claiming to be the same “Unemailer” (as the news media labeled him or her in reference to the Unabomber) launched a similar mass-subscription mail-bombing run in late December.

This attack is estimated to involve 10,139 listservs groups, 3 times greater than the one that took place in the summer, also at xchaotic’s instigation. If each mailing list in this attack sent the targeted individuals just a modest 10 letters to the subscribers’ computers those individuals would receive more than 100,000 messages. If each listing system sent 100 messages—and many do—then the total messages could tally 1,000,000.⁸⁵

In December, the attacker(s) sneered at list administrators for failing to use authentication before allowing subscriptions and wrote that they would continue their attacks until practices changed.⁸⁶

Partly as a result of the Unemailer’s depredations, list administrators did in fact change their practices—not that anyone thanked Johnny [x]chaotic for his method of persuasion.

2.13.2 2000: MafiaBoy. On February 8, 2000, Yahoo.com suffered a three-hour flood from a distributed denial-of-service (DDoS) attack and lost its capacity to serve Web pages to visitors. The next day, the same technique was extended to Amazon.com, eBay.com, Buy.com, and CNN.com.⁸⁷ Later information also showed that Charles Schwab, the online stock brokerage, had been seriously impeded in serving its customers because of the DDoS. Buy.com managers were particularly disturbed because the attack occurred on the day of their initial public offering. As a result of the attacks, a number of firms formed a consortium to fight DDoS attacks.⁸⁸

Investigation by the RCMP and the FBI located a 15-year-old child in west-end Montreal who used a modem to control zombies in his DDoS escapade:

On April 15, 2000, the RCMP arrested a Canadian juvenile known as Mafiaboy for the February 8th DDoS attack on CNN in Atlanta, Georgia. On August 3, 2000, Mafiaboy was charged with 64 additional counts. On January 18, 2001, Mafiaboy appeared before the Montreal Youth Court in Canada and pleaded guilty to 56 counts. These counts included mischief to property in excess of \$5,000 against Internet sites, including CNN.com, in relation to the February 2000 attacks. The other counts related to unauthorized access to several other Internet sites, including those of several US universities. On September 12, 2001, Mafiaboy appeared before the Montreal Youth Court in Canada and was sentenced to eight months “open custody,” one year probation, and restricted use of the Internet.⁸⁹

MafiaBoy’s name was not released by Canadian authorities because of Canadian laws protecting juveniles, although several U.S. reporters distributed his identity in their publications. His chief contribution to the history of computer crime was to demonstrate asymmetric warfare in cyberspace.⁹⁰ His actions showed that even an ignorant child with little knowledge of computing could use low-tech hardware and tools available to anyone on the Internet to cripple major organizations.

2.14 HACKER UNDERGROUND. Newcomers to the field of information assurance will encounter references to the computer underground in texts, articles, and

HACKER UNDERGROUND 2 · 23

discussions. The sections that follow provide thumbnail sketches of some of the key groups and events in the shadowy world of criminal hacking (known as *black hats*, in contrast to *white hats*, who are law enforcement and establishment security experts), and the intermediate range of well-intentioned rebels who use unorthodox means to challenge corporations and governments over what they see as security failings (these people are often called *gray hats*).

2.14.1 1981: Chaos Computer Club. On September 12, 1981, a group of German computer enthusiasts with a strong radical political orientation formed the Chaos Computer Club (CCC) in Hamburg.⁹¹ One of their first achievements was to demonstrate a serious problem in the Bundespost's (German post office) new Bilschirmtext (BTX) interactive videotext service in 1984, not long after the service was announced.⁹² The CCC used security flaws in BTX to transfer a sizable amount of money into their own bank account through a script that ran overnight as a demonstration to the press (returning the money publicly).

After the Legion of Underground (LoU) announced on January 1, 1999, that they would attack and disable the computer systems of the People's Republic of China and of Iraq, a coalition of hacker organizations including the CCC announced opposition to the move. "We strongly oppose any attempt to use the power of hacking to threaten or destroy the information infrastructure of a country, for any reason," the coalition said. "Declaring war against a country is the most irresponsible thing a hacker group could do. This has nothing to do with hacktivism or hacker ethics and is nothing a hacker could be proud of," the coalition said in the statement.

The CCC has, in general, challenged the general view that "hacker" necessarily means "criminal hacker."⁹³ Their annual Chaos Communications Conferences have proven to be a site of technology exchange and serious discussion of information security issues. Their continued commitment to the rule of law (except where their own activities are concerned), and their willingness to engage authorities in the courts when necessary has gained them an unusual degree of credibility and acceptance in the information security community as relatively pale-gray hats.⁹⁴

2.14.2 1982: The 414s. One morning in June 1982, a system administrator for a DEC VAX 11/780 minicomputer at the Memorial Sloan-Kettering Cancer Center in Manhattan found his system down. Investigation led to the discovery that his and dozens of other systems around the country were being hacked by Milwaukee-area teenagers and others aged 15 to 22. The youths called themselves the 414s after the Milwaukee area code.

Using home computers connected to ordinary telephone lines, they had been breaking into computers across the U.S. and Canada, including one at a bank in Los Angeles, another at a cement company in Montreal and, ominously, an unclassified computer at a nuclear weapons laboratory in Los Alamos, [New Mexico].⁹⁵

In March 1984, "two members of Milwaukee's 414 Gang ... pleaded guilty to misdemeanor charges of making obscene or harassing phone calls. Maximum sentence for each charge: six months in jail and a \$500 fine."⁹⁶

2.14.3 1984: Cult of the Dead Cow. Another influential criminal-hacker group is the Cult of the Dead Cow (cDc), which used to sport amusing (although intentionally offensive to some) cartoons such as that of a crucified cow.⁹⁷ The cDc

2 · 24 HISTORY OF COMPUTER CRIME

was noted for its consistent use of humor and parody; for example, “Swamp Rat’s” 1985 article on building “The infamous ... GERBIL FEED BOMB” included instructions such as “Light the fuse if you put one in. If you dropped a match into it, then go to the nearest phone, dial ‘911’ and tell the nice people that you have a large number of glass shards embedded in your lower body. An ambulance should be there soon.”⁹⁸

The cDc became important proponents of hactivism in the 1990s—the use of criminal hacking techniques for political purposes. They also released a number of hacking tools, of which Back Orifice (BO) and especially Back Orifice 2000 (BO2K) were notorious examples. BO2K was ostensibly a remote administration tool but was in fact a Trojan that ran in stealth mode and allowed remote control of infected machines.⁹⁹ Some observers felt that presenting BO2K as a legitimate tool was another instance of cDc’s satirical bent: The idea that anyone would consider software written by criminal hackers as a trustworthy administration tool struck them as ludicrous.

2.14.4 1984: 2600: The Hacker Quarterly. Eric Corley founded *2600: The Hacker Quarterly* in 1984. This publication has become a standard-bearer for proponents of criminal hacking. The magazine has published a steady stream of explanations of how to exploit specific vulnerabilities in a wide range of operating systems and application environments. In addition, the editor’s political philosophy has influenced more than one generation of black-hat and gray-hat hackers:

In the worldview of *2600*, the tiny band of technocrat brothers (rarely, sisters) are a besieged vanguard of the truly free and honest. The rest of the world is a maelstrom of corporate crime and high-level governmental corruption, occasionally tempered with well-meaning ignorance. To read a few issues in a row is to enter a nightmare akin to Solzhenitsyn’s, somewhat tempered by the fact that *2600* is often extremely funny.¹⁰⁰

2.14.5 1984: Legion of Doom. The DC Comics empire created an animated cartoon series called *Super Friends* that appeared in 1973; it starred various DC Comics heroes, such as Superman, Aquaman, Wonder Woman, and Batman.¹⁰¹ In a follow-up series called *Challenge of the Super Friends* that ran from 1978 through 1979, the archenemies of these heroes were a group known as the *Legion of Doom*, which included Lex Luthor, archenemy of Superman.¹⁰² A group of phone phreakers who later turned to criminal hacking called themselves the Legion of Doom (LOD); their founder called himself “Lex Luthor.” Another major member was Loyd Blankenship (“The Mentor”).

Bruce Sterling describes the LOD as an influential hacker underground group of the 1980s and one of the earliest to capitalize on regular publication of their findings of vulnerabilities and exploits in the phone system and then in computer networks:

LOD members seemed to have an instinctive understanding that the way to real power in the underground lay through covert publicity. LOD were flagrant. Not only was it one of the earliest groups, but the members took pains to widely distribute their illicit knowledge. Some LOD members, like “The Mentor,” were close to evangelical about it. *Legion of Doom Technical Journal* began to show up on boards throughout the underground.

LOD Technical Journal was named in cruel parody of the ancient and honored *AT&T Technical Journal*. The material in these two publications was quite similar—much of it, adopted from public journals and discussions in the telco community. And yet, the predatory attitude of LOD made even its most innocuous data seem deeply sinister; an outrage; a clear and present danger.¹⁰³

HACKER UNDERGROUND 2 · 25

In the later 1980s, the LOD actually helped law enforcement on occasion by restraining malicious hackers.

One of the best-known members was Chris Goggans, whose handle was “Erik Bloodaxe”; he was also an editor of *Phrack* and later became part of the Masters of Deception (MOD), which was involved in a conflict with LOD in 1990 and 1991 known in hacker circles as “The Great Hacker War.”¹⁰⁴

Another well-known hacker who started in LOD and moved to MOD was Mark Abene (“Phiber Optik”), who was eventually imprisoned for a year after pleading guilty in federal court to conspiracy and unauthorized access to federal-interest computers (a violation of 18 USC 1030(a), the Computer Fraud and Abuse Act of 1986).¹⁰⁵ Abene’s punishment was the subject of much protest in the hacker community and elsewhere.¹⁰⁶

2.14.6 1985: *Phrack*. *Phrack* began publishing in November 1985. With a new issue every month or two at first, the electronic magazine continued uninterrupted distribution of technical information and rants. The uncensored commentary provided a fascinating glimpse of some of the personalities and worldviews of its contributors and editors, including Taran King and Craig Neidorf (later to become famous as “Knight Lightning” and for his involvement in an abortive prosecution involving Bell-South documents). For example, *Phrack* published what became known as the “Hacker Manifesto”—held up by criminal hackers as a light unto the nations (“Written almost 15 years ago by The Mentor, this should be taped up next to everyone’s monitor to remind them who we are, this rang true with Hackers, but it now rings truth to the internet generation.”¹⁰⁷), but viewed with skepticism by security professionals. It read in part:

This is our world now ... the world of the electron and the switch, the beauty of the baud. We make use of a service already existing without paying for what could be dirt-cheap if it wasn’t run by profiteering gluttons, and you call us criminals. We explore ... and you call us criminals. We seek after knowledge ... and you call us criminals. We exist without skin color, without nationality, without religious bias ... and you call us criminals. You build atomic bombs, you wage wars, you murder, cheat, and lie to us and try to make us believe it’s for our own good, yet we’re the criminals.

Yes, I am a criminal. My crime is that of curiosity. My crime is that of judging people by what they say and think, not what they look like. My crime is that of outsmarting you, something that you will never forgive me for.

I am a hacker, and this is my manifesto. You may stop this individual, but you can’t stop us all ... after all, we’re all alike.¹⁰⁸

In the 1990s, publication frequency faltered, falling to once every three to six months until the editors announced the final issue, #63, for August 2005. However, publication resumed under new editorial leadership in May 2007 with issue 64; given that issue 65 did not come out until April 2008, the magazine’s heyday is presumably past.

2.14.7 1989: Masters of Deception. The Masters of Deception (MOD) were a New York hacker group active from about 1989 through 1992.¹⁰⁹ Among the most notorious criminal hackers in the group was “Phiber Optik” (Mark Abene, born in 1972), who was unusually visible in the media:

Phiber Optik in particular was to seize the day in 1990. A devotee of the 2600 circle and stalwart of the New York hackers’ group “Masters of Deception,” Phiber Optik was a splendid exemplar of the computer intruder as committed dissident. The eighteen-year-old Optik, a high-school dropout and part-time computer repairman, was young, smart, and ruthlessly obsessive, a sharp-dressing, sharp-talking digital dude who was utterly and airily contemptuous of anyone’s rules

2 · 26 HISTORY OF COMPUTER CRIME

but his own. By late 1991, Phiber Optik had appeared in *Harper's*, *Esquire*, the *New York Times*, in countless public debates and conventions, even on a television show hosted by Geraldo Rivera.¹¹⁰

2.14.8 1990: Operation Sundevil. After two years of investigation, on May 7, 8, and 9, 1990, 150 FBI agents, aided by state and local authorities, raided presumed criminal-hacker organizations allegedly involved in credit-card abuse and theft of telephone services. They seized 42 computers and 23,000 disks from locations in 14 cities. Targets were principally sites running discussion boards, some of which were classified as “hacker boards.” However, two years after the raid, there were only three indictments (resulting in three guilty pleas). Evidence began to accumulate that much of the evidence seized in the raids was useless.¹¹¹ Bruce Sterling spent a year and a half researching the operation and concluded that it was largely a propaganda effort:

... An unprecedented action of great ambition and size, Sundevil's motives can only be described as political. It was a public-relations effort, meant to pass certain messages, meant to make certain situations clear: both in the mind of the general public, and in the minds of various constituencies of the electronic community.

First—and this motivation was vital—a “message” would be sent from law enforcement to the digital underground. This very message was recited in so many words by Garry M. Jenkins, the Assistant Director of the US Secret Service, at the Sundevil press conference in Phoenix on May 9, 1990, immediately after the raids. In brief, hackers were mistaken in their foolish belief that they could hide behind the “relative anonymity of their computer terminals.” On the contrary, they should fully understand that state and federal cops were actively patrolling the beat in cyberspace—that they were on the watch everywhere, even in those sleazy and secretive dens of cybernetic vice, the underground boards.¹¹²

2.14.9 1990: Steve Jackson Games. Two months before the Operation Sundevil raids, but (contrary to popular conflation of the two) in a completely separate operation, a role-playing game company called Steve Jackson Games in Austin, Texas, was raided on March 1, 1990. The Secret Service seized computers and disks at the company’s offices and also at the home of one of their employees, Loyd Blankenship—“The Mentor,” formerly of the LOD. Blankenship was writing a role-playing game called GURPS Cyberpunk, which the agents interpreted as “a handbook for computer crime.” Some of the equipment seized in the raid was returned four weeks later; most but not all was returned four months later. The company nearly went bankrupt as a result of the sequestration of critical resources.¹¹³

Outrage in the computing community spread beyond the underground. Mitch Kapor, John Barlow, and John Gilmore founded the Electronic Frontier Foundation in part because of their outrage over the treatment of Steve Jackson Games:

... We got the attorneys involved, and then we asked them to look into what was going on with a variety of government investigations and prosecutions. We identified a couple of particular legal situations, like Craig Neidorf in Chicago and Steve Jackson Games, where there seemed to us to have been a substantial overstepping of bounds by the government and an infringement on rights of free speech and freedom of the press. We were in the process of deciding how to intervene when we also realized very clearly that we didn’t want to be a legal defense fund as that was too narrow. What was really needed was to somehow improve the discourse about how technology is going to be used by society; we need to do things in the area of public education and policy development.¹¹⁴

Steve Jackson Games sued the Secret Service for damages and were awarded \$50,000 in damages and more than \$25,000 in attorney’s fees.¹¹⁵ The case had a lasting effect

HACKER UNDERGROUND 2 · 27

on how law enforcement officials carried investigations of computer crimes and seizure of electronic evidence.

2.14.10 1992: L0pht Heavy Industries. In 1992, a group of computer enthusiasts arranged to store their spare equipment in some rented space in Boston. They collaborated on analysis of vulnerabilities, especially Microsoft product vulnerabilities, and gained a reputation for contributing serious research to the field and for appearing at security conferences. Their “L0phtCrack” program was adopted by many system administrators for testing password files to locate easy-to-guess passwords; members even testified before a Senate Subcommittee on Government Cybersecurity in 1998 (saying they could take down the Internet in half an hour).¹¹⁶ Famous handles from the group included “Brian Oblivion,” “Kingpin,” “Mudge,” “Space Rogue,” “Stefan von Neumann,” “Tan,” and “Weld Pond.”¹¹⁷

The group caused ripples in both the underground and aboveground security communities when their company, L0pht Heavy Industries, was purchased by security services firm @stake, Inc. in 2000. @stake was eventually bought by Symantec in 1994.¹¹⁸

2.14.11 2004: Shadowcrew. Stealing physical credit cards and creating fake ones are part of the criminal technique called “carding.” One of the significant successful investigations and prosecutions of an international credit card fraud ring of the 2000 decade began with the U.S. Secret Service’s *Operation Firewall* in late 2004. The investigators discovered a network of more than 4,000 members communicating through the Internet and conspiring to use phishing, spamming, forged identity documents (e.g., fake driver’s licenses), creation of fake plastic credit cards, resale of gift cards bought with fake credit cards, fencing of stolen goods via eBay, and interstate or international funds transfers using electronic money such as E-Gold and Web Money.

In October 2004, the Department of Justice indicted 19 of the leaders of Shadowcrew.¹¹⁹ By November 2005, 12 of these people had already pleaded guilty to charges of conspiracy and trafficking in stolen credit card numbers with losses of more than \$4 million.¹²⁰

In February 2006, Shadowcrew leader Kenneth J. Flury, 41, of Cleveland Ohio, was sentenced to 32 months in prison with three years of supervised release and \$300,000 in restitution to Citibank.¹²¹ In June 2006, cofounder Andrew Mantovani, 24, of Scottsdale, Arizona, was fined \$5,000 and also received 32 months of prison with three years of supervised release. Five other indicted Shadowcrew criminals were sentenced with him. By that time, a total of 18 of 28 indicted suspects had already pleaded guilty.¹²²

2.14.12 Late 2000s: Russian Business Network (RBN). The Russian Business Network (RBN) may have originated as a legitimate Web hosting company in 2006:

According to internet security company Verisign, which in June published an extensive investigation into the Russian outfit (tinyurl.com/ywvgpg), RBN was registered as an internet site in 2006.

Initially, much of its activity was legitimate. But apparently the founders soon discovered that it was more profitable to host illegitimate activities and started hiring its services to criminals. Verisign says simply that it is now “entirely illegal.” Since then its activities have been monitored by a number of organisations, including the London-based anti-spam group Spamhaus. “RBN is among the world’s worst spammer, child-pornography, malware, phishing

2 · 28 HISTORY OF COMPUTER CRIME

and cybercrime hosting networks,” says a spokesman. “It provides ‘bulletproof’ hosting, but is probably involved in the crime too.”¹²³

A researcher for the Internet Storm Center, “David Bizeul, spent the past three months researching the Russian Business Network (RBN). The RBN is a virtual safe house for Russian criminals responsible for malicious code attacks, phishing attacks, child pornography, and other illicit operations ...” Bizeul’s study is a 70-page report with extensive documentation about the criminal activities of the RBN.¹²⁴ The group has supported malware diffusion, spam, phishing, denial of service, distribution of cyberattack tools, pornography, and child pornography.

A 2011 report by David Goldman included the following useful insights:

“It’s not like the Mafia, it is a Mafia running these operations,” said Karim Hijazi, CEO of botnet monitoring company Unveillance. “The Russian Mafia are the most prolific cybercriminals in the world.”

Organized cybercrime is a truly international affair, but the most advanced attacks tend to stem from Russia. The Russian mob is incredibly talented for a reason: After the Iron Curtain lifted in the 1990s, a number of ex-KGB cyberspies realized they could use their expert skills and training to make money off of the hacked information they had previously been retrieving for government espionage purposes.

Former spies grouped together to form the Russian Business Network, a criminal enterprise that is capable of some truly scary attacks. It’s just one of many organized cybercriminal organizations, but it’s one of the oldest and the largest.

“The Russians have everyone nailed cold in terms of technical ability,” said Greg Hoglund, CEO of cybersecurity company HBGary. “The Russian crime guys have a ridiculous toolkit. They’re targeting end users in many cases, so they have to be sophisticated.”¹²⁵

2.14.13 Anonymous. In 2003, political activists with a penchant for computer skills formed a loose association calling itself *Anonymous* for collaboration in a range of cyberattacks on targets its members disliked. The philosophy of the group explicitly rejects any centralized controls; anyone can claim to be a member of Anonymous.

In 2008, self-identified members of the movement labeling their efforts *Chanol-
ogy*¹²⁶ attacked the Church of Scientology (readers interested in following the reference provided in the end note should be aware that the site is loaded with pornographic advertisements for pornography sites). Members also harassed organizations attempting to strengthen intellectual property laws and enforcement or antipiracy restrictions. Other targets of the nonorganization include the Epilepsy Foundation, hip-hop Web-sites, Sarah Palin’s political campaign, the government of Iran, the government of Australia, and the Tea Party chapter in Oregon.

One of the most publicized campaigns was in support of Julian Assange, leader of the WikiLeaks Foundation, whose group made public more than a million documents classified by the United States and other governments as restricted or secret and revealing embarrassing details of several wars and internal communications among diplomats.

In January 2013, members announced that they would release large amounts of U.S. Government-restricted information. They let the world know about their plans by posting their messages on a hacked U.S. Government Website.¹²⁷

2.14.14 2013: Unlimited Operations. In May 2013, eight criminal hackers, New York City area members of a much larger worldwide ring of cybercriminals calling themselves Unlimited Operations, were charged with theft of more than \$45 million from automated teller machines (ATMs) around the planet. The gang “used

INDUSTRIAL ESPIONAGE 2 · 29

sophisticated intrusion techniques to hack into the systems of global financial institutions, steal prepaid debit card data, and eliminate withdrawal limits. The stolen card data was then disseminated worldwide and used in making fraudulent ATM withdrawals on a massive scale across the globe ...”

In the first phase, the criminals broke into National Bank of Ras Al-Khaimah PSC (RAKBANK) in the United Arab Emirates. Using these compromised data, the criminal network completed more than 4,500 ATM transactions in 20 hours and stole more than \$5 million.

The second phase began “... on the afternoon of February 19 and lasted into the early morning of February 20, 2013. This operation again breached the network of a credit card processor that serviced MasterCard prepaid debit cards, this time issued by the Bank of Muscat, located in Oman.” Total losses from 36,000 transactions in 24 countries netted \$40 million in cash from ATMs.¹²⁸

2.15 INDUSTRIAL ESPIONAGE. Why spend money developing competitive products when you can steal the work once it’s ready to apply? Many firms in countries with little or no rule of law have taken advantage of poor security, outsourcing, and liberal immigration policies to steal intellectual property and compete at a discount with the originators of the ideas.

- In 2001, Junsheng Wang of Bell Imaging Technologies pled guilty to violation of 18 USC 132(a)(2) by stealing trade secrets from Acuson Corporation. The Counterintelligence News and Developments (CIND) report noted, “In pleading guilty, Wang admitted that prior to August 24, 2000, that he took without authorization and copied for Bell Imaging a document providing the architecture for the Sequoia ultrasound machine that contained the trade secrets of Acuson Corporation. According to Wang’s plea agreement, he had been able to obtain access to the Acuson trade secret materials because his wife was employed as an engineer at that company and because she had brought that document into their home. After he had copied the document, he took it with him on business trips to the People’s Republic of China, turning it over to Bell Imaging during 2000.”¹²⁹
- In May 2001, Federal authorities arrested two Lucent scientists and a third man described as their business partner on May 4, charging them with stealing source code for software associated with Lucent’s PathStar Access Server and sharing it with Datang Telecom Technology Co., a Beijing firm majority-owned by the Chinese government. The software is considered a “crown jewel” of the company. Chinese nationals Hai Lin and Kai Xu were regarded as “distinguished members” of Lucent’s staff up until their arrests. The motivation for the theft, according to court documents, was to build a networking powerhouse akin to the “Cisco of China.” The men faced charges of conspiracy to commit wire fraud, punishable by a maximum five years in prison and a \$250,000 fine.¹³⁰ In April 2002, the two were also charged with stealing secrets from four companies in addition to Lucent: Telenetworks, NetPlane Systems, Hughes Software Systems, and Ziatech. An additional Chinese national, Yong-Qing Cheng was also charged. They developed a joint venture with the Datang Telecom Technology Company of Beijing to sell a clone of Lucent’s Path Star data and voice transmission system to Internet providers in China.¹³¹
- In September 2002, the 3DGeo company in Mountain View, CA accused Shan Yanming, an employee of the China National Petroleum Corporation on loan to

2 · 30 HISTORY OF COMPUTER CRIME

the company, of industrial espionage for trying to steal the software designed for using seismic data to map oil deposits. He was caught trying to download corporate data to his personal computer and was arrested by FBI agents.¹³²

- In April 2003, the United States Attorney's Office for the Northern District of California announced that Tse Thow Sun pled guilty on April 9, 2003, to theft of trade secrets. He admitted that in early 2002, while working for a language translation company, he delivered a laptop computer and a hard drive that contained trade secrets and confidential proprietary information to a competitor and asked for \$3 million in payment. Mr. Sun, 32, a citizen of Singapore, was indicted by a federal Grand Jury on April 9, 2002. He was charged with theft of trade secrets, in violation of 18 U.S.C. §1832(a)(3); attempted theft of trade secrets, in violation of 18 U.S.C. §1832(a)(4); and interstate transportation of stolen goods, in violation of 18 U.S.C. §2314. Under the plea agreement, Mr. Sun pled guilty to theft of trade secrets.¹³³
- In May 2003, three Swedish employees of LM Ericsson were charged with espionage for allegedly stealing intellectual property and sending it to Russian spies. “[Afshin] Bavand was arrested Nov. 5, 2002, while talking to a Russian intelligence agent in a Stockholm suburb. Police searched the Russian, who wasn’t identified, and found \$4,000 in cash and Ericsson documents.”¹³⁴
- The series of attacks codenamed *Titan Rain* was discovered by Shawn Carpenter in late 2003. Carpenter noticed a flood of expert hacker activity focusing on data theft from a wide range of “the country’s most sensitive military bases, defense contractors and aerospace companies.” Carpenter discovered that “the attacks emanated from just three Chinese routers that acted as the first connection point from a local network to the Internet.” Carpenter worked with U.S. Army and FBI investigators to learn more about the attacks and the attackers. According to Thornburgh, various analysts judge that “Titan Rain is thought to rank among the most pervasive cyberespionage threats that U.S. computer networks have ever faced.”¹³⁵
- In July 2004, an Indian software engineer employed by a U.S. company’s software development center in India was accused of “zipping up” proprietary software source code for printing identification cards and uploading it to her personal e-mail account. Jolly Technologies shut down its Mumbai operations as a result of the breach of security.¹³⁶
- In 2005 and 2006, EMC filed lawsuits against several employees for allegedly stealing trade secrets.¹³⁷
- In December 2006, two Chinese nationals, Fei Ye and Ming Zhong, pleaded guilty in December 2006 to charges of economic espionage on behalf of the People’s Republic of China. They were arrested in November 2001 with stolen trade secrets in their luggage; the information was taken from Sun Microsystems and Transmeta Corporation. The agents were planning to design a competing microprocessor using the stolen designs; profits were to have been shared with the City of Hangzhou and the Province of Zhejiang. The agents’ company was funded in part by the National High Technology Research and Development Program of China.¹³⁸
- In April 2008, sleeper agent Chi Mak, a naturalized U.S. citizen who lived peacefully in Los Angeles for 20 years, was sentenced to 24.5 years in federal prison for industrial espionage. He stole detailed plans for U.S. Navy equipment including

CONCLUDING REMARKS 2 · 31

submarine propulsion systems and tried to send them to China via his brother and sister-in-law.¹³⁹

- In 2009, Siobhan Gorman, writing in *The Wall Street Journal*, reported as follows:

Cyberespionage have penetrated the U.S. electrical grid and left behind software programs that could be used to disrupt the system, according to current and former national-security officials. The spies came from China, Russia, and other countries, these officials said, and were believed to be on a mission to navigate the U.S. electrical system and its controls. The intruders haven't sought to damage the power grid or other key infrastructure, but officials warned they could try during a crisis or war. "The Chinese have attempted to map our infrastructure, such as the electrical grid," said a senior intelligence official. "So have the Russians." The espionage appeared pervasive across the U.S. and doesn't target a particular company or region, said a former Department of Homeland Security official. "There are intrusions, and they are growing," the former official said, referring to electrical systems. "There were a lot last year."¹⁴⁰

- The Office of the National Counterintelligence Executive (ONCIX) published its Report to Congress on Foreign Economic Collection and Industrial Espionage 2009–2011 with the title "Foreign Spies Stealing U.S. Economic Secrets in Cyberspace." The Executive Summary included this commentary:

Sensitive US economic information and technology are targeted by the intelligence services, private sector companies, academic and research institutions, and citizens of dozens of countries.

- Chinese actors are the world's most active and persistent perpetrators of economic espionage. U.S. private-sector firms and cybersecurity specialists have reported an onslaught of computer network intrusions that have originated in China, but the IC cannot confirm who was responsible.
- Russia's intelligence services are conducting a range of activities to collect economic information and technology from U.S. targets.
- Some U.S. allies and partners use their broad access to U.S. institutions to acquire sensitive U.S. economic and technology information, primarily through aggressive elicitation and other human intelligence (HUMINT) tactics. Some of these states have advanced cyber capabilities.¹⁴¹
- A March 2012 report detailed how a successful supervisory control and data acquisition (SCADA) software company, American Superconductor Corporation (AMSC), was practically destroyed economically by its major customer, the Chinese Sinovel company, which stole its proprietary wind-turbine software and then stopped paying for any further software services.¹⁴²
- By early 2013, Symantec's 2012 *Internet Security Threat Report*, Vol. 18 reported that small businesses were increasingly targeted for cyberattacks and industrial espionage: "In 2012, 50 percent of all targeted attacks were aimed at businesses with fewer than 2,500 employees. In fact, the largest growth area for targeted attacks in 2012 was businesses with fewer than 250 employees; 31 percent of all attacks targeted them."¹⁴³

2.16 CONCLUDING REMARKS. At some point, history becomes current events. At the time of writing (May 2013), the trends we were seeing dimly when the fifth edition of this work was published have become clearer. As the second decade

2 · 32 HISTORY OF COMPUTER CRIME

of the 21st century reaches its midpoint, organized crime has become an integral part of the computer-crime scene—and vice versa. The Russian criminal underworld has increasingly invested in high-technology forms of fraud and also relies on high-tech communications for marketing of criminal undertakings, such as international traffic in drugs, armaments, and slaves. Information warfare has become a real issue as China advances in technology by stealing industrial secrets and capitalizing on the savings in research and development—and seeks growing global power. Terrorist groups cannot ignore the power of asymmetric warfare and must be presumed to be planning attacks on critical infrastructures worldwide. As the global communications network spreads throughout the world, governments, corporations, and individuals will have to increase their collaboration and vigilance to defeat the growing army of computer criminals of every type.

2.17 FURTHER READING

- Anderson, N. *The Internet Police: How Crime Went Online, and the Cops Followed*. W. W. Norton & Company, 2013.
- Banks, M. A. *Web Psychos, Stalkers and Pranksters: How to Protect Yourself in Cyberspace*. Coriolis Group Books, 1997.
- Bequai, A. *Technocrimes: The Computerization of Crime and Terrorism*. Lexington Books, 1987.
- Freedman, D. H., and C. C. Mann. *@Large: The strange case of the world's biggest Internet invasion*. Simon & Schuster, 1997.
- Goodell, J. *The Cyberthief and the Samurai: The True Story of Kevin Mitnick—and the Man Who Hunted Him Down*. Dell, 1996.
- Hafner, K., and J. Markoff. *Cyberpunk: Outlaws and Hackers on the Computer Frontier*. Simon & Schuster, 1991.
- Hitchcock, J. A. *True Crime Online: Shocking Stories of Scamming, Stalking, Murder, and Mayhem*. Information Today, 2012.
- Johnson, M. *Cyber Crime, Security and Digital Intelligence*. Gower Publishing, 2013.
- Levy, S. *Hackers: Heroes of the Computer Revolution*. Doubleday, 1984.
- Littman, J. *The Watchman: The Twisted Life and Crimes of Serial Hacker Kevin Poulsen*. Little, Brown, 1997.
- Menn, J. *Fatal System Error: The Hunt for the New Crime Lords Who Are Bringing Down the Internet*. PublicAffairs, 2010.
- Mitnick, K. *Ghost in the Wires: My Adventures as the World's Most Wanted Hacker*. Little, Brown, 2011.
- Mungo, P. *Approaching Zero: The Extraordinary Underworld of Hackers, Phreakers, Virus Writers, and Keyboard Criminals*. Random House, 1993.
- Parker, D. B. *Fighting Computer Crime: A New Framework for Protecting Information*. John Wiley & Sons, 1998.
- Poulsen, K. *Kingpin: How One Hacker Took Over the Billion-Dollar Cybercrime Underground*. Crown, 2011.
- Power, R. *Tangled Web: Tales of Digital Crime from the Shadows of Cyberspace*. Indianapolis: Que, 2000.
- Shackelford, S. J. *Managing Cyber Attacks in International Law, Business, and Relations: In Search of Cyber Peace*. Cambridge University Press, 2013.
- Shimomura, T., and J. Markoff. *Takedown: The Pursuit and Capture of Kevin Mitnick, America's Most Wanted Computer Outlaw—by the Man Who Did It*. Hyperion, 1996.

NOTES 2 · 33

- Slatalla, M., and J. Quittner. *Masters of Deception: The Gang that Ruled Cyberspace*. HarperCollins, 1995.
- Sterling, B. *The Hacker Crackdown: Law and Disorder on the Electronic Frontier*. Bantam Doubleday Dell, 1992.
- Stoll, C. *The Cuckoo's Egg: Tracking a Spy through the Maze of Computer Espionage*. Simon & Schuster, 1989.
- Taylor, R. W., E. J. Fritsch, J. R. Liederbach, and T. J. Holt. *Digital Crime and Digital Terrorism*, 2nd ed. Prentice-Hall, 2010.
- Webb, W. *You've Been Hacked: 15 Hackers You Hope Your Computer Never Meets*. CreateSpace Independent Publishing Platform, 2013.
- Wells, J. T. *Computer Fraud Casebook: The Bytes that Bite*. Wiley, 2009.

2.18 NOTES

1. Some of the materials in this chapter use text from the author's prior publications, to which he holds the copyright. However, specific attributions or quotation marks in such cases are generally avoided, because changes are extensive and the typographical notations marking the changes would have been intrusive and disruptive.
2. Concordia University, "Who We Are: History," 2008, www.concordia.ca/about/whowere/ourhistory/sgw.php
3. T. Whiteside, *Computer Capers: Tales of Electronic Thievery, Embezzlement, and Fraud* (New York: New American Library, 1978).
4. J. Gehl and S. Douglas, "Survey Reveals Epidemic of Battered PCs," *NewsScan*, June 5, 2001.
5. M. Delio, "Battered Computers: An Epidemic," *Wired*, June 5, 2001, www.wired.com/culture/lifestyle/news/2001/06/44284
6. NIPC/DHS, "Physical Attack Still the Biggest Threat," *Daily Open-Source Threat Report*, April 11, 2003.
7. T. Fricke, "Physical Security of Electronic Voting Terminals," *RISKS* 23, No. 20 (2004), <http://catless.ncl.ac.uk/Risks/23.30.html>
8. Whiteside, *Computer Capers*.
9. S. Spielberg, director, *Catch Me If You Can*, 2002, www.imdb.com/title/tt0264464/
10. R. Bell, *Skywayman: The Story of Frank W. Abagnale, Jr.* (Crime Library: Criminal Minds and Methods, 2008), www.trutv.com/library/crime/criminal_mind/scams/frank_abagnale/index.html
11. Whiteside, *Computer Capers*.
12. T. C. Greene, "Chapter One: Kevin Mitnick's Story." *The Register*, January 13, 2003, www.theregister.co.uk/2003/01/13/chapter_one_kevin_mitnicks_story/
13. J. Littman, *The Fugitive Game: Online with Kevin Mitnick—The Inside Story of the Great Cyberchase* (Boston: Little, Brown, 1996), p. 30.
14. A. N. Mayorkas and T. Mrozek, "Kevin Mitnick Sentenced to Nearly Four Years in Prison; Computer Hacker Ordered to Pay Restitution to Victim Companies Whose Systems Were Compromised," Press Release, U.S. Department of Justice, United States Attorney's Office, Central District of California, August 9, 1999, www.usdoj.gov/criminal/cybercrime/mitnick.htm

2 · 34 HISTORY OF COMPUTER CRIME

15. P. Jacobus, "Mitnick Released from Prison," *CNET News*, September 21, 2000, http://news.cnet.com/Mitnick-released-from-prison/2100-1023_3-235933.html
16. K. D. Mitnick and W. L. Simon, *The Art of Intrusion: The Real Stories Behind the Exploits of Hackers, Intruders & Deceivers* (New York: John Wiley & Sons, 1995). K. D. Mitnick and W. L. Simon, *The Art of Deception: Controlling the Human Element of Security* (Hoboken, NJ: John Wiley & Sons, 2003). J. Long, J. Wiles, and K. D. Mitnick, *No Tech Hacking: A Guide to Social Engineering, Dumpster Diving, and Shoulder Surfing* (Syngress, 2008).
17. E. Corley, director (as "Emmanuel Goldstein"), *Freedom Downtime* (2002), www.imdb.com/title/tt0309614/
18. R. Davies, "Origins of Money and of Banking," 2005, www.projects.ex.ac.uk/RDavies/arian/origins.html
19. "Origin and History of Credit Cards," Financial Web: Credit Cards, 2008, www.finweb.com/banking-credit/origin-and-history-of-credit-cards.html or <http://tinyurl.com/5c2yhj>
20. J. Rosenberg, "The First Credit Card," About.com: 20th-Century History, 2008, <http://history1900s.about.com/od/1950s/a/firstcreditcard.htm> or tinyurl.com/6en9kg
21. B. Hutchins, "Notes on the Fair Credit Billing Act (FCBA)," 2002, www.ftc.gov/os/comments/dncpapercomments/04/lsap7.pdf
22. L. S. Fox, ed., *The Federal Reserve System: Purposes & Functions*, 9th ed. (Washington, DC: Board of Governors of the Federal Reserve System, 2005), www.federalreserve.gov/pf/pdf/pf_1.pdf; Chapter 2, "Consumer and Community Affairs," p. 78 (p. 4 of PDF file). www.federalreserve.gov/pf/pdf/pf_6.pdf
23. "Origin and History of Credit Cards."
24. K. Shorter, "Plastic Payments: Trends in Credit Card Fraud," FBI Law Enforcement Bulletin (June 1997), www.fbi.gov/publications/leb/1997/june971.htm
25. www.usdoj.gov, news release of August 5, 2008.
26. Lynn Langton, "Identity Theft Reported by Households, 2005–2010," Bureau of Justice Statistics, November 30, 2011, www.bjs.gov/index.cfm?ty=pb&detail&iid=2207
27. Javelin Strategy & Research, "2013 IDENTITY FRAUD REPORT: Data Breaches Becoming a Treasure Trove for Fraudsters." Javelin, 2013, www.javelinstrategy.com/brochure/276
28. B. Sterling, *The Hacker Crackdown: Law and Disorder on the Electronic Frontier* (New York: Bantam, 1992). Available free online: www.mit.edu/hacker/hacker.html
29. E. McCracken, "Dial-Tone Phreak," *New York Times*, December 30, 2007, www.nytimes.com/2007/12/30/magazine/30joybubbles-t.html?ex=1356584400&en=8d26486125a53d83&ei=5124&partner=permalink&exprod=permalink or <http://tinyurl.com/5s49cu>
30. John T. Draper home page, www.webcrunchers.com
31. S. Wozniak and G. Smith, *iWoz: Computer Geek to Cult Icon: How I Invented the Personal Computer, Co-Founded Apple, and Had Fun Doing It* (New York: Norton, 2006).
32. R. Rosenbaum, "Secrets of the Little Blue Box," *Esquire Magazine* (October 1971). Reprinted at www.slate.com/articles/technology/the_spectator/2011/10/the_article_that_inspired_steve_jobs_secrets_of_the_little_blue_.html

NOTES 2 · 35

33. The text displayed was available on Poulsen's Website until at least April 5, 2001, according to the Internet Archive. Sometime after that date, the biography was shortened and then sometime on or before December 4, 2002, it disappeared altogether and was replaced by a redirect to search for the string "By Kevin Poulsen" in Google.
34. K. Poulsen, "MySpace Predator Caught by Code," *Wired*, October 16, 2006, www.wired.com/science/discoveries/news/2006/10/71948
35. *Wired* magazine, "Threat Level" blog, www.wired.com/threatlevel
36. B. Trumbore, "Ray Dirks and the Equity Funding Scandal," *Wall Street History*, February 6, 2004, www.stocksandnews.com/wall-street-history.php?aid=MTU3M19XUw==
37. M. Kabay, "Crime, Use of Computers in." In H. Bidgoli, ed., *Encyclopedia of Information Systems*, vol. 1 (New York: Academic Press, 2003), www2.norwich.edu/mkabay/overviews/crime_use_of_computers_in.pdf or <http://tinyurl.com/3wqfxc>
38. D. E. Hoffman, "CIA Slipped Bugs to Soviets: Memoir Recounts Cold War Technological Sabotage," *Washington Post*, February 27, 2004, www.msnbc.msn.com/id/4394002
39. E. D. Shaw, K. G. Ruby, and J. M. Post, "The Insider Threat to Information Systems," *Security Awareness Bulletin* No. 2-98, Department of Defense Security Institute (September 1998), www.ntc.doe.gov/cita/CI_Awareness_Guide/Treason/Infosys.htm
40. CNN.com, Y2K Archive, "Looking at the Y2K Bug," 2000, www.cnn.com/TECH/specials/y2k/ (URL inactive)
41. Gardica-Feijoo, L., and J. R. Wingender, "Y2K: Myth or Reality." *Quarterly Journal of Business and Economics* (Summer 2007), http://findarticles.com/p/articles/mi_qa5466/is_200707/ai_n21295780/pg_1 or <http://tinyurl.com/64w5jm> (URL inactive)
42. United States Department of Homeland Security, "Look Before You Click: Trojan Horses and Other Attempts to Compromise Networks," *Joint Information Bulletin*, December 21, 2005, www.us-cert.gov/reading_room/JIB-Trojan122105.pdf or <http://tinyurl.com/6zwmes> (URL inactive)
43. D. Izenberg, "Trojan Horse Masterminds Being Extradited to Israel," *Jerusalem Post*, January 18, 2006. Available for purchase online: <http://pqasb.pqarchiver.com/jpost/access/972012371.html?did=972012371:972012371&FMT=ABS&FMTS=ABS:FT&type=current&date=Jan+18%2C+2006&author=DAN+IZENBERG&pub=Jerusalem+Post&edition=&startpage=04&desc=%27Trojan+horse%27+heads+extradited+to+Israel> or <http://tinyurl.com/5wlsgz> (URL inactive)
44. W. K. Haskins, "Married Couple Indicted for Corporate Espionage," SCI-TECH TODAY.com, March 7, 2006, www.sci-tech-today.com/story.xhtml?story_id=12100DICT7FG&page=1 or <http://tinyurl.com/3qant> (URL inactive)
45. L. Leyden, "Spyware-for-Hire Couple Plead Guilty: Israeli Prison Looms for Haephratis," *The Register*, March 15, 2006, www.theregister.co.uk/2006/03/15/spyware_trojan_guilty_plea/
46. "Court Hands Hefty Fine and Jail Sentence to Israeli Spyware Couple, Reports Sophos," *Sophos*, March 27, 2006, www.sophos.com/pressoffice/news/articles/2006/03/israelspyduo.html or <http://tinyurl.com/4gx38p>

2 · 36 HISTORY OF COMPUTER CRIME

47. J. Leyden, "Israeli Spyware-for-Hire PIs Jailed," *The Register*, April 29, 2008, www.theregister.co.uk/2008/04/29/spyware-for-hire/
48. R. Stiennon, "Four Private Investigators in the Israeli Trojan Fiasco Sentenced. Finally," *Network World*, "Stiennon on Security," April 30, 2008, www.networkworld.com/community/node/27387
49. J. L. Tkacik, "Trojan Dragon: China's Cyber Threat," Heritage Foundation Backgrounder #2106, February 8, 2008, www.heritage.org/Research/asiaandthepacific/bg2106.cfm
50. T. Claburn, "Operation 'Cisco Raider' Nets \$76 in Fake Gear: The Multiyear Effort to Curb the Flow of Counterfeit Network Hardware into the U.S. and Canada Reflects a Steady Escalation in the War on Intellectual Property Crime," *InformationWeek*, February 29, 2008, www.informationweek.com/operation-cisco-raider-nets-76-million-i/206901053
51. J. Markoff, "Trojan Horse Threat Stalks Pentagon after Bogus Hardware Purchase," *CIO TODAY*, May 12, 2008, www.cio-today.com/story.xhtml?story_id=103006ROXFYH or <http://tinyurl.com/5tvz32> (URL inactive)
52. "Ms Smith," "DHS: Imported Tech Tainted with Backdoor Attack Tools," *NetworkWorld | Privacy and Security Fanatic*, July 12, 2011, www.networkworld.com/community/blog/dhs-imported-tech-tainted-backdoor-attack-too
53. For a detailed and personal view of malware history, see virus expert Roger Thompson's "Malicious Code," Chapter 2 in S. Bosworth and M. E. Kabay, eds. *Computer Security Handbook*, 4th ed. (Hoboken, NJ: John Wiley & Sons, 2002). Also, see Chapter 16 in this *Handbook*.
54. R. H. Zakon, "Hobbes' Internet Timeline v8.2," 1996, www.zakon.org/robert/internet/timeline/
55. "Virus Encyclopedia: History of Malware." Viruslist.com, 2008, www.viruslist.com/en/viruses/encyclopedia?chapter=153310937
56. "Virus Encyclopedia: History of Malware."
57. Clausthal University of Technology homepage (English), www.tu-clausthal.de>Welcome.php.en (URL inactive)
58. Thompson, "Malicious Code."
59. R. Patterson, "Re: IBM Christmas Virus," *Risks Forum Digest* 5, No. 80 (December 21, 1987):1.1, catless.ncl.ac.uk/Risks/5.80.html~subj1.1
60. C. Schmidt and T. Darby "The What, Why and How of the 1988 Internet Worm," 1995, snowplow.org/tom/worm/worm.html
61. Robert Morris MIT faculty biography, www.csail.mit.edu/user/972
62. D. Emm, "Changing Threats, Changing Solutions: A History of Viruses and Antivirus," Viruslist.com (now Securelist.com), April 14, 2008, www.securelist.com/en/analysis?pubid=204791996
63. "The WildList Organization International: Frequently Asked Questions." WildList Organization International, 2008, www.wildlist.org/faq.htm
64. R. Thompson, personal communication, May 25, 2008.
65. M. E. Kabay, "INFOSEC Year in Review 1999," 1999, www2.norwich.edu/mkabay/iyir/1999.PDF (URL inactive)
66. CERT Advisory CA-2000-04 Love Letter Worm. CERT/CC, May 9, 2000, www.cert.org/advisories/CA-2000-04.html

NOTES 2 · 37

67. D. I. Hopper, “Focus of ‘ILOVEYOU’ Investigation Turns to Owner of Apartment,” CNN.com, May 10, 2000, <http://archives.cnn.com/2000/TECH/computing/05/10/i.love.you.03/index.html> or <http://tinyurl.com/4elq2l>
68. “Suspected Creator of ‘ILOVEYOU’ Virus Chats Online,” CNN.com chat transcript, September 26, 2000, <http://archives.cnn.com/2000/TECH/computing/09/26/guzman.chat/>
69. M. Landler, “A Filipino Linked to ‘Love Bug’ Talks about His License to Hack,” *New York Times*, October 21, 2000, <http://query.nytimes.com/gst/fullpage.html?res=990DE5D8113EF932A15753C1A9669C8B63> or <http://tinyurl.com/4b826p>
70. R. G. Smith, “Impediments to the Successful Investigation of Transnational High Tech Crime,” *Trends & Issues in Crime and Criminal Justice*, No. 285 (December 13, 2004), www.crime-research.org/articles/trends-and-issues-in-criminal-justice/ or <http://tinyurl.com/44pn4s>
71. Jaikumar Vijayan, “Stuxnet Renews Power Grid Security Concerns: First Known SCADA Malware Program to Target Control Systems Prompts New Questions about Security of U.S. Power Grid,” *NetworkWorld*, July 26, 2010, www.networkworld.com/news/2010/072610-stuxnet-renews-power-grid-security.html
72. Jim Finkle, “Researchers say Stuxnet was deployed against Iran in 2007.” Reuters, February 26, 2013, www.reuters.com/article/2013/02/26/us-cyberwar-stuxnet-idUSBRE91P0PP20130226
73. A. Lawrence, “Internet Growing Pains—The Canter & Siegel Story,” *Computer Business Review* (June 1994), www.coin.org.uk/roadshow/presentation/canter.html
74. K. K. Campbell, “A NET.CONSPIRACY SO IMMENSE.... Chatting with Martha Siegel of the Internet’s Infamous Canter & Siegel,” 1994, <http://lcs.www.media.mit.edu/people/foner/Essays/Civil-Liberties/Project/green-card-lawyers.html> or <http://tinyurl.com/45f3fe> (URL inactive)
75. D. R. Hilton, “Green Card Lottery—Last Call,” 1994, <http://groups.google.com/group/misc.legal/msg/3416cd3d6cfcdbe> (URL inactive)
76. L. Flynn, “‘Spamming’ on the Internet,” *New York Times*, October 16, 1994, www.nytimes.com/1994/10/16/business/sound-bytes-spamming-on-the-internet.html
77. A. Craddock, “Spamming Lawyer Disbarred,” *Wired*, July 10, 1997, www.wired.com/politics/law/news/1997/07/5060
78. N. Swidey, “Spambusters: Cyberwarriors of many stripes have joined the battle against junk email. But the enemy is wily, elusive—and multiplying,” *Boston Globe*, October 5, 2003, www.boston.com/news/globe/magazine/articles/2003/10/05/spambusters?mode=PF or <http://tinyurl.com/4y3chj>
79. C. Garretson, “The Summer of Spam: Record Growth, Record Irritation,” *Network World*, August 16, 2007, www.networkworld.com/news/2007/081607-spam-summer.html or <http://tinyurl.com/6xoda3> (URL inactive)
80. J. Leyden, “Most Spam Comes from Just Six Botnets,” *The Register*, February 29, 2008, www.theregister.co.uk/2008/02/29/botnet_spam_deluge/
81. See NetworkWorld’s “Security Research Center” for up-to-date news about spam and phishing: www.networkworld.com/topics/security.html

2 · 38 HISTORY OF COMPUTER CRIME

82. T. Espiner, "Police Maintain Uneasy Relations with Cybervigilantes," *CNET News*, January 17, 2007, http://news.cnet.com/Police-maintain-uneasy-relations-with-cybervigilantes/2100-7348_3-6150817.html or <http://tinyurl.com/6fjykr>
83. "The Net's Most Wanted," *CNET News*, August 16, 1996, <http://news.cnet.com/2100-1023-221580.html>
84. L. Z. Koch, "Jacking in from the 'Spam in the Stocking' Port: Unemailer Delivers Christmas Grief," CyberWire Dispatch, December 26, 1996, www.petting-zoo.net/~deadbeef/archive/2122.html
85. Koch, "Jacking in from the 'Spam in the Stocking' Port."
86. "Unemailer Explains Bombings," *CNET News*, December 30, 1996, http://news.cnet.com/Unemailer-explains-bombings/2100-1017_3-258247.html or <http://tinyurl.com/422kgc>
87. M. Richtel and S. Robinson, "Several Web Sites Are Attacked on Day after Assault Shut Yahoo," *New York Times*, February 9, 2000, www.nytimes.com/library/tech/00/02/biztech/articles/09hack.html
88. E. Messmer, "Web Sites Unite to Fight Denial-of-Service War," *NetworkWorld*, September 25, 2000, www.networkworld.com/news/2000/0925userdefense.html?nf&_ref=858966935 or <http://tinyurl.com/4cuvsf>
89. "Today's FBI: Facts and Figures," 2003, www.fbi.gov/libref/factsfigure/factsfiguresapri2003.htm
90. See "The RMA Debate" for resources about "The Revolution in Military Affairs," www.comw.org/rma/fulltext/asymmetric.html (URL inactive)
91. For German-speakers or those with automated translation programs, see "FAQ—Über den Chaos Computer Club," May 27, 2004, www.ccc.de/de/faq
92. T. von Randow, "Bildschirmtext: A Blow against the System," translation from *Die Zeit*, November 30, 1984, www.textfiles.com/news/boh-20f8.txt
93. H. Nissenbaum, "Hackers and the Battle for Cyberspace," *Dissent* (Fall 2002), www.dissentmagazine.org/article/hackers-and-the-battle-for-cyberspace
94. Chaos Computer Club press release, "Chaos Computer Club Takes Legal Proceedings against the Voting Computer in Hesse," January 7, 2008, www.ccc.de/updates/2008/wahlcomputer-hessen
95. P. Elmer-Dewitt, "The 414 Gang Strikes Again: Pranksters disrupt a hospital, and nobody is laughing," *Time* (August 29, 1983), www.time.com/time/magazine/article/09171949797,00.html
96. P. Elmer-Dewitt, "Cracking Down: Hackers face tough new laws," *Time* (May 14, 1984), www.time.com/time/magazine/article/09171955290,00.html
97. At the time of writing (May 2008), the group's site (www.cultdeadcow.com) simply showed the words "BE BACK REAL SOON! / -xXx- cDc loves you with the fervor of a THOUSAND SUNS!! -xXx-" and a link to a YouTube video of a teenager playing a ukulele and singing. Consult the Internet Archive for historical snapshots of the site; web.archive.org/web/*http://www.cultdeadcow.com/ (URL inactive)
98. S. Rat, "The infamous ... GERBIL FEED BOMB: Striking fear into the hearts of model citizens everywhere ..." cDc communications, 1985, http://web.archive.org/web/20050212092311/http://www.cultdeadcow.com/cDc_files/cDc-0001.html or <http://tinyurl.com/44yyth>

NOTES 2 · 39

99. E. Messmer, "Bad Rap for Back Orifice 2000?" CNN.com, July 21, 1999, www.cnn.com/TECH/computing/9907/21/badrap.idg/
100. Sterling, *Hacker Crackdown*.
101. Hanna-Barbera, *Super Friends*, 1973, Internet Movie Database, www.imdb.com/title/tt0069641/
102. Hanna-Barbera, *Challenge of the Super Friends*, 1978, Internet Movie Database, www.imdb.com/title/tt0076994/
103. Sterling, *Hacker Crackdown*.
104. M. Slatalla and J. Quittner, "Gang War in Cyberspace," *Wired*, 2.12 (December 1994), www.wired.com/wired/archive/2.12/hacker.html
105. Datastream Cowboy, "MOD Indicted," *Phrack* 4, No. 40 (July 8, 1992): 13, www.phrack.com/issues.html?issue=40&id=13
106. J. Dibbell, "The Prisoner: Phiber Optik Goes Directly to Jail," *Village Voice* (January 12, 1994), www.juliandibbell.com/texts/phiber.html
107. J. Barone, "Manifesto." TechnoZen.com, 2000, www.technozen.com/manifesto.htm
108. The Mentor, "The Conscience of a Hacker," *Phrack* 1, No. 7(1986): 3; www.phrack.com/issues.html?issue=7&id=3#article
109. M. Slatalla and J. Quittner, *Masters of Deception: The Gang that Ruled Cyberspace* (New York: HarperCollins, 1995).
110. Sterling, *Hacker Crackdown*.
111. D. Charles, "'Innocent' Hackers Want Their Computers Back," *New Scientist*, No. 1820, May 9, 1992, p. 9; www.newscientist.com/article/mg13418201.400-innocent-hackers-want-their-computers-back-.html or <http://tinyurl.com/3vw26e>
112. Sterling, *Hacker Crackdown*.
113. S. Jackson, "SJ Games vs. the Secret Service," 2008, www.sjgames.com/SS/
114. D. Gans and K. Goffman, "Mitch Kapor & John Barlow Interview," Electronic Frontier Foundation, August 5, 1990, http://w2.eff.org/Misc/Publications/John_Perry_Barlow/HTML/barlow_and_kapor_in_wired_interview.html or <http://tinyurl.com/4pgskr>
115. S. Sparks, "Judge's Decision in Steve Jackson Games v. United States Secret Service," March 12, 1993, www.sjgames.com/SS/decision-text.html
116. D. Fisher, "The Long, Strange Trip of the L0pht," SearchSecurity.com, March 17, 2008, <http://searchsecurity.techtarget.com/news/1305880/The-long-strange-trip-of-the-L0pht>
117. M. Fitzgerald, "L0pht in Transition," CSO, April 17, 2007, www.csionline.com/article/print/221192
118. Symantec News Release, "Symantec to Acquire @stake," September 16, 2004, www.symantec.com/press/2004/n040916b.html
119. United States Department of Justice, Computer Crime & Intellectual Property Section, "Shadowcrew Organization Called 'One-Stop Online Marketplace for Identity Theft': Nineteen Individuals Indicted in Internet 'Carding' Conspiracy," October 28, 2004, www.usdoj.gov/criminal/cybercrime/mantovaniIndict.htm (URL inactive)

2 · 40 HISTORY OF COMPUTER CRIME

120. United States Department of Justice, Computer Crime & Intellectual Property Section, "Six Defendants Plead Guilty in Internet Identity Theft and Credit Card Fraud Conspiracy," November 17, 2005, www.usdoj.gov/criminal/cybercrime/mantovaniPlea.htm (URL inactive)
121. G. A. White and R. W. Kern, "Cleveland, Ohio Man Sentenced to Prison for Bank Fraud and Conspiracy," U.S. Department of Justice, Eastern District of Pennsylvania, February 28, 2006, www.usdoj.gov/criminal/cybercrime/flurySent.htm
122. United States Attorney's Office, District of New Jersey, "'Shadowcrew' Identity Theft Ringleader Gets 32 Months in Prison," June 29, 2006, www.usdoj.gov/usao/nj/press/files/mant0629_r.htm (URL inactive)
123. Peter Warren, "Hunt for Russia's Web Criminals: The Russian Business Network—Which Some Blame for 60% of All Internet Crime—Appears To Have Gone To Ground. But, Asks Peter Warren, Has It Really Disappeared?" *The Guardian*, November 15, 2007, www.guardian.co.uk/technology/2007/nov/15/news.crime
124. Bizuel, David. "Russian Business Network study," [bizuel.org](http://bizuel.org/files/RBN_study.pdf), November 20, 2007, [www.bizuel.org/files/RBN_study.pdf](http://bizuel.org/files/RBN_study.pdf)
125. David Goldman, "The Cyber Mafia Has Already Hacked You," CNNMoney, July 7, 2011, http://money.cnn.com/2011/07/27/technology/organized_cybercrime/index.htm
126. Anonymous, "Portal: Anonymous/Chanology," *Encyclopedia Dramatica*, October 29, 2012, <https://encycopediadramatica.se/Portal:Anonymous/Chanology>
127. Gregory Ferenstein, "Anonymous Threatens Massive WikiLeaks-Style Exposure, Announced On Hacked Gov Site," *TechCrunch*, January 26, 2013. <http://techcrunch.com/2013/01/26/anonymous-threatens-massive-wikileaks-style-exposure-announced-on-hacked-gov-site/>
128. U.S. Department of Justice, "Eight Members of New York Cell of Cybercrime Organization Indicted in \$45 million Cybercrime Campaign," U.S. Department of Justice | Eastern District of New York, May 9, 2013, www.justice.gov/usao/nye/pr/2013/2013may09.html
129. "Guilty Pleas in Trade Secret Case," *San Francisco Business Times*, April 27, 2001, www.bizjournals.com/eastbay/stories/2001/04/23/daily42.html
130. Andrew Backover, "Feds: Trio Stole Lucent's Trade Secrets," *USA Today*, May 3, 2001, <http://usatoday30.usatoday.com/life/cyber/tech/2001-05-03-lucent-scientists-china.htm>
131. "Trade-Secret Case Is Expanded," *New York Times*, April 12, 2002, www.nytimes.com/2002/04/12/technology/12LUCE.html
132. John Markoff, "Silicon Valley Concern Says It Thwarted Software Theft," *New York Times*, September 20, 2002, www.nytimes.com/2002/09/20/technology/20SOFT.html
133. U.S. Department of Justice, "Chicago, Illinois Man Pleads Guilty to Theft of Trade Secrets, Offered to Sell Online Interpreter's Information," U.S. Department of Justice | Northern District of California, April 11, 2003, www.justice.gov/criminal/cybercrime/press-releases/2003/sunPlea.htm
134. "Three Charged in Ericsson Spy Investigation in Sweden," *USA Today*, May 8, 2003, http://usatoday30.usatoday.com/tech/news/2003-05-08-ericsson_x.htm

NOTES 2 · 41

135. Nathan Thornburgh, "Inside the Chinese Hack Attack," *Time*, August 25, 2005, www.time.com/time/nation/article/085991098371,00.html
136. John Ribeiro, "Source Code Stolen from U.S. Software Company in India: Jolly Technologies Blamed an Insider for the Theft," *Computerworld*, August 5, 2004, www.computerworld.com/s/article/95045/Source_code_stolen_from_U.S._software_company_in_India?taxonomyId=082
137. Todd Wallack, "EMC Sues Ex-Employees To Guard Trade Secrets," *Boston Business Journal*, October 9, 2006, www.bizjournals.com/boston/stories/2006/10/09/story4.html?page=all
138. U.S. Department of Justice, "Two Men Plead Guilty to Stealing Trade Secrets from Silicon Valley Companies to Benefit China," U.S. Department of Justice | Northern District of California, December 14, 2006, www.justice.gov/criminal/cybercrime/press-releases/2006/yePlea.htm
139. Joby Warrick and Carrie Johnson, "Chinese Spy 'Slept' In U.S. for 2 Decades: Espionage Network Said To Be Growing," *Washington Post*, April 3, 2008, www.washingtonpost.com/wp-dyn/content/story/2008/04/02/ST2008040204050.html
140. Siobhan Gorman, "Electricity Grid in U.S. Penetrated By Spies," *Wall Street Journal*, April 8, 2009, <http://online.wsj.com/article/SB123914805204099085.html>
141. Office of the National Counterintelligence Executive (ONCIX), "Foreign Spies Stealing US Economic Secrets in Cyberspace: Report to Congress on Foreign Economic Collection and Industrial Espionage, 2009–2011," Office of the National Counterintelligence Executive, November 3, 2011, www.ncix.gov/publications/reports/fecie_all/Foreign_Economic_Collection_2011.pdf
142. Michael Riley and Ashlee Vance, "Inside the Chinese Boom in Corporate Espionage," *Bloomberg Businessweek*, March 15, 2012, www.businessweek.com/articles/2012-03-14/inside-the-chinese-boom-in-corporate-espionage
143. Symantec, "Internet Security Threat Report 2013," Symantec, April 15, 2013, www.symantec.com/content/en/us/enterprise/other_resources/b-istr_main-report_v18_2012_21291018.en-us.pdf

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 3

TOWARD A NEW FRAMEWORK FOR INFORMATION SECURITY*

Donn B. Parker, CISSP

3.1 PROPOSAL FOR A NEW INFORMATION SECURITY FRAMEWORK	3·1	3.4.1 Complete List of Information Loss Acts	3·11
3.2 SIX ESSENTIAL SECURITY ELEMENTS	3·3	3.4.2 Examples of Acts and Suggested Controls	3·14
3.2.1 Loss Scenario 1: Availability	3·4	3.4.3 Physical Information and Systems Losses	3·17
3.2.2 Loss Scenario 2: Utility	3·4	3.4.4 Challenge of Complete Lists	3·18
3.2.3 Loss Scenario 3: Integrity	3·5		
3.2.4 Loss Scenario 4: Authenticity	3·5	3.5 FUNCTIONS OF INFORMATION SECURITY	3·19
3.2.5 Loss Scenario 5: Confidentiality	3·6	3.6 SELECTING SAFEGUARDS USING A STANDARD OF DUE DILIGENCE	3·20
3.2.6 Loss Scenario 6: Possession	3·7		
3.2.7 Conclusions about the Six Elements	3·8	3.7 THREATS, ASSETS, VULNERABILITIES MODEL	3·20
3.3 WHAT THE DICTIONARIES SAY ABOUT THE WORDS WE USE	3·9	3.8 CONCLUSION	3·20
3.4 COMPREHENSIVE LISTS OF SOURCES AND ACTS CAUSING INFORMATION LOSSES	3·10	3.9 FURTHER READING	3·23

3.1 PROPOSAL FOR A NEW INFORMATION SECURITY FRAMEWORK.

Information security, historically, has been limited by the lack of a comprehensive, complete, and analytically sound framework for analysis and improvement. The persistence of the classic triad of CIA (confidentiality, integrity, availability) is inadequate to describe what security practitioners include and implement when doing their jobs. We need a new information security framework that is complete, correct, and consistent

*This chapter is a revised excerpt from Donn B. Parker, *Fighting Computer Crime* (New York: John Wiley & Sons, 1998), Chapter 10, “A New Framework for Information Security,” pp. 229–255.

3 · 2 TOWARD A NEW FRAMEWORK FOR INFORMATION SECURITY

to express, in practical language, the means for information owners to protect their information from any adversaries and vulnerabilities.

The current focus on computer systems security is attributable to the understandable tendency of computer technologists to protect what they know best—the computer and network systems rather than the application of those systems. With a technological hammer in hand, everything looks like a nail. The primary security challenge comes from people misusing or abusing information, and often—but not necessarily—using computers and networks. Yet the individuals who currently dominate the information security folk art are neither criminologists nor computer application specialists.

This chapter presents a comprehensive new information security framework that resolves the problems of the existing models. The chapter demonstrates the need for six security elements—availability, utility, integrity, authenticity, confidentiality, and possession—to replace incomplete CIA security (which does not even seem to include security for information that is not confidential) in the new security framework. This new framework is used to list all aspects of security at a basic level. The framework is also presented in another form, the *Threats, Assets, Vulnerabilities Model*, which includes detailed descriptors for each topic in the model. This model supports the new security framework, demonstrating its contribution to advance information security from its current technological stage, and as a folk art, into the basis for an engineering and business art in cyberspace.

The new security framework model incorporates six essential parts:

1. Security elements of information to be preserved are:

- Availability
- Utility
- Integrity
- Authenticity
- Confidentiality
- Possession

2. Sources of loss of these security elements of information:

- Abusers and misusers
- Accidental occurrences
- Natural physical forces

3. Acts that cause loss:

- Destruction
- Interference with use
- Use of false data
- Modification or replacement
- Misrepresentations or repudiation
- Misuse or failure to use
- Location
- Disclosure
- Observation
- Copying

SIX ESSENTIAL SECURITY ELEMENTS 3 · 3

- Taking
- Endangerment

4. Safeguard functions to protect information from these acts:

- Audit
- Avoidance
- Deterrence
- Detection
- Prevention
- Mitigation
- Transference
- Investigation
- Sanctions and rewards
- Recovery

5. Methods of safeguard selection:

- Use due diligence
- Comply with regulations and standards
- Enable business
- Meet special needs

6. Objectives to be achieved by information security:

- Avoid negligence
- Meet requirements of laws and regulations
- Engage in successful commerce
- Engage in ethical conduct
- Protect privacy
- Minimize impact of security on performance
- Advance an orderly and protected society

In summary, this model is based on the goal of meeting owners' needs to protect the desired *security elements* of their information from sources of loss that engage in harmful *acts* and events by applying *safeguard functions* that are selected by accepted *methods* to achieve desired *objectives*. The sections of the model are explained next. It is important to note that security risk, return on security investment (ROSI), and net present value (NPV) based on unknown future losses and enemies and their intentions are not identified in this model, since they are not measurable and, hence, not manageable.

3.2 SIX ESSENTIAL SECURITY ELEMENTS. Six security elements in the proposed framework model are essential to information security. If any one of them is omitted, information security is deficient in protecting information owners. Six scenarios of information losses, all derived from real cases, are used to demonstrate this contention. We show how each scenario involves violation of one, and only one, element of information security. Thus, if we omit that element from information security, we also remove that scenario from the concerns of information security, which would

3 · 4 TOWARD A NEW FRAMEWORK FOR INFORMATION SECURITY

be unacceptable. It is likely that information security professionals will agree that all of these scenarios fall well within the range of the abuse and misuse that we need to protect against.

3.2.1 Loss Scenario 1: Availability. A rejected contract programmer, intent on sabotage, removed the name of a data file from the file directories in a credit union's computer. Users of the computer and the data file no longer had the file available to them, because the computer operating system recognizes the existence of information available for users only if it is named in the file directories. The credit union was shut down for two weeks while another programmer was brought in to find and correct the problem so that the file would be available. The perpetrator was eventually convicted of computer crime.

Except for availability, the other elements of information security—utility, integrity, authenticity, confidentiality, and possession—do not address this loss, and their state does not change in the scenario. The owner of the computer (the credit union) retained possession of the data file. Only the availability of the information was lost, but it is a loss that clearly should have been prevented by information security. Thus, the preservation of availability must be accepted as a purpose of information security.

It is true that good security practice might have prevented the disgruntled programmer from having use of the credit union application system, and credit union management could have monitored his work more carefully. They should not have depended on the technical capabilities and knowledge of only one person, and they should have employed several controls to preserve or restore the availability of data files in the computer, such as by maintaining a backup directory with the names of erased files and pointers to their physical location. The loss might have been prevented, or minimized, through good backup practices, good usage controls for computers and specific data files, use of more than one name to identify and find a file, and the availability of utility programs to search for files by content or to mirror file storage. These safeguards would at least have made the attack more difficult and would have confronted the programmer with the malfeasance of his act.

The severity of availability loss can vary considerably. A perpetrator may destroy copies of a data file in a manner that eliminates any chance of recovery. In other situations, the data file may be partially usable, with recovery possible for a moderate cost, or the user may have inconvenienced or delayed use of the file for some period of time, followed by complete recovery.

3.2.2 Loss Scenario 2: Utility. In this case, an employee routinely encrypted the only copy of valuable information stored in his organization's computer and then accidentally erased the encryption key. The usefulness of the information was lost and could be restored only through difficult cryptanalysis.

Although this scenario can be described as a loss of availability or authenticity of the encryption key, the loss focuses on the usefulness of the information rather than on the key, since the only purpose of the key was to facilitate encryption. The information in this scenario is available, but in a form that is not useful. Its integrity, authenticity, and possession are unaffected, and its confidentiality, unfortunately, is greatly improved.

To preserve utility of information in this case, management should require mandatory backup copies of all critical information and should control the use of powerful protective mechanisms such as cryptography. Management should require security walk-through tests during application development to limit unusable forms of information. It should minimize the adverse effects of security on information use and should

SIX ESSENTIAL SECURITY ELEMENTS 3 · 5

control the types of activities that enable unauthorized persons to reduce the usefulness of information.

The loss of utility can vary in severity. The worst-case scenario would be the total loss of usefulness of the information, with no possibility of recovery. Less severe cases may range from a partially useful state with the potential for full restoration of usefulness at moderate cost.

3.2.3 Loss Scenario 3: Integrity. In this scenario, a software distributor purchased a copy (on DVD) of a program for a computer game from an obscure publisher. The distributor made copies of the DVD and removed the name of the publisher from the DVD copies. Then, without informing the publisher or paying any royalties, the distributor sold the DVD copies in a foreign country. Unfortunately, the success of the program sales was not deterred by the lack of an identified publisher on the DVD or in the product promotional materials.

Because the DVD copies of the game did not identify the publisher that created the program, the copies lacked integrity. (“Integrity” means a state of completeness, wholeness, and soundness, or adhering to a code of moral values.) However, the copies did not lack authenticity, since they contained the genuine game program and only lacked the identity of the publisher, which was not necessary for the successful use of the product. Information utility of the DVD was maintained, and confidentiality and availability were not at issue. Possession also was not at issue, since the distributor bought the original DVD. But copyright protection was violated as a consequence of the loss of integrity and unauthorized copying of the otherwise authentic program.

Several controls can be applied to prevent the loss of information integrity, including using and checking sequence numbers, checksums, and/or hash totals to ensure completeness and wholeness for a series of items. Other controls include performing manual and automatic text checks for required presence of records, subprograms, paragraphs, or titles, and testing to detect violations of specified controls.

The severity of information integrity loss also varies. Significant parts of the information can be missing or misordered (but still available), with no potential for recovery. Or missing or misordered information can be restored, with delay and at moderate cost. In the least severe cases, an owner can recover small amounts of misordered or mislocated information in a timely manner at low cost.

3.2.4 Loss Scenario 4: Authenticity. In a variation of the preceding scenario, another software distributor obtained the program (on DVD) for a computer game from an obscure publisher. The distributor changed the name of the publisher on the DVD and in title screens to that of a well-known publisher, then made copies of the DVD. Without informing either publisher, the distributor then proceeded to distribute the DVD copies in a foreign country. In this case, the identity of a popular publisher on the DVDs and in the promotional materials significantly added to the success of the product sales.

Because the distributor misrepresented the publisher of the game, the program did not conform to reality: It was not an authentic game from the well-known publisher. Availability and utility are not at issue in this case. The game had integrity because it identified a publisher and was complete and sound. (Certainly the distributor lacked *personal* integrity because his acts did not conform to ethical practice, but that is not the subject of the scenario.) The actual publisher did not lose possession of the game, even though copies were deceptively represented as having come from a different publisher.

3 · 6 TOWARD A NEW FRAMEWORK FOR INFORMATION SECURITY

And, although the distributor undoubtedly tried to keep his actions secret from both publishers, confidentiality of the content of the game was not at issue.

What if someone misrepresents your information by claiming that it is his? Violation of CIA does not include this act. A stockbroker in Florida cheated his investors in a Ponzi (pyramid sales) scheme. He stole \$50 million by claiming that he used a super-secret computer program on his giant computer to make profits of 60 percent per day by arbitrage, a stock trading method in which the investor takes advantage of a small difference in prices of the same stock in different markets. He showed investors the mainframe computer at a Wall Street brokerage firm and falsely claimed that it and the information stored therein were his, thereby lending believability to his claims of successful trading.

This stockbroker's scheme was certainly a computer crime, but the CIA elements do not address it as such because its definition of integrity does not include misrepresentation of information. "Integrity" means only that information is whole or complete; it does not address the validity of information. Obviously, confidentiality and availability do not cover misrepresentation either. The best way to extend CIA to include misrepresentation is to use the more general term "authenticity." We can then assign the correct English meaning to the phrase "integrity of information": wholeness, completeness, and good condition. Dr. Peter Neumann at SRI International is correct when he says that information with integrity means that the information is what you expect it to be. This does not, however, necessarily mean that the information is valid (you may expect it to be invalid). "Authenticity" is the word that means conformance to reality.

A number of controls can be applied to ensure authenticity of information. These include confirming transactions, names, deliveries, and addresses; validating products; checking for out-of-range or incorrect information; and using digital signatures and watermarks to authenticate documents.

The severity of authenticity loss can take several forms, including lack of conformance to reality with no recovery possible; moderately false or deceptive information with delayed recovery at moderate cost; or factually correct information with only annoying discrepancies. If the CIA elements included authenticity, with misrepresentation of information as an important associated threat, Kevin Mitnick (the notorious criminal hacker who used deceit as his principal tool for penetrating security barriers) might have faced a far more difficult challenge in perpetrating his crimes. The computer industry might have understood the need to prove computer operating system updates and Web sites genuine, to avoid misrepresentation with fakes before their customers used those fakes in their computers.

3.2.5 Loss Scenario 5: Confidentiality. A thief deceptively obtained information from a bank's technical maintenance staff. He used a stolen key to open the maintenance door of an automated teller machine (ATM) and secretly inserted a radio transmitter that he purchased from a Radio Shack store. The radio received signals from the touch-screen display in the ATM that customers use to enter their personal identification numbers (PINs) and to receive account balance information. The radio device broadcast the information to the thief's radio receiver in his nearby car, which recorded the PINs and account balances on tape in a modified videocassette recorder. The thief used the information to loot the customers' accounts from other ATMs. The police and the Federal Bureau of Investigation caught the thief after elaborate detective and surveillance efforts. He was sentenced to 10 years in a federal prison.

The thief violated the secrecy of the customers' PINs and account balances, and he violated their privacy. Availability, utility, integrity, and authenticity were unaffected in

SIX ESSENTIAL SECURITY ELEMENTS 3 · 7

this violation of confidentiality. The customers' and the bank's exclusive possession of the PINs and account balance information was lost, but not possession per se, because they still held and owned the information. Therefore, this was primarily a case of lost confidentiality.

According to most security experts, confidentiality deals with disclosure, but confidentiality also can be lost by observation, whether that observation is voluntary or involuntary, and whether the information is disclosed or not disclosed. For example, if you leave sensitive information displayed on an unattended computer monitor screen, you have disclosed it and it may or may not lose its confidentiality. If you turn the monitor off, leaving a blank screen, you have not disclosed sensitive information, but if someone turns the monitor on and reads its contents without permission, then confidentiality is lost by observation. We must prevent both disclosure and observation in order to protect confidentiality.

Controls to maintain confidentiality include using cryptography, training employees to resist deceptive social engineering attacks intended to obtain their technical knowledge, and controlling the use of computers and computer devices. Good security also requires that the cost of resources for protection not exceed the value of what may be lost, especially with low incidence. For example, protecting against radio frequency emanations in ATMs (as in this scenario) is probably not advisable, considering the cost of shielding and the paucity of such high-tech attacks.

The severity of loss of confidentiality can vary. The worst-case scenario loss is when a party with the intent and ability to cause harm observes a victim's sensitive information. In this case, unrecoverable damage may result. But information also may be known to several moderately harmful parties, with a moderate loss effect, or be known to one harmless, unauthorized party with short-term recoverable effect.

3.2.6 Loss Scenario 6: Possession. A gang of burglars aided by a disgruntled, recently fired operations supervisor broke into a computer center and stole tapes and disks containing the company's master files. They also raided the backup facility and stole all backup copies of the files. They then held the materials for ransom in an extortion attempt against the company. The burglary resulted in the company's losing possession of all copies of the master files as well as the media on which they were stored. The company was unable to continue business operations. The police eventually captured the extortionists with help from the company during the ransom payment, and they recovered the stolen materials. The burglars were convicted and served long prison sentences.

Loss of possession occurred in this case. The perpetrators delayed availability, but the company could have retrieved the files at any time by paying the ransom. Alternatively, the company could have re-created the master files from paper documents, but at great cost. Utility, integrity, and authenticity were not issues in this situation. Confidentiality was not violated because the burglars had no reason to read or disclose the files. Loss of ownership and permanent loss of possession would have been accomplished if the perpetrators had never returned the materials or if the company had stopped trying to recover them.

The security model must include protecting the possession of information so as to prevent theft, whether the information is confidential or not. Confidentiality, by definition, deals only with secret information that people may possess. Our increasing use of computers magnifies this difference; huge amounts of information are possessed for automated use and not necessarily held confidentially for only specified people to know. Computer object programs are examples of proprietary but not confidential

3 · 8 TOWARD A NEW FRAMEWORK FOR INFORMATION SECURITY

information we do not know but possess by selling, buying, bartering, giving, receiving, and trading until we ultimately control, transport, and use them. We have incorrectly defined possession if we include only the protective efforts for confidential material.

We protect the possession of information by preventing people from unauthorized taking, from making copies, and from holding or controlling it—whether confidentiality is involved or not. The loss of possession of information also includes the loss of control of it, and may allow the new possessor to violate its confidentiality at will. Thus, loss of confidentiality may accompany loss of possession. But we must treat confidentiality and possession separately to determine what actions criminals might take and what controls we need to apply to prevent their actions. Otherwise, we may overlook a particular threat or an effective control. The failure to anticipate a threat and vulnerability is one of the greatest dangers we face in security.

Controls that can protect the possession of information include using copyright laws, implementing physical and logical usage limitations, preserving and examining computer audit logs for evidence of stealing, inventorying tangible and intangible assets, using distinctive colors and labels on media containers, and assigning ownership to enforce accountability of organizational information assets.

The severity of loss of possession varies with the nature of the offense. In a worst-case scenario, a criminal may take information, as well as all copies of it, and there may be no means of recovery—either from the perpetrator or from other sources such as paper documentation. In a less harmful scenario, a criminal might take information for some period of time but leave some opportunity for recovery at a moderate cost. In the least harmful situation, an owner could possess more than one copy of information, leaving open the possibility of recovery from other sources (e.g., backup files) within a reasonable period of time.

3.2.7 Conclusions about the Six Elements. We need to understand some important differences between integrity and authenticity. For one, integrity deals with the intrinsic condition of information, while authenticity deals with the extrinsic value or meaning relative to external sources and uses. Integrity does not deal with the meaning of the information with respect to external sources, that is, whether the information is timely and not obsolete. Authenticity, in contrast, concerns the question of whether information is genuine or valid and not out of date with respect to its potential use. A user who enters false information into a computer possibly has violated authenticity, but as long as the information remains unchanged, it has integrity. An information security technologist who designs security into computer operating systems is concerned only with application information integrity because the designer cannot know if any user is entering false information. In this case, the security technologist's job is to ensure that both true and false information remain whole and complete. It is the information owner, with guidance from an information security advisor, who has the responsibility of ensuring that the information conforms to reality—in other words, that it has authenticity.

Some types of loss that information security must address require the use of all six elements of the framework model to determine the appropriate security to apply. Each of the six elements can be violated independently of the others, with one important exception: A violation of confidentiality always results in loss of exclusive possession, at the least. Loss of possession, however—even exclusive possession—does not necessarily result in loss of confidentiality.

Other than that exception, the six elements are unique and independent, and often require different security controls. Maintaining the availability of information does

WHAT THE DICTIONARIES SAY ABOUT THE WORDS WE USE 3 · 9

not necessarily maintain its utility; information may be available but useless for its intended purpose, and vice versa. Maintaining the integrity of information does not necessarily mean that the information is valid, only that it remains the same or, at least, whole and complete. Information can be invalid and, therefore, without authenticity, yet it may be present and identical to the original version and, thus, have integrity. Finally, who is allowed to view and know information and who possesses it are often two very different matters.

Unfortunately, the written information security policies of many organizations do not acknowledge the need to address many kinds of information loss. This is because their policies are limited to achieving CIA. To define information security completely, the policies must address all six elements presented. Moreover, to eliminate (or at least reduce) security threats adequately, all six elements need to be considered to ensure that nothing is overlooked in applying appropriate controls. These elements are also useful for identifying and anticipating the types of abusive actions that adversaries may take—before such actions are undertaken.

For simplification and ease of reference, we can pair the six elements into three double elements, which should be used to identify threats and select proper controls, and we can associate them with synonyms so as to facilitate recall and understanding:

availability and utility	→ usability and usefulness
integrity and authenticity	→ completeness and validity
confidentiality and possession	→ secrecy and control

Availability and utility fit together as the first double element. Controls common to these elements include secure location, appropriate form for secure use, and usability of backup copies. Integrity and authenticity also fit together; one is concerned with internal structure and the other with conformance to external facts or reality. Controls for both include double entry, reasonableness checks, use of sequence numbers and checksums or hash totals, and comparison testing. Control of change applies to both as well. Finally, confidentiality and possession go together because, as discussed, they are interrelated. Commonly applied controls for both include copyright protection, cryptography, digital signatures, escrow, and secure storage.

The order of the elements here is logical, since availability and utility are necessary for integrity and authenticity to have value, and these first four elements are necessary for confidentiality and possession to have material meaning.

3.3 WHAT THE DICTIONARIES SAY ABOUT THE WORDS WE USE. CIA would be adequate for security purposes if the violation of confidentiality were defined to be anything done *with* information, if integrity were defined to be anything done *to* information, and if availability were to include utility, but these definitions would be incorrect and are not understood by many people. Information professionals are already defining the term “integrity” incorrectly, and we would not want to make matters worse. These definitions of security and the elements are relevant abstractions from *Webster’s Third New International Dictionary* and *Webster’s Collegiate Dictionary*, 10th edition.

Security—freedom from danger, fear, anxiety, care, uncertainty, doubt; basis for confidence; measures taken to ensure against surprise attack, espionage, observation, sabotage; resistance of a cryptogram to cryptanalysis usually measured by the time and effort needed to solve it.

Availability—present or ready for immediate use.

3 · 10 TOWARD A NEW FRAMEWORK FOR INFORMATION SECURITY

Utility—useful, fitness for some purpose.

Integrity—unimpaired or unmarred condition; soundness; entire correspondence with an original condition; adherence to a code of moral, artistic or other values; the quality or state of being complete or undivided; material wholeness.

Authenticity—quality of being authoritative, valid, true, real, genuine, worthy of acceptance or belief by reason of conformity to fact and reality.

Confidentiality—quality or state of being private or secret; known only to a limited few, containing information whose unauthorized disclosure could be prejudicial to the national interest.

Possession—act or condition of having or taking into one's control or holding at one's disposal; actual physical control of property by one who holds for himself, as distinguished from custody; something owned or controlled.

We lose credibility and confuse information owners if we do not use words precisely and consistently. When defined correctly, the six words are independent (with the exception that information possession is always violated when confidentiality is violated). They are also consistent, comprehensive, and complete. In other words, the six elements themselves possess integrity and authenticity, and therefore they have great utility. This does not mean that we will not find new elements or replace some of them as our insights develop and technology advances. (I first presented this demonstration of the need for the six elements in 1991 at the 14th U.S. National Security Agency/National Institute of Standards and Technology National Computer Security Conference in Baltimore.)

My definitions of the six elements are considerably shorter and simpler than the dictionary definitions, but appropriate for information security.

Availability—usability of information for a purpose

Utility—usefulness of information for a purpose

Integrity—completeness, wholeness, and readability of information and quality being unchanged from a previous state

Authenticity—validity, conformance, and genuineness of information

Confidentiality—limited observation and disclosure of knowledge

Possession—holding, controlling, and having the ability to use information

3.4 COMPREHENSIVE LISTS OF SOURCES AND ACTS CAUSING INFORMATION LOSSES. The losses that we are concerned about in information security come from people who engage in unauthorized and harmful acts against information, communications, and systems, such as embezzlers, fraudsters, thieves, saboteurs, and criminal hackers. They engage in harmful using, taking, misrepresenting, observing, and every other conceivable form of human misbehavior. Natural physical forces such as air and earth movements, heat and cold, electromagnetic energy, living organisms, gravity and projectiles, and water and gases also are threats to information, as are inadvertent human errors.

Extensive lists of losses found in information security often include fraud, theft, sabotage, and espionage, along with disclosure, usage, repudiation, and copying. The first four losses in this list are criminal justice terms at a different level of abstraction from the last four and require an understanding of criminal law, which many information owners and security specialists lack. For example, fraud includes theft only if it is

COMPREHENSIVE LISTS OF SOURCES AND ACTS 3 · 11

performed using deception, and larceny includes burglary and theft from a victim's premises. What constitutes "premises" in an electronic network environment? This is a legal issue.

Many important types of information-related acts, such as false data entry, failure to perform, replacement, deception, misrepresentation, prolongation of use, delay of use, and even the obvious taking copies of information, are frequently omitted from lists of adverse incidents. Each of these losses may require different prevention and detection controls. This may be easily overlooked if our list of potential acts is incomplete—even though the acts that we typically omit are among the most common reported in actual loss experience. The people who cause losses often are aware that information owners have not provided adequate security and have not considered the full array of possible acts. It is, therefore, essential to include all types of potential harmful acts in our lists, especially when unique safeguards are applicable. Otherwise, we are in danger of being negligent, and those to whom we are accountable will view information security as incomplete or poorly conceived and implemented when a loss does occur.

The complete list of information loss acts in the next section is a comprehensive, nonlegalistic list of potential acts resulting in losses to or with information that I compiled from my 35 years in research about computer crime and security. I have simplified it to a single, low level of abstraction to facilitate understanding by information owners and to enable them to select effective controls. The list makes no distinction among the causes of the losses; as such, it applies equally well to accidental and intentional acts. Cause is largely irrelevant at this level of security analysis, as is the underlying intent or lack thereof. (Identifying cause is important at another level of security analysis. We need to determine the sources and motivation of threats in order to identify appropriate avoidance, deterrence, correction, and recovery controls.) In addition, the list makes no distinction between electronic and physical causes of loss, or among spoken, printed, or electronically recorded information.

The acts in the list are grouped to correspond to the six elements of information security outlined previously (e.g., availability and utility, etc.). Some types of acts in one element grouping may have a related effect in another grouping as well. For example, if no other copies of information exist, destroying the information (under *availability*) also may cause loss of possession, and taking (under *possession*) may cause loss of availability. Yet loss of possession and loss of availability are quite different, and may require different controls. I have placed acts in the most obvious categories, where a loss prevention analyst is likely to look first.

Here is an abbreviated version of the complete loss list for convenient use in the information security framework model:

- Destroy
- Interfere with use
- Introduce false data
- Modify or replace
- Misrepresent or repudiate

3.4.1 Complete List of Information Loss Acts

Availability and Utility Losses

- Destruction, damage, or contamination
- Denial, prolongation, acceleration, or delay in use or acquisition

3 · 12 TOWARD A NEW FRAMEWORK FOR INFORMATION SECURITY

- Movement or misplacement
- Conversion or obscuration

Integrity and Authenticity Losses

- Insertion, use, or production of false or unacceptable data
- Modification, replacement, removal, appending, aggregating, separating, or re-ordering
- Misrepresentation
- Repudiation (rejecting as untrue)
- Misuse or failure to use as required

Confidentiality and Possession Losses

- Locating
- Disclosing
- Observing, monitoring, and acquiring
- Copying
- Taking or controlling
- Claiming ownership or custodianship
- Inferring
- Exposing to all of the other losses
- Endangering by exposing to any of the other losses
- Failure to engage in or to allow any of the other losses to occur when instructed to do so

Users may be unfamiliar with some of the words in the lists of acts, at least in the context of security. For example, “repudiation” is a word that we seldom hear or use outside of the legal or security context. According to dictionaries, it means to refuse to accept acts or information as true, just, or of rightful authority or obligation. Information security technologists became interested in repudiation when the Massachusetts Institute of Technology (MIT) developed a secure network operating system for its internal use. The system was named Kerberos, taking the name of the three-headed dog that guarded the underworld in Greek mythology. Kerberos provides a means of forming secure links and paths between users and the computers serving them. Unfortunately, however, in early versions it allowed users to falsely deny using the links. This did not present any particular problems in the academic environment, but it did make Kerberos inadequate for business, even though its other security aspects were attractive. As the use of Kerberos spread into business, repudiation became an issue, and nonrepudiation controls became important.

Repudiation is an important issue in electronic transactions such as in electronic banking, purchases, and auctions used by so many people to automate their purchasing functions and Internet commerce, which require digital signatures, escrow, time stamps, and other authentication controls. I could, for example, falsely claim that I never ordered merchandise and that the order form or electronically transmitted ordering information that the merchant possesses is false. Repudiation is also a growing problem because of the difficulty of proving authorship or the source of electronic missives. And the inverse of repudiation—claiming that an act that did not happen actually did happen, or claiming that false information is true—is also important to security, although it

COMPREHENSIVE LISTS OF SOURCES AND ACTS 3 · 13

is often overlooked. Repudiation and its inverse are both types of misrepresentation, but I include both “repudiation” and “misrepresentation” on the list because they may require different types of controls.

Other words in the list of acts may seem somewhat obscure. For example, we seldom think of prolonging or delaying use as a loss of availability or a denial of use, yet they are losses that are often inflicted by computer virus attacks.

I use the word “locate” in the list rather than “access” because access can be confusing with regard to information security. Although it is commonly used in computer terminology, its use frequently causes confusion, as it did in the trial of Robert T. Morris for releasing the Internet worm of November 2, 1988, and in computer crime laws. For example, access may mean just knocking on a door or opening the door but not going in. How far “into” a computer must you go to “access” it? A perpetrator can cause a loss simply by locating information, because the owner may not want to divulge possession of such information. In this case, no access is involved. For these reasons, I prefer to use the terms “entry,” “intrusion,” and “usage”—as well as “locate”—to refer to a computer as the object of the action. I have a similar problem with the use of the word “disclosure” and ignoring observation as I indicated earlier. “Disclose” is a verb that means to divulge, reveal, make known, or report knowledge to others. We can disclose knowledge by:

- Broadcasting
- Speaking
- Displaying
- Showing
- Leaving it in the presence and view of another person
- Leaving it in possible view where another person is likely to be
- Handing or sending it to another person

Disclosure is what an owner or potential victim might do inadvertently or intentionally, not what a perpetrator does, unless it is the second act after stealing, such as selling stolen intellectual property to another person. Disclosure can be an abuse if a person authorized to know information discloses it to an unauthorized person, or if an unauthorized person discloses knowledge to another person without permission. In any case, confidentiality is lost or is potentially lost, and the person disclosing the information may be accused of negligence, violation of privacy, conspiracy, or espionage.

Loss of confidentiality also can occur by observation, whether the victim or owner disclosed knowledge, resisted disclosure, or did nothing either to protect or to disclose it. Observing is an abuse of listening, spying by eavesdropping, shoulder surfing (looking over another person’s shoulder or overhearing), looking at or listening to a stolen copy of information, or even by tactile feeling, as in the case of reading Braille. We should think about loss of confidentiality as a loss caused by inadvertent disclosure by the victim, observation by the perpetrator, and disclosure by the perpetrator who passes information to a third party. Disclosure and observation of information that is not knowledge converts it into knowledge if cognition takes place. Disclosure always results in loss of confidentiality by putting information into a state where there is no longer any secrecy, but observation results in loss of confidentiality only if cognition or use to the detriment of the owner takes place. Privacy is a right that is a whole other topic that I do not cover here. (This issue is discussed in Chapter 69.)

3 · 14 TOWARD A NEW FRAMEWORK FOR INFORMATION SECURITY

Loss of possession of information (including knowledge) is the loss from the unintended or regretful giving or taking of information. At a higher level of crime description, we call it larceny (theft or burglary) or fraud (when deceit is involved). Possession seems to be most closely associated with confidentiality. The two are placed together in the list because they share the common losses of taking and copying (loss of exclusive possession). I could have used “ownership” of information, since it is a synonym for possession, but “ownership” seems to be not as broad, because someone may rightly or wrongly possess information that is rightfully owned by another. The concepts of owner or possessor of information, along with user, provider, or custodian of information, are important distinctions in security for assigning asset accountability. This provides another reason for including possession in the list.

The act of *endangerment* is quite different from, but applies to, the other losses. It means putting information in harm’s way, or that a person has been remiss (and possibly negligent) by not applying sufficient protection to information, such as leaving sensitive or valuable documents in an unlocked office or open trash bin. Leaving a computer unnecessarily connected to the Internet is another example. Endangerment of information may lead to charges of negligence or criminal negligence and civil liability suits that may be more costly than direct loss incidents. My objectives of security in the framework model invoke a standard of due diligence to deal with this exposure.

The last act in the list—failure to engage in or allow any of the other acts when instructed to do so—may seem odd at first glance. It means that an information owner may require an act resulting in any of the other acts to be carried out. Or the owner may wish that an act be allowed to occur, or information to be put into danger of loss. There are occasions when information should be put in harm’s way for testing purposes or to accomplish a greater good. For example, computer programmers and auditors often create information files that are purposely invalid for use as input to a computer to make sure that the controls to detect or mitigate a loss are working correctly. A programmer bent on crime might remove invalid data in a test input file to avoid testing a control that the perpetrator has neutralized or has avoided implementing for nefarious purposes. The list would surely be incomplete without this type of loss, yet I have never seen it included or discussed in any other information security text.

The acts in the list are described at the appropriate level for deriving and identifying appropriate security controls. At the next lower level of abstraction (e.g., read, write, and execute), the losses would not be so obvious and would not necessarily suggest important controls. At the level that I choose, there is no attempt to differentiate acts that make no change to information from those that do, since these differences are not important for identifying directly applicable controls or for performing threat analyses. For example, an act of modification changes the information, while an act of observation does not, but encryption is likely to be employed as a powerful primary control against both acts.

3.4.2 Examples of Acts and Suggested Controls. The next examples illustrate the relationships between acts and controls in threat analysis. Groups of acts are followed by examples of the losses and applicable controls.

3.4.2.1 Destroy, Damage, or Contaminate. Perpetrators or harmful forces can damage, destroy, or contaminate information by electronically erasing it, writing other data over it, applying high-energy radio waves to damage delicate electronic circuits, or physically damaging the media (e.g., paper, flash memory, or disks) containing it.

COMPREHENSIVE LISTS OF SOURCES AND ACTS 3 · 15

Controls include disaster prevention safeguards such as locked facilities, safe storage of backup copies, and write-usage authorization requirements.

3.4.2.2 Deny, Prolong, Delay Use or Acquisition. Perpetrators can make information unavailable by hiding it or denying its use through encryption and not revealing the means to restore it, or by keeping critical processing units busy with other work, such as in a denial-of-service attack. Such actions would not necessarily destroy the information. Similarly, a perpetrator may prolong information use by making program changes that slow the processing in a computer or by slowing the display of the information on a screen. Such actions might cause unacceptable timing for effective use of the information. Information acquisition may be delayed by requiring too many passwords to retrieve it or by slowing retrieval. These actions can make the information obsolete by the time it becomes available.

Controls include making multiple copies available from different sources, preventing overload of processing by selective allowance of input, or preventing the activation of harmful mechanisms such as computer viruses by using antiviral utilities.

3.4.2.3 Enter, Use, or Produce False Data. Data diddling, my term for false data entry and use, is a common form of computer crime, accounting for much of the financial and inventory fraud. Losses may be either intentional, such as those resulting from the use of Trojan horses (including computer viruses), or unintentional, such as those from input errors.

Most internal controls such as range checks, audit trails, separation of duties, duplicate data entry detection, program proving, and hash totals for data items protect against these threats.

3.4.2.4 Modify, Replace, or Reorder. These acts are often intelligent changes rather than damage or destruction. *Reordering*, which is actually a form of modification, is included separately because it may require specific controls that could otherwise be overlooked. Similarly, *replacement* is included because users might not otherwise include the idea of replacing an entire data file when considering modification. Any of these actions can produce a loss inherent in the threats of entering and modifying information, but including all of them covers modifying data both before entry and after entry, since each requires different controls.

Cryptography, digital signatures, usage authorization, and message sequencing are examples of controls to protect against these acts, as are detection controls to identify anomalies.

3.4.2.5 Misrepresent. The claim that information is something different from what it really is or has a different meaning from what was intended arises in counterfeiting, forgery, fraud, impersonation (of authorized users), and many other deceptive activities. Hackers use misrepresentation in social engineering to deceive people into revealing information needed to attack systems. Misrepresenting old data as new information is another act of this type.

Controls include user and document authentication methods such as passwords, digital signatures, and data validity tests. Making trusted people more resistant to deception by reminders and training is another control.

3.4.2.6 Repudiate. This type of loss, in which perpetrators generally deny having made transactions, is prevalent in electronic data interchange (EDI) and Internet

3 · 16 TOWARD A NEW FRAMEWORK FOR INFORMATION SECURITY

commerce. Oliver North's denial of the content of his email messages is a notable example of repudiation, but as I mentioned earlier, the inverse of repudiation also represents a potential loss.

Repudiation can be controlled most effectively through the use of digital signatures and public key cryptography. Trusted third parties, such as certificate authorities with secure computer servers, provide the independence of notary publics to resist denial of truthful information as long as they can be held liable for their failures.

3.4.2.7 Misuse or Fail to Use as Required. Misuse of information is clearly an act resulting in many information losses. Misuse by failure to perform duties such as updating files or backing up information is not so obvious and needs explicit identification. Implicit misuse by conforming exactly to inadequate or incorrect instructions is a sure way to sabotage systems.

Information usage control and internal application controls that constrain the modification or use of trusted software help to avoid these problems. Keeping secure logs of routine activities can help catch operational vulnerabilities.

3.4.2.8 Locate. Unauthorized use of someone's computer or data network to locate and identify information is a crime under most computer crime statutes—even if there is no overt intention to cause harm. Such usage is a violation of privacy, and trespass to engage in such usage is a crime under other laws.

Log-on and usage controls are major features in many operating systems such as Microsoft Windows and some versions of UNIX as well as in add-on security utilities such as RACF and ACF2 for large IBM computers and many security products for personal computers.

3.4.2.9 Disclose. Preventing information from being revealed to people not authorized to know it is the purpose of business, personal, and government secrecy. Disclosure may be verbal, by mail, or by transferring messages or files electronically or on disks, flash memories, or tape. Disclosure can result in loss of privacy and trade secrets.

Military organizations have advanced protection of information confidentiality to an elaborate art form.

3.4.2.10 Observe or Monitor. Observation, which requires action on the part of a perpetrator, is the inverse of disclosure, which results from actions of a possessor. Workstation display screens, communication lines, and monitoring devices such as recorders and audit logs are common targets of observation and monitoring. Observation of output from printers is another possible source, as is shoulder surfing—the technique of watching screens of other computer users.

Physical entry protection for input and output devices represents the major control to prevent this type of loss. Preventing wiretapping and eavesdropping is also important.

3.4.2.11 Copy. Copy machines and the software *copy* command are the major sources of unauthorized copying. Copying is used to violate exclusive possession and privacy. Copying can destroy authenticity, as when used to counterfeit money or other business instruments.

Location and use controls are effective against copying, as are unique markings such as those used on U.S. currency and watermarks on paper and in computer files.

COMPREHENSIVE LISTS OF SOURCES AND ACTS 3 · 17

3.4.2.12 Take. Transferring data files in computers or networks constitutes taking. So does taking small computers and DVDs or documents for the value of the information stored in them. Perpetrators can easily take copies of information without depriving the owner of possession or confidentiality.

A wide range of physical and logical location controls applies to these losses; most are based on common sense and a reasonable level of due care.

3.4.2.13 Endanger. Putting information into locations or conditions in which others may cause loss in any of the previously described ways clearly endangers the information, and the perpetrator may be accused of negligence, at least.

Physical and logical means of preventing information from being placed in danger are important. Training people to be careful, and holding them accountable for protecting information, may be the most effective means of preventing endangerment.

3.4.3 Physical Information and Systems Losses. Information also can suffer from physical losses such as those caused by floods, earthquakes, radiation, and fires. Although these losses may not directly affect the information itself (e.g., knowledge of operating procedures held in the minds of operators), they can damage or destroy the media and the environment that contain representations of the information. Water, for example, can destroy printed pages and damage magnetic disks; physical shaking or radio frequency radiation can short-out electronic circuits, and fires can destroy all types of media. Overall, physical loss may occur in seven natural ways by application of:

1. Extreme temperature
2. Gases
3. Liquids
4. Living organisms
5. Projectiles
6. Movements
7. Energy anomalies

Each way, of course, comes from specific sources of loss (e.g., smoke or water). And the various ways can be broken down further, to identify the underlying cause of the source of loss. For example, the liquid that destroys information may be water flowing from a plumbing break above the computer workstation, caused in turn by freezing weather. The next list presents examples of each of the seven major sources of physical loss.

1. **Extreme temperature.** Heat or cold. Examples: sunlight, fire, freezing, hot weather, and the breakdown of air-conditioning equipment.
2. **Gases.** War gases, commercial vapors, humid or dry air, suspended particles. Examples: sarin nerve gas, PCBs from exploding transformers, release of Freon from air conditioners, smoke and smog, cleaning fluid, and fuel vapors.
3. **Liquids.** Water, chemicals. Examples: floods, plumbing failures, precipitation, fuel leaks, spilled drinks, acid and base chemicals used for cleaning, and computer printer fluids.

3 · 18 TOWARD A NEW FRAMEWORK FOR INFORMATION SECURITY

4. **Living organisms.** Viruses, bacteria, fungi, plants, animals, and human beings. Examples: sickness of key workers, molds, contamination from skin oils and hair, contamination and electrical shorting from defecation and release of body fluids, consumption of information media such as paper or of cable insulation, and shorting of microcircuits from cobwebs.
5. **Projectiles.** Tangible objects in motion, powered objects. Examples: meteorites, falling objects, cars and trucks, airplanes, bullets and rockets, explosions, and windborne objects.
6. **Movement.** Collapse, shearing, shaking, vibration, liquefaction, flows, waves, separation, slides. Examples: dropping or shaking fragile equipment, earthquakes, earth slides, lava flows, sea waves, and adhesive failures.
7. **Energy anomalies.** Electric surge or failure, magnetism, static electricity, aging circuitry; radiation, sound, light, radio, microwave, electromagnetic, atomic. Examples: electric utility failures, proximity of magnets and electromagnets, carpet static, electromagnetic pulses (EMP) from nuclear explosions, lasers, loudspeakers, high-energy radio frequency (HERF) guns, radar systems, cosmic radiation, and explosions.

Although meteorites, for example, clearly pose little danger to computers, it is nonetheless important to include all such unlikely events in a thorough analysis of potential threats. In general, include every possible act included in a threat analysis. Then consider it carefully; if it is too unlikely, document the consideration and discard the item. It is better to have thought of a source of loss and to have discarded it than to have overlooked an important one. Invariably, when you present a threat analysis to others, someone will try to surprise the developer with another source of loss that has been overlooked.

Inensitive practitioners have ingrained inadequate loss lists in the body of knowledge from the very inception of information security. Proposing a major change at this late date is a bold action that may take significant time to accomplish. However, we must not perpetuate our past inadequacies by using the currently accepted destruction, disclosure, use, and modification (DDUM) as a complete list of losses. We must not underrate or simplify the complexity of our subject at the expense of misleading information owners. Our adversaries are always looking for weaknesses in information security, but our strength lies in anticipating sources of threats and having plans in place to prevent the losses that they may cause.

It is impossible to collect a truly complete list of the sources of information losses that can be caused by the intentional or accidental acts of people. We really have no idea what people may do—now or in the future. We base our lists on experience, but until we can conceive of an act, or until a threat actually surfaces or occurs, we cannot include it on the list. And not knowing the threat means that we cannot devise a plan to protect against it. This is one of the reasons that information security is still a folk art rather than a science.

3.4.4 Challenge of Complete Lists. I believe that my lists of physical sources of loss and information losses are complete, but I am always interested in expanding them to include new sources of loss that I may have overlooked.

While I was lecturing in Australia, for example, a delegate suggested that I had omitted an important category. His computer center had experienced an invasion of field mice with a taste for electrical insulation. The intruders proceeded to chew through the computer cables, ruining them. Consequently, I had to add rodents to my list of

FUNCTIONS OF INFORMATION SECURITY 3 · 19

sources. I then heard about an incident in San Francisco in which the entire evening shift of computer operations workers ate together in the company cafeteria to celebrate a birthday. Then they all contracted food poisoning, leaving their company without sufficient operations staff for two weeks. I combined the results of these two events into a category named “Living Organisms.”

3.5 FUNCTIONS OF INFORMATION SECURITY. The model for information security that I have proposed includes 12 security functions instead of the 3 (prevention, detection, and recovery) included in previous models. These functions describe the activities that information security practitioners and information owners engage in to protect information, as well as the objectives of the security controls that they use. Every control serves one or more of these functions.

Although some security specialists add other functions to the list, such as quality assurance and reliability, I consider these to be outside the scope of information security; other specialized fields deal with them. Reliability is difficult to relate to security except as endangerment when perpetrators destroy the reliability of information and systems, which is a violation of security. Thus, security must preserve a state of reliability but need not necessarily attempt to improve it. Security must protect the auditability of information and systems while, at the same time, security itself must be reliable and auditable. I believe that my security definitions include destruction of the reliability and auditability of information at a high level of abstraction. For example, reliability is reduced when the authenticity of information is put into question by changing it from a correct representation of fact.

Similarly, I do not include such functions as authentication of users and verification in my lists, since I consider these to be control objectives to achieve the 12 functions of information security.

There is a definite logic to the order in which I present the 12 functions in my list. A methodical information security practitioner is likely to apply the functions in this order when resolving security vulnerabilities.

1. Information security must first be independently *audited* in an adversarial manner in order to document its state and to identify its weaknesses and strengths.
2. The practitioner must determine if a security problem can be *avoided* altogether.
3. If the problem cannot be avoided, the practitioner needs to try to *deter* potential abusers or forces from misbehaving.
4. If the threat cannot be avoided or deterred, the practitioner attempts to *detect* its activation.
5. If detection is not assured, then the practitioner tries to *prevent* the act from occurring.
6. If prevention fails and an act occurs, then the practitioner needs to stop it or minimize its harmful effects through *mitigation*.
7. The practitioner needs to determine if *transferring* the responsibility to another individual or department might be more effective at resolving the situation resulting from the attack, or if another party (e.g., an insurer) might be held accountable for the cost of the loss.
8. After a loss occurs, the practitioner needs to *investigate* and search for the individual(s) or force(s) that caused or contributed to the incident as well as for any parties that played a role in it—positively or negatively.
9. When identified, all parties should be *sanctioned or rewarded* as appropriate.

3 · 20 TOWARD A NEW FRAMEWORK FOR INFORMATION SECURITY

10. After an incident is concluded, the victim needs to *recover* or assist with recovery.
11. The stakeholders should take *corrective* actions to prevent the same type of incident from occurring again.
12. The stakeholders must learn from the experience in order to advance their knowledge of information security and *educate* others.

3.6 SELECTING SAFEGUARDS USING A STANDARD OF DUE DILIGENCE.

Information security practitioners usually refer to the process of selecting safeguards as risk assessment, risk analysis, or risk management. Selecting safeguards based on risk calculations can be a fruitless and expensive process. Although many security experts and associations advocate using risk assessment methods, many organizations ultimately find that using a standard of due diligence (or care) is far superior and more practical. Often one sad experience of using security risk assessment is sufficient to convince information security departments and corporate management of their limitations. Security risk is a function of probability or frequency of occurrence of rare loss events and their impact, and neither is sufficiently measurable or predictable for investment in security. Note that risk applies only to rare events. Events such as computer virus attacks or credit card fraud are occurring continuously and are not risks; they are certainties and can be measured, controlled, and managed.

The standard of due diligence approach is simple and obvious; it is the default process that I recommend and that is commonly used today instead of more elaborate “scientific” approaches. The standard of due diligence approach is recognized and accepted by many legal documents and organizations and is documented in numerous business guides. The 1996 U.S. federal statute on protecting trade secrets (18 USC §1831), for example, states in (3)(a) that the owner of information must take “reasonable measures to keep such information secret” for it to be defined as a trade secret. (See Chapter 45.)

3.7 THREATS, ASSETS, VULNERABILITIES MODEL. Pulling all of the aspects together in one place is a useful way to analyze security threats and vulnerabilities and to create effective scenarios to test real information systems and organizations. The model illustrated in Exhibit 3.1 is designed to help readers do this. Users can outline a scenario or analyze a real case by circling and connecting the appropriate descriptors in each column of the model.

In this version of the model, the Controls column lists only the subject headings of control types; a completed model would contain hundreds of controls. If the model is being used to conduct a review, I suggest that the Vulnerabilities section of the model be renamed to Recommended Controls.

3.8 CONCLUSION. The security framework proposed in this chapter represents an attempt to overcome the dominant technologist view of information security by focusing more broadly on all aspects of security, including the information that we are attempting to protect, the potential sources of loss, the types of loss, the controls that we can apply to avoid loss, the methods for selecting those controls, and our overall objectives in protecting information. This broad focus should have two beneficial effects: advancing information security from a narrow folk art to a broad-based discipline and—most important—helping to reduce many of the losses associated with information, wherever it exists.

Exhibit 3.1 Threats, Assets, and Vulnerabilities Model

Threats				Assets			
Offenders Have/Acquire	Abuse/ Misuse	Methods	Losses	Control Objectives	Controls (Types)	Control Guides	
Skills learning technology people	Errors Omissions Negligence Recklessness	External heat, cold gases, air water chemical bacteria viruses people animals insects	Availability and Utility destroy damage contaminate deny prolong accelerate delay move misplace convert obscure	Information spoken printed magnetic electronic optical radio biological	Avoidance Deterrence Prevention Detection Mitigation Sanction Transfer Investigate Recovery Correction	Organization Physical Development Automation Operation Voice Network Access Training Motivation Management Applications Printing Audit	Cost effective Due care Complete Consistent Performance Sustain Automatic Tolerated Consequences Override Failsafe Default Instrument Auditable Nonrepudiate Secrecy Universal Independent Unpredictable Tamperproof Compartment
Knowledge direct indirect	Delinquency Civil Disputes Conspiracy Nature Disruption Desecration Theft Privacy Trespass			Computer			
Resources computer services transport financial				CommLines			
				Networks			
				Facilities			
				Buildings			
				Transport			
				People			
Authority employment contract ownership possession custodian right other							

(continued)

Exhibit 3.1 Threats, Assets, and Vulnerabilities Model (Continued)

Offenders Have/Acquire	Abuse/ Misuse	Methods	Threats		Assets			Vulnerabilities (Missing and Deficient Controls)		
			Losses	Control Objectives	Assets Lost	Controls (Types)	Control Guides			
Motives no intent negligence errors and omissions	Embezzlement Bribery Extortion Racketeering Infringement Plagiarism Piracy Espionage Antitrust Contract Securities Employment Kickbacks Laundering Libel Drugs Pornography Harassment Assault Sex attack Kidnapping Murder Suicide	Radio Atomic Massquerade impersonate spoof Programmed Trojan virus bomb bypass trapdoor Authority violation Active deny service false data entry Passive browse observe Failure omit duty Indirect crime use	append reorder misrepresent Repudiate Fail to use Confidential and Possession locate disclose observe monitor acquire copy take control own infer	Depth Isolate Least Accountability Trust Multifunction Deception Positional Transparent						
Intentional problem solving gain higher ethic										
Extreme Advocacy social political religious										

Source: Donn B. Parker, *Fighting Computer Crime* (New York, NY: John Wiley & Sons, 1998).

FURTHER READING 3 · 23

3.9 FURTHER READING

Parker, D. B. *Fighting Computer Crime: A New Framework for Protecting Information*. Wiley (ISBN 978-0471163787), 1998. 528 pp.

Parker, D. B. "What's wrong with information security and how to fix it. Lecture at the Naval Postgraduate School (2005-04-28)." YouTube. www.youtube.com/watch?v=RW9hOBCSy0g

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 4

HARDWARE ELEMENTS OF SECURITY

Sy Bosworth and Stephen Cobb

4.1	INTRODUCTION	4·2	4.8.5	Dirt and Dust	4·12
4.2	BINARY DESIGN	4·2	4.8.6	Radiation	4·12
4.2.1	Pulse Characteristics	4·2	4.8.7	Downtime	4·12
4.2.2	Circuitry	4·2			
4.2.3	Coding	4·3			
4.3	PARITY	4·4	4.9	DATA COMMUNICATIONS	4·13
4.3.1	Vertical Redundancy Checks	4·4	4.9.1	Terminals	4·13
4.3.2	Longitudinal Redundancy Checks	4·4	4.9.2	Wired Facilities	4·14
4.3.3	Cyclical Redundancy Checks	4·5	4.9.3	Wireless Communications	4·16
4.3.4	Self-Checking Codes	4·5			
4.4	HARDWARE OPERATIONS	4·6	4.10	CRYPTOGRAPHY	4·16
4.5	INTERRUPTS	4·7	4.11	BACKUP	4·17
4.5.1	Types of Interrupts	4·7	4.11.1	Personnel	4·18
4.5.2	Trapping	4·8	4.11.2	Hardware	4·18
4.6	MEMORY AND DATA STORAGE	4·8	4.11.3	Power	4·19
4.6.1	Main Memory	4·8	4.11.4	Testing	4·20
4.6.2	Read-Only Memory	4·8			
4.6.3	Secondary Storage	4·9			
4.7	TIME	4·10	4.12	RECOVERY PROCEDURES	4·20
4.7.1	Synchronous	4·10	4.13	MICROCOMPUTER CONSIDERATIONS	4·20
4.7.2	Asynchronous	4·11	4.13.1	Accessibility	4·20
4.8	NATURAL DANGERS	4·11	4.13.2	Knowledge	4·21
4.8.1	Power Failure	4·11	4.13.3	Motivation	4·21
4.8.2	Heat	4·11	4.13.4	Opportunity	4·21
4.8.3	Humidity	4·11	4.13.5	Threats to Microcomputers	4·21
4.8.4	Water	4·12	4.13.6	Maintenance and Repair	4·24
			4.14	CONCLUSION	4·25
			4.15	HARDWARE SECURITY CHECKLIST	4·25
			4.16	FURTHER READING	4·27

4 · 2 HARDWARE ELEMENTS OF SECURITY

4.1 INTRODUCTION. Computer hardware has always played a major role in computer security. Over the years, that role has increased dramatically, due to both the increases in processing power, storage capacity, and communications capabilities as well as the decreases in cost and size of components. The ubiquity of cheap, powerful, highly connected computing devices poses significant challenges to computer security. At the same time, the challenges posed by large, centralized computing systems have not diminished. An understanding of the hardware elements of computing is thus essential to a well-rounded understanding of computer security.

Chapter 1 of this *Handbook* has additional history of the evolution of information technology.

4.2 BINARY DESIGN. Although there are wide variations among computer architectures and hardware designs, all have at least one thing in common: They utilize a uniquely coded series of electrical impulses to represent any character within their range. Like the Morse code with its dots and dashes, computer pulse codes may be linked together to convey alphabetic or numeric information. Unlike the Morse code, however, computer pulse trains may also be combined in mathematical operations or data manipulation.

In 1946, Dr. John von Neumann, at the Institute for Advanced Study of Princeton University, first described in a formal report how the binary system of numbers could be used in computer implementations. The binary system requires combinations of only two numbers, 0 and 1, to represent any digit, letter, or symbol and, by extension, any group of digits, letters, or symbols. In contrast, the conventional decimal system requires combinations of 10 different numbers, from 0 to 9, letters from a to z, and a large number of symbols, to convey the same information. Von Neumann recognized that electrical and electronic elements could be considered as having only two states, on and off, and that these two states could be made to correspond to the 0 and 1 of the binary system. If the turning on and off of a computer element occurred at a rapid rate, the voltage or current outputs that resulted would best be described as pulses. Despite 60 years of intensive innovation in computer hardware, and the introduction of some optically based methods of data representation, the nature of these electrical pulses and the method of handling them remain the ultimate measure of a computer's accuracy and reliability.

4.2.1 Pulse Characteristics. Ideally, the waveform of a single pulse should be straight-sided, flat-topped, and of an exactly specified duration, amplitude, and phase relationship to other pulses in a series. It is the special virtue of digital computers that they can be designed to function at their original accuracy despite appreciable degradation of the pulse characteristics. However, errors will occur when certain limits are exceeded, and thus data integrity will be compromised. Because these errors are difficult to detect, it is important that a schedule of preventive maintenance be established and rigidly adhered to. Only in this way can operators detect degraded performance before it is severe enough to affect reliability.

4.2.2 Circuitry. To generate pulses of desirable characteristics, and to manipulate them correctly, requires components of uniform quality and dependability. To lower manufacturing costs, to make servicing and replacement easier, and generally to improve reliability, computer designers try to use as few different types of components as possible and to incorporate large numbers of each type into any one machine.

BINARY DESIGN 4 · 3

First-generation computers used as many as 30,000 vacuum tubes, mainly in a half-dozen types of logic elements. The basic circuits were flip-flops, or gates, that produced an output pulse whenever a given set of input pulses was present. However, vacuum tubes generated intense heat, even when in a standby condition. As a consequence, the useful operating time between failures was relatively short.

With the development of solid state diodes and transistors, computers became much smaller and very much cooler than their vacuum-tube predecessors. With advances in logic design, a single type of gate, such as the not-and (NAND) circuit, could replace all other logic elements. The resulting improvements in cost and reliability have been accelerated by the use of monolithic integrated circuits. Not least in importance is their vastly increased speed of operation. Since the meantime between failures of electronic computer circuitry is generally independent of the number of operations performed, it follows that throughput increases directly with speed; speed is defined as the rate at which a computer accesses, moves, and manipulates data. The ultimate limitation on computer speed is the time required for a signal to move from one physical element to another. At a velocity of 299,792,458 meters per second (186,282 miles per second) in vacuum, an electrical signal travels 3.0 meters or 9.84 feet in 10 nanoseconds (0.000.000.01 seconds). If components were as large as they were originally, and consequently as far apart, today's nanosecond computer speeds would clearly be impossible, as would be the increased throughput and reliability now commonplace.

4.2.3 Coding. In a typical application, data may be translated and retranslated automatically into a half-dozen different codes thousands of times each second. Many of these codes represent earlier technology retained for backward compatibility and economic reasons only. In any given code, each character appears as a specific group of pulses. Within each group, each pulse position is known as a *bit*, since it represents either of the two *binary digits*, 0 or 1. Exhibit 4.1 illustrates some of the translations that may be continuously performed as data move about within a single computer.

A *byte* is the name originally applied to the smallest group of bits that could be read and written (accessed or addressed) as a unit. Today a byte is always considered by convention to have 8 bits. In modern systems, a byte is viewed as the storage unit for a single American Standard Code for Information Interchange (ASCII) character, although newer systems such as Unicode, which handle international accented characters

EXHIBIT 4.1 Common Codes for Numeral 5

Code	Bits	Typical Use	Bit Pattern for "5"
Hexadecimal	4	Console switches	0101
Baudot	5	Paper tape	00001
Binary-Coded-Decimal (BCD)	6	Console indicators	000101
Transcode	6	Data transmission	110101
USASCII	7	Data transmission	0110101
EBCDIC	8	Buffer	11110101
EBCDIC, zoned decimal	8	Main memory	11000101
EBCDIC, packed decimal	8	Arithmetic logic unit	01011100
USASCII-8	8	Data transmission	01010101
Hollerith	12	Card reader/punch	0000000010000
Binary, halfword	16	Arithmetic logic unit	0000000000000101

4 · 4 HARDWARE ELEMENTS OF SECURITY

and many other symbols, use up to 4 bytes per character. By convention, most people use metric prefixes (kilo-, mega-, giga-, tera-) to indicate collections of bytes; thus *KB* refers to *kilobytes* and is usually defined as 1,024 bytes. Outside the data processing field, *K* would normally indicate the multiplier 1,000. Because of the ambiguity in definitions, the United States National Institute of Standards and Technology (NIST) proposed, and in 2000 the International Electrotechnical Commission established, a new set of units for information or computer storage. These units are established by a series of prefixes to indicate powers of 2; in this scheme, KB means *kibibytes* and refers exclusively to 1,024 (2^{10} or $\sim 10^3$) bytes. However, *kibibytes*, *mebibytes* (2^{20} or $\sim 10^6$), *gibibytes* (2^{30} or $\sim 10^9$), and *tebibytes* (2^{40} or $\sim 10^{12}$) are terms that have not yet become widely used.

Because translations between coding systems are accomplished at little apparent cost, any real incentive to unify the different systems is lacking. However, all data movements and translations increase the likelihood of internal error, and for this reason parity checks and validity tests have become indispensable.

4.3 PARITY. Redundancy is central to error-free data processing. By including extra bits in predetermined locations, certain types of errors can be detected immediately by inspection of these *metadata* (data about the original data). In a typical application, data move back and forth many times, among primary memory, secondary storage, input and output devices, as well as through communications links. During these moves, the data may lose integrity by dropping 1 or more bits, by having extraneous bits introduced, and by random changes in specific bits. To detect some of these occurrences, parity bits are added before data are moved and are checked afterward.

4.3.1 Vertical Redundancy Checks. In this relatively simple and inexpensive scheme, a determination is initially made as to whether there should be an odd or an even number of “1” bits in each character. For example, using the binary-coded decimal representation of the numerical “5,” we find that the 6-bit pulse group 000101 contains two 1s, an even number. Adding a seventh position to the code group, we may have either type of parity. If odd parity has been selected, a 1 would be added in the leftmost checkbit position:

Odd parity	1000101 three 1s
Even parity	0000101 two 1s

After any movement the number of 1 bits would be counted, and if not an odd number, an error would be assumed and processing halted. Of course, if 2 bits, or any even number, had been improperly transmitted, no error would be indicated since the number of “1” bits would still be odd.

To compound the problem of nonuniformity illustrated in Exhibit 4.1, each of the 4-, 5-, 6-, 7-, 8-, and 16-bit code groups may have an additional bit added for parity checking. Furthermore, there may be inconsistency of odd or even parity between manufacturers, or even between different equipment from a single supplier.

4.3.2 Longitudinal Redundancy Checks. Errors may not be detected by a vertical redundancy check (VRC) alone, for reasons just discussed. An additional safeguard, of particular use in data transmission and media recording such as tapes and disks, is the longitudinal redundancy check (LRC). With this technique, an extra

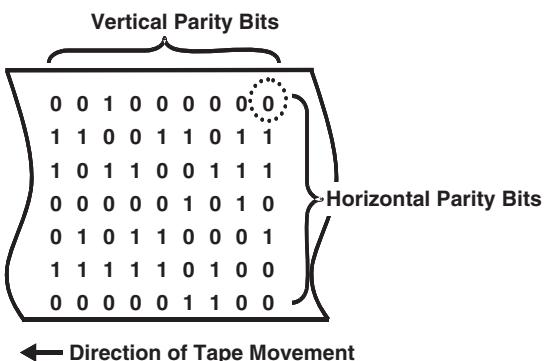
PARITY 4 · 5

EXHIBIT 4.2 Vertical and Longitudinal Parity,
Seven-Track Magnetic Tape

character is generated after some predetermined number of data characters. Each bit in the extra character provides parity for its corresponding row, just as the vertical parity bits do for their corresponding columns. Exhibit 4.2 represents both types as they would be recorded on 7-track magnetic tape. One bit has been circled to show that it is ambiguous. This bit appears at the intersection of the parity row and the parity column, and must be predetermined to be correct for one or the other, as it may not be correct for both. In the illustration, the ambiguous bit is correct for the odd parity requirement of the VRC character column; it is incorrect for the even parity of the LRC bit row.

In actual practice, the vertical checkbits would be attached to each character column as shown, but the longitudinal bits would follow a block of data that might contain 80 to several hundred characters. Where it is possible to use both LRC and VRC, any single data error in a block will be located at the intersection of incorrect row and column parity bits. The indicated bit may then be corrected automatically. The limitations of this method are: (1) multiple errors cannot be corrected, (2) an error in the ambiguous position cannot be corrected, and (3) an error that does not produce both a VRC and LRC indication cannot be corrected.

4.3.3 Cyclical Redundancy Checks. Where the cost of a data error could be high, the added expense of cyclical redundancy checks (CRCs) is warranted. In this technique, a relatively large number of redundant bits is used. For example, each 4-bit character requires 3 parity bits, while a 32-bit computer word needs 6 parity bits. Extra storage space is required in main and secondary memory, and transmissions take longer than without such checks. The advantage, however, is that any single error can be detected, whether in a data bit or a parity bit, and its location can be positively identified. By a simple electronic process of complementation, an incorrect 0 is converted to a 1, and vice versa.

4.3.4 Self-Checking Codes. Several types of codes are in use that inherently contain a checking ability similar to that of the parity system. Typical of these is the 2-of-5 code, in which every decimal digit is represented by a bit configuration containing exactly two 1s and three 0s. Where a parity test would consist of counting 1s to see if their number was odd or even, a 2-of-5 test would indicate an error whenever the number of 1s was more or less than 2.

4 · 6 HARDWARE ELEMENTS OF SECURITY

4.4 HARDWARE OPERATIONS. Input, output, and processing are the three essential functions of any computer. To protect data integrity during these operations, several hardware features are available.

- **Read-after-write.** In disk and tape drives, it is common practice to read the data immediately after they are recorded and to compare them with the original values. Any disagreement signals an error that requires rewriting.
- **Echo.** Data transmitted to a peripheral device, to a remote terminal (see Section 4.9.1), or to another computer can be made to generate a return signal. This echo is compared with the original signal to verify correct reception. However, there is always the risk that an error will occur in the return signal and falsely indicate an error in the original transmission.
- **Overflow.** The maximum range of numerical values that any computer can accommodate is fixed by its design. If a program is improperly scaled, or if an impossible operation such as dividing by zero is called for, the result of an arithmetic operation may exceed the allowable range, producing an overflow error. Earlier computers required programmed instructions to detect overflows, but this function now generally is performed by hardware elements at the machine level. Overflows within application programs still must be dealt with in software. (Indeed, failure to do so can render software susceptible to abuse by malicious parties.)
- **Validation.** In any one computer coding system, some bit patterns may be unassigned, and others may be illegal. In the IBM System/360 Extended Binary Coded Decimal Interchange Code (EBCDIC), for example, the number 9 is represented by 11111001, but 11111010 is unassigned. A parity check would not detect the second group as being in error, since both have the same number of 1 bits. A validity check, however, would reject the second bit configuration as invalid.

Similarly, certain bit patterns represent assigned instruction codes while others do not. In one computer, the instruction to convert packed-decimal numbers to zoned-decimal numbers is 11110011, or F3 in hexadecimal notation; 11110101, or F5, is unassigned, and a validity check would cause a processing halt whenever that instruction was tested.

- **Replication.** In highly sensitive applications, it is good practice to provide backup equipment on site, for immediate changeover in the event of failure of the primary computer. For this reason, it is sometimes prudent to retain two identical, smaller computers rather than to replace them with a single unit of equivalent or even greater power. Fault-tolerant, or fail-safe, computers use two or more processors that operate simultaneously, sharing the load and exchanging information about the current status of duplicate processes running in parallel. If one of the processors fails, another continues the processing without pause.

Many sensitive applications, such as airline reservation systems, have extensive data communications facilities. It is important that all of this equipment be duplicated as well as the computers themselves. (The failure of an airline reservation system, if permitted to extend beyond a relatively small number of hours, could lead to failure of the airline itself.)

Replacements should also be immediately available for peripheral devices. In some operating systems, it is necessary to inform the system that a device is down and to reassign its functions to another unit. In the more sophisticated systems, a malfunctioning device is automatically cut out and replaced. For example, the New York Stock Exchange operates and maintains two identical trading systems so that failure of the primary system should not result in any interruption to trading.

INTERRUPTS 4 · 7

4.5 INTERRUPTS. The sequence of operations performed by a computer system is determined by a group of instructions: a program. However, many events that occur during operations require deviations from the programmed sequence. *Interrupts* are signals generated by hardware elements that detect changed conditions and initiate appropriate action. The first step is immediately to store the status of various elements in preassigned memory locations. The particular stored bit patterns, commonly called program status words, contain the information necessary for the computer to identify the cause of the interrupt, to take action to process it, and then to return to the proper instruction in the program sequence after the interrupt is cleared.

4.5.1 Types of Interrupts. Five types of interrupts are in general use. Each of them is of importance in establishing and maintaining data processing integrity.

4.5.1.1 Input/Output Interrupts. Input/output (I/O) interrupts are generated whenever a device or channel that had been busy becomes available. This capability is necessary to achieve error-free use of the increased throughput provided by buffering, overlapped processing, and multiprogramming.

After each I/O interrupt, a check is made to determine whether the data have been read or written without error. If so, the next I/O operation can be started; if not, an error-recovery procedure is initiated. The number of times that errors occur should be recorded so that degraded performance can be detected and corrected.

4.5.1.2 Supervisor Calls. The supervisor, or monitor, is a part of the operating system software that controls the interactions between all hardware and software elements.

Every request to read or write data is scheduled by the supervisor when called upon to do so. I/O interrupts also are handled by supervisor calls that coordinate them with read/write requests. Loading, executing, and terminating programs are other important functions initiated by supervisor calls.

4.5.1.3 Program Check Interrupts. Improper use of instructions or data may cause an interrupt that terminates the program. For example, attempts to divide by zero and operations resulting in arithmetic overflow are voided. Unassigned instruction codes, attempts to access protected storage, and invalid data addresses are other types of exceptions that cause program check interrupts.

4.5.1.4 Machine Check Interrupts. Among the exception conditions that will cause machine check interrupts are parity errors, bad disk sectors, disconnection of peripherals in use, and defective circuit modules. It is important that proper procedures be followed to clear machine checks without loss of data or processing error.

4.5.1.5 External Interrupts. External interrupts are generated by timer action, by pressing an Interrupt key, or by signals from another computer. When two central processing units are interconnected, signals that pass between them initiate external interrupts. In this way, control and synchronization are continuously maintained while programs, data, and peripheral devices may be shared and coordinated.

In mainframes, an electronic clock generally is included in the central processor unit for time-of-day entries in job logs and for elapsed-time measurements. As an interval timer, the clock can be set to generate an interrupt after a given period. This feature should be used as a security measure, preventing sensitive jobs from remaining on the computer long enough to permit unauthorized manipulation of data or instructions.

4 · 8 HARDWARE ELEMENTS OF SECURITY

4.5.2 Trapping. Trapping is a type of hardware response to an interrupt. Upon detecting the exception, an unconditional branch is taken to some predetermined location. An instruction there transfers control to a supervisor routine that initiates appropriate action.

4.6 MEMORY AND DATA STORAGE. Just as the human mind is subject to aberrations, so is computer memory. In the interests of data security and integrity, various therapeutic measures have been developed for the several types of storage.

4.6.1 Main Memory. Random access memory (RAM), and its derivatives, such as dynamic RAM (DRAM), synchronous DRAM (SDRAM, introduced in 1996 and running at 133 megaHertz [MHz]), and DDR-3 (Double Data Rate 3 SDRAM, announced in 2005 and running at 800 MHz), share the necessary quality of being easily and quickly accessed for reading and writing of data. Unfortunately, this necessary characteristic is at the same time a potential source of difficulty in maintaining data integrity against unwanted read/write operations. The problems are greatly intensified in a multiprogramming environment, especially with dynamic memory allocation, where the possibility exists that one program will write improperly over another's data in memory. Protection against this type of error must be provided by the operating system. Chapter 24 in this *Handbook* discusses operating system security in more detail.

One form of protection requires that main memory be divided into blocks or *pages*; for example, 2,048 eight-bit bytes each. Pages can be designated as read-only when they contain constants, tables, or programs to be shared by several users. Additionally, pages that are to be inaccessible except to designated users may be assigned a lock by appropriate program instructions. If a matching key is not included in the user's program, access to that page will be denied. Protection may be afforded against writing only or against reading and writing.

4.6.2 Read-Only Memory. One distinguishing feature of main memory is the extremely high speed at which data can be entered or read out. The set of sequential procedures that accomplishes this and other functions is the program, and the programmer has complete freedom to combine any valid instructions in a meaningful way. However, certain operations, such as system start-up, or *booting*, are frequently and routinely required, and they may be performed automatically by a preprogrammed group of memory elements. These elements should be protected from inadvertent or unauthorized changes.

For this purpose, a class of memory elements has been developed that, once programmed, cannot be changed at all, or require a relatively long time to do so. These elements are called *read-only memory*, or ROM; the process by which sequential instructions are set into these elements is known as *microprogramming*. The technique can be used to advantage where data integrity is safeguarded by eliminating the possibility of programmer error.

Variations of the principle include programmable read-only memories (PROM) and erasable, programmable read-only memory (EPROM), all of which combine microprogramming with a somewhat greater degree of flexibility than read-only memory itself. The data on these chips can be changed through a special operation often referred to as *flashing* (literally exposure to strong ultraviolet light; this is different from *flash memory* used today for storage of such data as digital music files and digital photographs—we will return to the subject of *flash memory* in the next section).

MEMORY AND DATA STORAGE 4 · 9

4.6.3 Secondary Storage. The term “secondary storage” traditionally has been used to describe storage such as magnetic disks, diskettes, tapes, and tape cartridges. Although the 1.44 megabyte (MB) magnetic floppy disk is obsolete, the magnetic hard drive, with capacities up to terabytes, remains an essential element of virtually all computers, and terabyte-capacity external hard drives the size of a paperback book are now available off-the-shelf for a few hundred dollars.

A more recent development are optical drives such as the removable, *compact disc read-only memory* (CD-ROM), originally made available in the early 1980s, which are useful for long-term archival storage of around 700 MB per disc. Hybrid forms of this type exist as well, such as CD-Rs, which can be written to once, and CD-RWs, which accommodate multiple reads and writes. The digital video disc (DVD), or as it has been renamed, the digital versatile disc, was introduced in 1997 and provides capacities ranging from 4.7 gigabytes (GB) per disc up to 30 GB for data archiving. The higher-capacity optical discs use Blu-ray technology introduced in 2002 and can store 25 GB per side; they typically are used for distributing movies, but BD-R (single use) and BD-RE (rewritable) discs hold much potential for generalized data storage.

The newest addition to secondary storage is RAM that simulates hard disks, known as *flash memory*. Derived from electrical EPROMs (EEPROMs) and introduced by Toshiba in the 1980s, this kind of memory now exists in a huge variety of formats, including relatively inexpensive Universal Serial Bus (USB) tokens with storage capacities now in the gigabyte range. These devices appear as external disk drives when plugged into a plug-and-play personal computer. Another flash memory format is small cards, many the size of postage stamps, that can be inserted into mobile phones, cameras, printers, and other devices as well as computers.

Hardware safeguards described earlier, such as redundancy, validity, parity, and read-after-write, are of value in preserving the integrity of secondary storage. These safeguards are built into the equipment and are always operational unless disabled or malfunctioning. Other precautionary measures are optional, such as standard internal labeling procedures for drives, tapes, and disks. Standard internal labels can include identification numbers, record counts, and dates of creation and expiration. Although helpful, external plastic or paper labels on recordable media are not an adequate substitute for computer-generated labels, magnetically inscribed on the medium itself and automatically checked by programmed instructions.

Another security measure sometimes subverted is write-protection on removable media. Hardware interlocks prevent writing to them. These locks should be activated immediately when the media are removed from the system. Failure to do so will cause the data to be destroyed if the same media are improperly used on another occasion.

Hard drives, optical discs, and flash memory cards are classified as direct access storage devices (DASDs). Unlike magnetic tapes with their exclusively sequential processing, DASDs may process data randomly as well as in sequence. This capability is essential to online operations, where it is not possible to sort transactions prior to processing. The disadvantage of direct access is that there may be less control over entries and more opportunity to degrade the system than exists with sequential batch processing.

One possible source of DASD errors arises from the high rotational velocity of the recording medium and, except on head-per-track devices, the movement of heads as well. To minimize this possibility, areas on the recording surface have their addresses magnetically inscribed. When the computer directs that data be read from or into a particular location, the address in main memory is compared with that read from the DASD. Only if there is agreement will the operation be performed.

4 · 10 HARDWARE ELEMENTS OF SECURITY

Through proper programming, the integrity of data can be further assured. In addition to the address check, comparisons can be made on identification numbers or on key fields within each record. Although the additional processing time is generally negligible, there can be a substantial improvement in properly posting transactions.

Several other security measures often are incorporated into DASDs. One is similar to the protection feature in main memory and relies on determining “extents” for each data set. If these extents, which are simply the upper and lower limits of a data file’s addresses, are exceeded, the job will terminate.

Another safety measure is necessitated by the fact that defective areas on a disk surface may cause errors undetectable in normal operations. To minimize this possibility, disks should be tested and certified prior to use and periodically thereafter. Further information is provided by operating systems that record the number of disk errors encountered. Reformatting or replacement must be ordered when errors exceed a predetermined level. Many personal computer hard drives now have some form of Self-Monitoring, Analysis, and Reporting Technology (SMART). Evolved from earlier technology such as IBM’s Predictive Failure Analysis (PFA) and Intellisafe by computer manufacturer Compaq, and disk drive manufacturers Seagate, Quantum, and Conner, SMART can alert operators to potential drive problems. Unfortunately, the implementation of SMART is not standardized, and its potential for preventive maintenance and failure prediction is often overlooked.

Note that SMART is different from the range of technologies used to protect hard drives from head crashes. A head crash occurs when the component that reads data from the disk actually touches the surface of the disk, potentially damaging it and the data stored on it. Many hard drives have systems in place to withdraw heads from the disk before such contact occurs. These protective measures have reached the point where an active hard drive can be carried around in relative safety as part of a music and video player (e.g., Apple iPod or Microsoft Zune).

4.7 TIME. Within the computer room and in many offices, a wall clock is usually a dominant feature. There is no doubt that this real-time indicator is of importance in scheduling and regulating the functions of people and machines, but the computer’s internal timings are more important for security.

4.7.1 Synchronous. Many computer operations are independent of the time of day but must maintain accurate relationships with some common time and with each other. Examples of this synchronism include the operation of gates, flip-flops, and registers, and the transmission of data at high speeds. Synchronism is obtained in various ways. For gates and other circuit elements, electronic clocks provide accurately spaced pulses at a high-frequency rate, while disk and tape drives are maintained at rated speed by servomotor controls based on power-line frequency.

Of all computer errors, the ones most difficult to detect and correct are probably those caused by timing inconsistencies. Internal clocks may produce 1 billion pulses per second (known as 1 gigahertz [GHz]), or more, when the computer is turned on. The loss of even a single pulse, or its random deformation or delay, can cause undetected errors. More troublesome is the fact that even if errors are detected, their cause may not be identified unless the random timing faults become frequent or consistent.

An example of the insidious nature of timing faults is the consequence of electrical power fluctuations when voltage drops below standard. During these power transients, disk drives may slow down; if sectors are being recorded, their physical size will be

NATURAL DANGERS 4 · 11

correspondingly smaller. Then, when the proper voltage returns, the incorrect sector sizes can cause data errors or loss.

4.7.2 Asynchronous. Some operations do not occur at fixed time intervals and therefore are termed “asynchronous.” In this mode, signals generated by the end of one action initiate the following one. As an example, low-speed data transmissions such as those using ordinary modems are usually asynchronous. Coded signals produced by the random depression of keyboard keys are independent of any clock pulses.

4.8 NATURAL DANGERS. To preserve the accuracy and timeliness of computer output, computers must be protected against environmental dangers. Chapters 22 and 23 of this *Handbook* discuss such threats in extensive detail.

4.8.1 Power Failure. Probably the most frequent cause of computer downtime is power failure. Brownouts and blackouts are visible signs of trouble; undetected voltage spikes are far more common, although hardly less damaging.

Lightning can produce voltage spikes on communications and power lines of sufficient amplitude to destroy equipment or, at the very least, to alter data randomly. Sudden storms and intense heat or cold place excessive loads on generators. The drop in line voltage that results can cause computer or peripheral malfunction. Even if it does not, harmful voltage spikes may be created whenever additional generators are switched in to carry higher loads.

Where warranted, a recording indicator may be used to detect power-line fluctuations. Such monitoring often is recommended when computer systems show unexplained, erratic errors. At any time that out-of-tolerance conditions are signaled, the computer output should be checked carefully to ensure that data integrity has not been compromised. If such occurrences are frequent, or if the application is a sensitive one, auxiliary power management equipment should be considered. These range from simple voltage regulators and line conditioners to uninterruptible power supplies (UPSs).

4.8.2 Heat. Sustained high temperatures can cause electronic components to malfunction or to fail completely. Air conditioning (AC) is therefore essential, and all units must be adequate, reliable, and properly installed. If backup electrical power is provided for the computer, it must also be available for the air conditioners. For example, after the San Francisco earthquake of 1989, the desktop computers and network servers in at least one corporate headquarters were damaged by a lack of synchronization between air conditioning and power supply. The AC was knocked out by the quake, and the building was evacuated, but the computers were left on. Many of them failed at the chip and motherboard level over the next few days because the temperature in the uncooled offices got too high. A frequently unrecognized cause of overheating is obstruction of ventilating grilles. Printouts, tapes, books, and other objects must not be placed on cabinets where they can prevent free air circulation. A digital thermometer is a good investment for any room in which computers are used. Today, many electronic devices include thermostats that cut off the power if internal temperatures exceed a danger limit.

4.8.3 Humidity. Either extreme of humidity can be damaging. Low humidity—below about 20 percent—permits buildup of static electricity charges that

4.12 HARDWARE ELEMENTS OF SECURITY

may affect data pulses. Because this phenomenon is intensified by carpeting, computer room floors should either be free of carpeting or covered with antistatic carpet.

High humidity—above about 80 percent—may lead to condensation that causes shorts in electrical circuits or corrodes metal contacts. To ensure operation within acceptable limits, humidity controls should be installed and a continuous record kept of measured values.

4.8.4 Water. Water introduced by rain, floods, bursting pipes, and overhead sprinklers probably has been responsible for more actual computer damage than fire or any other single factor. Care taken in locating computer facilities, in routing water pipes, and in the selection of fire-extinguishing agents will minimize this significant danger.

The unavailability of water—following a main break, for example—will cause almost immediate shutdown of water-cooled mainframes. Mission-critical data centers should be prepared for this contingency. As an example, when the Des Moines River flooded in 1993, it caused the skyscraper housing the headquarters of the Principal Financial Group to be evacuated, but not because of water in the building. The building stayed high and dry, but the flood forced the city water plant to shut down, depriving the building of the water necessary for cooling. After the flood, the company installed a 40,000-gallon water tank in the basement, to prevent any recurrence of this problem.

4.8.5 Dirt and Dust. Particles of foreign matter can interfere with proper operation of magnetic tape and disk drives, printers, and other electromechanical devices. All air intakes must be filtered, and all filters must be kept clean. Cups of coffee seem to become especially unstable in a computer environment; together with any other food or drink, they should be banned entirely. Throughout all areas where computer equipment is in use, good housekeeping principles should be rigorously enforced.

4.8.6 Radiation. Much has been written about the destructive effect of magnetic fields on tape or disk files. However, because magnetic field strength diminishes rapidly with distance, it is unlikely that damage actually could be caused except by large magnets held very close to the recorded surfaces. For example, storing a CD or DVD by attaching it to a filing cabinet with a magnet is not a good idea, but simply walking past a refrigerator decorated with magnets while holding a CD or DVD is unlikely to do any damage.

The proliferation of wireless signals can expose data to erroneous pulses. Offices should be alert for possible interference from and between cordless phones, mobile phones, wireless Internet access points and peripherals, and microwave ovens.

Radioactivity may be a great threat to personnel but not to the computer or its recording media.

4.8.7 Downtime. It is essential to the proper functioning of a data center that preventive maintenance be performed regularly and that accurate records be kept of the time and the reason that any element of the computer is inoperative. The more often the computer is down, the more rushed operators will be to catch up on their scheduled workloads. Under such conditions, controls are bypassed, shortcuts are taken, and human errors multiply.

Downtime records should be studied to detect unfavorable trends and to pinpoint equipment that must be overhauled or replaced before outages become excessive.

DATA COMMUNICATIONS 4 · 13

If unscheduled downtime increases, preventive maintenance should be expanded or improved until the trend is reversed.

4.9 DATA COMMUNICATIONS. One of the most dynamic factors in current computer usage is the proliferation of devices and systems for data transmission. These range from telephone modems to wired networks, from Internet-enabled cell phones to 802.11 wireless Ethernet, and include Bluetooth, infrared, personal digital assistants (PDAs), music players, and new technologies that appear almost monthly. Computers that do not function at least part time in a connected mode may well be rarities. For fundamentals of data communications, see Chapter 5 of this *Handbook*.

The necessity for speeding information over great distances increases in proportion to the size and geographic dispersion of economic entities; the necessity for maintaining data integrity and security, and the difficulty of doing so, increases even more rapidly. Major threats to be guarded against include human and machine errors, unauthorized accession, alteration, and sabotage. The term “accession” refers to an ability to read data stored or transmitted within a computer system; it may be accidental or purposeful. “Alteration” is the willful entering of unauthorized or incorrect data. “Sabotage” is the intentional act of destroying or damaging the system or the data within it. For each of these threats, the exposure and the countermeasures will depend on the equipment and the facilities involved.

4.9.1 Terminals. In these discussions, a *terminal* is any input/output device that includes facilities for receiving, displaying, composing, and sending data. Examples include personal computers and specialized devices such as credit card validation units.

Data communications are carried on between computers, between terminals, or between computers and terminals. The terminals themselves may be classified as *dumb* or *intelligent*. Dumb terminals have little or no processing or storage capability and are largely dependent on a host computer for those functions. Intelligent terminals generally include disk storage and capabilities roughly equivalent to those of a personal computer. In addition to vastly improved communications capabilities, they are capable of stand-alone operation.

In the simplest of terminals, the only protection against transmission errors lies in the inability to recognize characters not included in the valid set and to display a question mark or other symbol when one occurs. Almost any terminal can be equipped to detect a vertical parity error. More sophisticated terminals are capable of detecting additional errors through longitudinal and cyclical redundancy characters, as well as by vertical parity and validity checks. Of course, error detection is only the first step in maintaining data integrity. Error correction is by far the more important part, and retransmission is the most widely used correction technique.

Intelligent terminals and personal computers are capable of high-speed transmission and reception. They can perform complicated tests on data before requesting retransmission, or they may even be programmed to correct errors internally. The techniques for self-correction require forward-acting codes, such as the Hamming cyclical code. These are similar to the error-detecting cyclic redundancy codes, except that they require even more redundant bits. Although error correction is more expensive and usually slower than detection with retransmission, it is useful under certain circumstances. Examples include simplex circuits where no return signal is possible, and half-duplex circuits where the time to turn the line around from transmission to reception is too long. Forward correction is also necessary where errors are so

4 · 14 HARDWARE ELEMENTS OF SECURITY

numerous that retransmissions would clog the circuits, with little or no useful information throughput.

A more effective use of intelligent terminals and personal computers is to preserve data integrity by encryption, as described in this chapter and in Chapter 7. Also, they may be used for compression or compaction. Reducing the number of characters in a message reduces the probability of an error as well as the time required for transmission. One technique replaces long strings of spaces or zeroes with a special character and a numerical count; the procedure is reversed when receiving data.

Finally, the intelligent terminal or microprocessor may be used to encode or decipher data when the level of security warrants cryptography.

All terminals, of every type, including desktop and notebook personal computers (PCs), have at least one thing in common: the need to be protected against sabotage or unauthorized use. Although the principles for determining proper physical location and the procedures for restricting access are essentially the same as those that apply to a central computer facility, the actual problems of remote terminals are even more difficult. Isolated locations, inadequate supervision, and easier access by more people all increase the likelihood of compromised security.

4.9.2 Wired Facilities. Four types of wired facilities are in widespread use: dial-up access, leased lines, digital subscriber lines (DSL), and cable carriers. Both common carriers and independent systems may employ various media for data transmission. The increasing need for higher speed and better quality in data transmission has prompted utilization of coaxial and fiber optic cables, while microwave stations and communication satellites often are found as wireless links within wired systems.

Generally, decisions as to the choice of service are based on the volume of data to be handled and on the associated costs, but security considerations may be even more important.

4.9.2.1 Dial-Up Lines. Still widely used for credit and debit card terminals, dial-up lines have been replaced for many other applications by leased lines, DSL lines, and cables carrying Internet traffic (using the TCP/IP protocol discussed in Chapter 5 of this *Handbook*). Dial-up connections are established between modems operating over regular voice lines sometimes referred to as *plain old telephone service* (POTS).

Where dial-up access to hardware still exists, for example, for maintenance of certain equipment, proper controls are essential to protect both the equipment and the integrity of other systems to which it might be connected. Dial-up ports may be reached by anyone with a phone, anywhere on the planet, and the practice of *war-dialing* to detect modems is still used by those seeking unauthorized access to an organization's network. (War dialing involves dialing blocks of numbers to find which ones respond as modems or fax machines. These numbers are recorded and may be dialed later in an attempt to gain unauthorized access to systems or services.) It is advisable to:

- Compile a log of unauthorized attempts at entry, and use it to discourage further efforts.
- Compile a log of all accesses to sensitive data, and verify their appropriateness.
- Equip all terminals with internal identification generators or answer-back units, so that even a proper password would be rejected if sent from an unauthorized

DATA COMMUNICATIONS 4 · 15

terminal. This technique may require the availability of an authorized backup terminal in the event of malfunction of the primary unit.

- Provide users with personal identification in addition to a password if the level of security requires it. The additional safeguard could be a magnetically striped or computerized plastic card to be inserted into a special reader. The value of such cards is limited, since they can be used by anyone, whether authorized or not. For high-security requirements, other hardware-dependent biometric identifiers, such as handprints and voiceprints, should be considered.
- Where appropriate, utilize call-back equipment that prevents a remote station from entering a computer directly. Instead, the device dials the caller from an internal list of approved phone numbers to make the actual connection.

With proper password discipline, problems of accession, alteration, and data sabotage can be minimized. However, the quality of transmissions is highly variable. Built into the public telephone system is an automatic route-finding mechanism that directs signals through uncontrollable paths. The distance and the number of switching points traversed, and the chance presence of cross-talk, transients, and other noise products will have unpredictable effects on the incidence of errors. Parity systems, described earlier, are an effective means of reducing such errors.

4.9.2.2 Leased Lines. Lines leased from a common carrier for the exclusive use of one subscriber are known as *dedicated lines*. Because they are directly connected between predetermined points, normally they cannot be reached through the dial-up network. Traditionally, leased lines were copper, but point-to-point fiber optic and coaxial cable lines can also be leased.

Wiretapping is a technically feasible method of accessing leased lines, but it is more costly, more difficult, and less convenient than dialing through the switched network. Leased lines are generally more secure than those that can be readily war-dialed.

To this increased level of security for leased lines is added the assurance of higher-quality reception. The problems of uncertain transmission paths and switching transients are eliminated, although other error sources are not. In consequence, parity checking remains a minimum requirement.

4.9.2.3 Digital Subscriber Lines. Falling somewhere in between a leased line and POTS, a digital subscriber line offers digital transmission locally over ordinary phone lines that can be used simultaneously for voice transmission. This is possible because ordinary copper phone lines can carry, at least for short distances, signals that are in a much higher-frequency range than the human voice. A DSL modem is used by a computer to reach the nearest telephone company switch, at which point the data transmission enters the Internet backbone. Computers connected to the Internet over DSL communicate using TCP/IP and are said to be hosts rather than terminals. They are prone to compromise through a wide range of exploits. However, few if any of these threats are enabled by the DSL itself. As with leased lines, wiretapping is possible, but other attacks, such as exploiting weaknesses in TCP/IP implementations on host machines, are easier.

4.9.2.4 Cable Carriers. Wherever cable television (TV) is available, the same optical fiber or coaxial cables that carry the TV signal also can be used to provide high-speed data communications. The advantages of this technology include download

4 · 16 HARDWARE ELEMENTS OF SECURITY

speeds that can, in the case of coaxial cables, exceed 50 megabits per second, or in the case of fiber optic cable, exceed 100 gigabits per second.

The disadvantages arise from the fact that connections to the carrier may be shared by other subscribers in the same locality. Unless the service provider limits access, perhaps in accordance with a quality-of-service agreement, multiple subscribers can be online simultaneously and thus slow down transmission speeds. Even more serious is the possibility of security breaches, since multiple computers within a neighborhood may be sharing part of a virtual local area network, and thus each is potentially accessible to every other node on that network. For this reason alone, cable connections should be firewalled. For details of firewalls and their uses, see Chapter 26 in this *Handbook*. Another reason for using firewalls is that cable connections are always on, providing maximal opportunity for hackers to access an unattended computer.

4.9.3 Wireless Communications. Data transfers among multinational corporations have been growing very rapidly, and transoceanic radio and telephone lines have proved too costly, too slow, too crowded, and too error-prone to provide adequate service. An important alternative is the communications satellite. Orbiting above Earth, the satellite reflects ultra-high-frequency radio signals that can convey a television program or computer data with equal speed and facility.

For communications over shorter distances, the cost of common-carrier wired services has been so high as to encourage competitive technologies. One of these, the microwave radio link, is used in many networks. One characteristic of such transmissions is that they can be received only on a direct line-of-sight path from the transmitting or retransmitting antenna. With such point-to-point ground stations, it is sometimes difficult to position the radio beams where they cannot be intercepted; with satellite and wireless broadcast communications, it is impossible. This is a significant issue with wireless local area network technology based on the IEEE 802.11 standards and commonly known as Wi-Fi (a brand name owned by the Wi-Fi Alliance; the term is short for *wireless fidelity*). The need for security is consequently greater, and scramblers or cryptographic encoders are essential for sensitive data transfers.

Because of the wide bandwidths at microwave frequencies, extremely fast rates of data transfer are possible. With vertical, longitudinal, and cyclical redundancy check characters, almost all errors can be detected, yet throughput remains high.

4.10 CRYPTOGRAPHY. Competitive pressures in business, politics, and international affairs continually create situations where morality, privacy, and the laws all appear to give way before a compelling desire for gain. Information, for its own sake or for the price it brings, is an eagerly sought after commodity. We are accustomed to the sight of armored cars and armed guards transporting currency, yet often invaluable data are moved with few precautions. When the number of computers and competent technicians was small, the risk in careless handling of data resources was perhaps not great. Now, however, a very large population of knowledgeable computer people exists, and within it are individuals willing and able to use their knowledge for illegal ends. Others find stimulation and satisfaction in meeting the intellectual challenge that they perceive in defeating computer security measures.

Acquiring information in an unauthorized manner is relatively easy when data are communicated between locations. One method of discouraging this practice, or rendering it too expensive to be worth the effort, is cryptographic encoding of data prior to transmission. This technique is also useful in preserving the security of files within data storage devices. If all important files were stored on magnetic or optical media

BACKUP 4 · 17

in cryptographic cipher only, the incidence of theft and resale would unquestionably be less.

Many types of ciphers might be used, depending on their cost and the degree of security required. Theoretically, any code can be broken, given enough time and equipment. In practice, if a cipher cannot be broken fairly quickly, the encoded data are likely to become valueless. However, since the key itself can be used to decipher later messages, it is necessary that codes or keys be changed frequently.

For further information on cryptography, refer to Chapter 7 in this *Handbook*.

4.11 BACKUP. As with most problems, the principal focus in computer security ought to be on prevention rather than on cure. No matter how great the effort, however, complete success can never be guaranteed. There are four reasons for this being so:

1. Not every problem can be anticipated.
2. Where the cost of averting a particular loss exceeds that of recovery, preventive measures may not be justified.
3. Precautionary measures, carried to extremes, can place impossible constraints on the efficiency and productivity of an operation. It may be necessary, therefore, to avoid such measures aimed at events whose statistical probability of occurrence is small.
4. Even under optimum conditions, carefully laid plans may go astray. In the real world of uncertainty and human fallibility, where there is active or inadvertent interference, it is almost a certainty that at one time or another, the best of precautionary measures will prove to be ineffective.

Recognizing the impossibility of preventing all undesired actions and events, it becomes necessary to plan appropriate means of recovering from them. Such plans must include backup for personnel, hardware, power, physical facilities, data, and software. Data backups are discussed more fully in Chapter 57 of this *Handbook*.

Responding to emergencies is described in Chapters 56 of this *Handbook* and business continuity planning and disaster recovery are discussed in Chapter 58 and 59.

Backup plans should be evaluated with respect to:

- The priorities established for each application, to ensure that they are properly assigned and actually observed.
- The time required to restore high-priority applications to full-functioning status.
- The degree of assurance that plans actually can be carried out when required. For important applications, alternative plans should be available in the event that the primary plan cannot be implemented.
- The degree of security and data integrity that will exist if backup plans actually are put into effect.
- The extent to which changing internal or external conditions are noted, and the speed with which plans are modified to reflect such changes.

The assignment of priorities in advance of an actual emergency is an essential and critically important process. In most organizations, new applications proliferate, while old ones are rarely discarded. If backup plans attempt to encompass all jobs,

4 · 18 HARDWARE ELEMENTS OF SECURITY

they are likely to accomplish none. Proper utilization of priorities will permit realistic scheduling, with important jobs done on time and at acceptable costs.

4.11.1 Personnel. The problems of everyday computer operation require contingency plans for personnel on whose performance hardware functioning depends. Illnesses, vacations, dismissals, promotions, resignations, overtime, and extra shifts are some of the reasons why prudent managers are continuously concerned with the problem of personnel backup. The same practices that work for everyday problems can provide guidelines for emergency backup plans. This subject is covered more fully in Chapter 45 of this *Handbook*.

4.11.2 Hardware. Hardware backup for data centers can take several forms:

- Multiple processors at the same site to protect against loss of service due to breakdown of one unit
- Duplicate installations at nearby facilities of the same company
- Maintaining programs at a compatible service bureau, on a test or standby basis
- A contract for backup at a facility dedicated to disaster recovery
- A reciprocal agreement with a similar installation at another company

The probability of two onsite processors both being down at the same time due to internal faults is extremely small. Consequently, most multiple installations rarely fall behind on mission-critical applications. However, this type of backup offers no protection against power failure, fire, vandalism, or any disaster that could strike two or more processors at once. The disasters of September 11, 2001, proved that even a highly unlikely event actually could occur. With duplicate processors at different but commonly owned sites, there is little chance of both being affected by the same forces. Although the safety factor increases with the distance separating them, the difficulty of transporting people and data becomes greater. An alternate site must represent a compromise between these conflicting objectives. Furthermore, complete compatibility of hardware and software will have to be preserved, even though doing so places an undue operational burden on one of the installations. Shortly after September 11, a number of New York financial firms were back in operation with their alternative computer sites across the Hudson River.

The backup provided by service bureaus can be extremely effective, particularly if the choice of facility is carefully made. Although progressive service bureaus frequently improve both hardware and software, they almost never do so in a way that would cause compatibility problems for their existing customers. Once programs have been tested, they can be stored offline on tape or disk at little cost. Updated masters can be rotated in the service bureau library, providing offsite data backup as well as the ability to become fully operational at once.

Effective hardware backup is also available at independent facilities created expressly for that purpose. In one type of facility, there are adequate space, power, air conditioning, and communication lines to accommodate a very large system. Most manufacturers are able to provide almost any configuration on short notice when disaster strikes a valued customer. The costs for this type of base standby facility are shared by a number of users so that expenses are minimal until an actual need arises. However, if two or more sharers are geographically close, their facilities may be rendered inoperative by the same fire, flood, or power failure. Before contracting for such a facility,

BACKUP 4 · 19

it is necessary to analyze this potential problem; the alternative is likely to be a totally false sense of security. Several firms whose facilities were damaged or destroyed on September 11 were provided with complete replacement equipment by their vendors within a short time.

Another type of backup facility is already equipped with computers, disk and tape drives, printers, terminals, and communications lines so that it can substitute instantly for an inoperative system. The standby costs for this service are appreciably more than for a base facility, but the assurance of recovery in the shortest possible time is far greater. Here, too, it would be prudent to study the likelihood of more than one customer requiring the facility at the same time and to demand assurance that one's own needs will be met without fail. Several companies successfully availed themselves of this type of backup and disaster recovery after September 11.

Backup by reciprocal agreement was for many years an accepted practice, although not often put to the test. Unfortunately, many managers still rely on this outmoded safeguard. One has only to survive a single major change of operating system software to realize that when it occurs, neither the time nor the inclination is available to modify and test another company's programs. Even the minor changes in hardware and software that continuously take place in most installations could render them incompatible. At the same time, in accordance with Parkinson's Law, workloads always expand to fill the available time and facilities. In consequence, many who believe that they have adequate backup will get little more than an unpleasant surprise, should they try to avail themselves of the privilege.

4.11.3 Power. The one truly indispensable element of any data processing installation is electric power. Backing up power to PCs and small servers by uninterruptible power supplies is reasonable in cost and quite effective. For mainframes and large servers, several types of power backup are available. The principal determinant in selection should be the total cost of anticipated downtime and reruns versus the cost of backup to eliminate them. Downtime and rerun time may be extrapolated from records of past experience.

Problems due to electrical power may be classified by type and by the length of time that they persist. Power problems as they affect computers consist of variations in amplitude, frequency, and waveform, with durations ranging from fractions of a millisecond to minutes or hours. Long-duration outages usually are due to high winds, ice, lightning, vehicles that damage power lines, or equipment malfunctions that render an entire substation inoperative. For mainframes in data centers, it is usually possible, although costly, to contract for power to be delivered from two different substations, with one acting as backup.

Another type of protection is afforded by gasoline or diesel motor generators. Controls are provided that sense a power failure and automatically start the motor. Full speed is attained in less than a minute, and the generator's output can power a computer for days if necessary.

The few seconds' delay in switching power sources is enough to abort programs running on the computer and to destroy data files. To avoid this, the "uninterruptible" power supply was designed. In one version, the AC power line feeds a rectifier that furnishes direct current to an inverter. The inverter in turn drives a synchronous motor coupled to an alternator whose AC output powers the computer. While the rectifier is providing DC to the inverter, it also charges a large bank of heavy-duty batteries. As soon as a fault is detected on the main power line, the batteries are instantaneously and automatically switched over to drive the synchronous motor. Because the huge drain

4 · 20 HARDWARE ELEMENTS OF SECURITY

on the batteries may deplete them in a few minutes, a diesel generator must also be provided. The advantages of this design are:

- Variations in line frequency, amplitude, and waveform do not get through to the computer.
- Switchover from power line to batteries is undetectable by the computer. Programs keep running, and no data are lost.
- Millisecond spikes and other transients that may be responsible for equipment damage, and undetected data loss are completely suppressed.

A fuller treatment of physical threats is presented in Chapters 22 and 23 of this *Handbook*.

4.11.4 Testing. The most important aspect of any backup plan is its effectiveness. Will it work? It would be a mistake to wait for an emergency to find out. The only sensible alternative is systematic testing.

One form of test is similar to a dress rehearsal, with the actual emergency closely simulated. In this way the equipment, the people, and the procedures can all be exercised, until practice assures proficiency. Periodically thereafter the tests should be repeated, so that changes in hardware, software, and personnel will not weaken the backup capability.

4.12 RECOVERY PROCEDURES. The procedures required to recover from any system problem will depend on the nature of the problem and on the backup measures that were in place. Hardware recovery ranges from instantaneous and fully automatic, through manual repair or replacement of components, to construction, equipping, and staffing of an entirely new data center. Chapters 58 and 59 of this *Handbook* provide extensive information about these issues.

Almost every data center is a collection of equipment, with options, modifications, additions, and special features. Should it become necessary to replace the equipment, a current configuration list must be on hand and the procedures for reordering established in advance. An even better practice would be to keep a current list of *desired* equipment that could be used as the basis for replacement. Presumably, the replacements would be faster and more powerful, but additional time should be scheduled for training and conversion.

4.13 MICROCOMPUTER CONSIDERATIONS. Four factors operate to intensify the problems of hardware security as they relate to small computers:

1. Accessibility
2. Knowledge
3. Motivation
4. Opportunity

4.13.1 Accessibility. *Accessibility* is a consequence of operating small computers in a wide-open office environment rather than in a controlled data center. No security guards, special badges, man-traps, cameras, tape librarians, or shift supervisors limit access to equipment or data media in the office, as they do in a typical data center.

MICROCOMPUTER CONSIDERATIONS 4 · 21

4.13.2 Knowledge. *Knowledge* and its lack are equally dangerous. On one hand, as personal computers pervade the office environment, technical knowledge becomes widely disseminated. Where once this knowledge was limited to relatively few computer experts who could be controlled rather easily, its growing universality now makes control extremely difficult, if not impossible. On the other hand, when computers are operated by people with minimal knowledge and skills, the probability of security breaches through error and inadvertence is greatly increased.

4.13.3 Motivation. *Motivation* exists in numerous forms. It is present wherever valuable assets can be diverted for personal gain; it arises when real or fancied injustice creates a desire for revenge; and it can simply be a form of self-expression.

The unauthorized diversion of corporate assets always has provided opportunities for theft; now, with many employees owning computers at home, the value of stolen equipment, programs, and data can be realized without the involvement of third parties. When a third party is added to the equation and the thriving market in purloined personal data is factored in, the potential for data theft, a low-risk/high-return crime, is greatly increased.

Computers and networks are also a target for sabotage as well as data theft. The reliance upon such systems by governments, the military, large corporations, and other perceived purveyors of social or economic ills means that criminal acts are likely to continue. Because personal computers are now part of these systems, they are also a link to any policy or practice of which one or more groups of people disapprove. The motivation for sabotaging personal computers is more likely in the near term to increase than it is to disappear.

A third motivation for breaching computer security is the challenge and excitement of doing so. Whether trying to overcome technical hurdles, to break the law with impunity, or merely to trespass on forbidden ground, some hackers find these challenges irresistible, and they become criminal hackers. To view such acts with amused tolerance or even mild disapproval is totally inconsistent with the magnitude of the potential damage and the sanctity of the trust barriers that are crossed. Since the technology exists to lock out all but the most determined and technically proficient criminal hacker, failure to protect sensitive systems is increasingly viewed as negligence.

4.13.4 Opportunity. With so many personal computers in almost every office, with virtually no supervision during regular hours, and certainly none at other times, opportunities are plentiful for two types of security breaches: intentional by those with technical knowledge and unintentional by those without.

4.13.5 Threats to Microcomputers. Among the most significant threats to microcomputers are those pertaining to:

- Physical damage
- Theft
- Electrical power
- Static electricity
- Data communications
- Maintenance and repair

4 · 22 HARDWARE ELEMENTS OF SECURITY

4.13.5.1 Physical Damage. Microcomputers and their peripheral devices are not impervious to damage. Disk drives are extremely susceptible to failure through impact; keyboards cannot tolerate dirt or rough handling. It is essential that computers be recognized as delicate instruments and that they be treated accordingly.

Even within an access-controlled data center, where food and drinks are officially banned, it is not uncommon for a cup of coffee to be spilled when set on or near operating equipment. In an uncontrolled office environment, it is rare that one does not see personal computers in imminent danger of being doused with potentially damaging liquids. The problem is compounded by the common practice of leaving unprotected media such as CDs and DVDs lying about on the same surface where food and drink could easily reach them. Although it may not be possible to eliminate these practices entirely, greater discipline will protect data media and equipment from contamination.

As mentioned in the section on heat, damage also can result from blocking vents necessary for adequate cooling. Such vents can be rendered ineffective by placing the equipment too close to a wall or, in the case of laptops, on soft surfaces, such as carpets, that block vents on the base of the machine. Vents on top of computer housings and cathode ray tube-style displays are too often covered by papers or books that prevent a free flow of cooling air. As a result, the internal temperature of the equipment increases, so that marginal components malfunction, intermittent contacts open, errors are introduced, and eventually the system malfunctions or halts.

4.13.5.2 Theft. The opportunities for theft of personal computers and their data media are far greater than for their larger counterparts. Files containing proprietary information or expensive programs are easily copied to removable media as small as a postage stamp and taken from the premises without leaving a trace. External disk drives are small enough to be carried out in a purse or an attaché case, and new thumb-size USB drives look like key fobs to the uninitiated. (For more information about removable, miniaturized, file storage, see Chapter 1 in this *Handbook*.) The widespread practice of taking portable computers home for evening or weekend work eventually renders even the most conscientious guards indifferent. In offices without guards, the problem is even more difficult. Short of instituting a police state of perpetual surveillance, what is to be done to discourage theft? Equipment can be chained or bolted to desks, or locked within cabinets built for the purpose. Greater diligence in recording and tracking serial numbers, more frequent inventories, and a continuing program of education can help. Most of all, it is essential that the magnitude of the problem be recognized at a sufficiently high management level so that adequate resources are applied to its solution. Otherwise, there will be a continuing drain of increasing magnitude on corporate profitability.

4.13.5.3 Power. Even in a controlled data center, brownouts, blackouts, voltage spikes, sags and surges, and other electrical power disturbances represent a threat. The situation is much worse in a typical office, where personal computers are plugged into existing outlets with little or no thought to the consequences of bad power.

Some of the rudimentary precautions that should be taken are:

- Eliminating, or at least controlling, the use of extension cords, cube taps, and multiple outlet strips. Each unit on the same power line may reduce the voltage available to all of the others, and each may introduce noise on the line.

MICROCOMPUTER CONSIDERATIONS 4 · 23

- Providing line voltage regulators and line conditioners where necessary to maintain power within required limits.
- Banning the use of vacuum cleaners or other electrical devices plugged into the same power line as computers or peripheral devices. Such devices produce a high level of electrical noise, in addition to voltage sags and surges.
- Connecting all ground wires properly. This is especially important in older offices equipped with two-prong outlets that require adapter plugs. The third wire of the plug must be connected to a solid earth ground for personnel safety, as well as for reduction of electrical noise.

In addition, the use of UPSs is highly recommended for all computers and ancillary equipment. These devices are available in capacities from about 200 watts for PCs to virtually unlimited sizes for mainframes. While the power line is operational, a UPS is capable of conditioning the line by removing electrical noise, sags, spikes, and surges. When line voltage drops below a preset value, or when power is completely lost, the UPS converts DC from its internal batteries to the AC required to supply the associated equipment.

Depending on its rating and the load, the UPS may provide standby power for several minutes to several hours. This is enough time to shut down a computer normally, or in the case of large installations, to have a motor generator placed online.

The services of a qualified electrician should be utilized wherever there is a possibility of electric power problems.

4.13.5.4 Static Electricity. After one walks across a carpeted floor on a dry day, the spark that leaps from fingertip to computer may be mildly shocking to a person, but to the computer it can cause serious loss of memory, degradation of data, and even component destruction. These effects are even more likely when people touch components inside a computer without proper grounding.

To prevent this, several measures are available:

- Use a humidifier to keep the humidity above 20 percent relative.
- Remove ordinary carpeting. Replace, if desired, with static-free types.
- Use an antistatic mat beneath chairs and desks.
- Use a grounding strip near each keyboard.
- Wear a grounding bracelet when installing or repairing the components of any electronic equipment.

Touching the grounding strip before operating the computer will drain any static electricity charge through the attached ground wire, as will spraying the equipment periodically with an antistatic spray.

Some combination of these measures will protect personnel, equipment, and data from the sometimes obscure, but always real, dangers of static electricity.

4.13.5.5 Data Communications. Although personal computers perform significant functions in a stand-alone mode, their utility is greatly enhanced by communications to mainframes, to information utilities, and to other small computers, remotely via phone lines or the Internet, or through local area networks. All of the security

4 · 24 HARDWARE ELEMENTS OF SECURITY

issues that surround mainframe communications apply to personal computers, with added complications.

Until the advent of personal computers, almost all terminals communicating with mainframes were “dumb.” That is, they functioned much like teletype machines, with the ability only to key in or print out characters, one at a time. In consequence, it was much more difficult to breach mainframe security, intentionally or accidentally, than it is with today’s fully intelligent personal computers.

The image of thousands of dedicated hackers dialing up readily available computer access numbers or probing Internet addresses, for illicit fun and games or for illegal financial gain, is no less disturbing than it is real. Countermeasures are available, including:

- Two-way encryption (see Chapter 7)
- Frequent password changes (see Chapter 28)
- Automatic call-back before logging on
- Investigation of unsuccessful logons
- Monitoring of hackers’ bulletin boards (see Chapters 12 and 15)
- Firewalls to restrict traffic into and out of the computer (see Chapter 26)
- Antivirus software (see Chapter 41)

Legislation that makes directors and senior officers personally liable for any corporate losses that could have been prevented should have a marked effect on overcoming the current inertia. Prudence dictates that preventive action be taken before, rather than corrective action after, such losses are incurred.

4.13.6 Maintenance and Repair. A regular program of preventive maintenance should be observed for every element of a personal computer system. This should include scheduled cleaning of disk drives and their magnetic heads, keyboards, and printers. A vital element of any preventive maintenance program is the frequent changing of air filters in every piece of equipment. If this is not done, the flow of clean, cool air will be impeded, and failure will almost surely result.

Maintenance options for personal computers, in decreasing order of timeliness, include:

- Onsite management by regular employees
- Onsite maintenance by third parties under an annual agreement
- On-call repair, with or without an agreement
- Carry-in service
- Mail-in service

As personal computers are increasingly applied to functions that affect the very existence of a business, their maintenance and repair will demand more management attention. Redundant equipment and onsite backup will always be effective, but the extended time for offsite repairs will no longer be acceptable. For most business applications, “loaners” or “swappers” should be immediately available, so that downtime will be held to an absolute minimum. Management must assess the importance of each functioning personal computer and select an appropriate maintenance and repair policy.

HARDWARE SECURITY CHECKLIST 4 · 25

Accessibility, knowledge, motivation, and opportunity are the special factors that threaten every personal computer installation. Until each of these factors has been addressed, no system can be considered secure.

4.14 CONCLUSION. This chapter has dealt principally with the means by which hardware elements of a data processing system affect the security and integrity of its operations. Many safeguards are integral parts of the equipment itself; others require conscious effort, determination, and commitment.

An effective security program—one that provides both decreased likelihood of computer catastrophe and mitigation of the consequences of damage—cannot be designed or implemented without considerable expenditures of time and money. As with other types of loss avoidance, the premium should be evaluated against the expected costs. Once a decision has been made, however, this equivalent to an insurance policy should not be permitted to lapse. The premiums must continue to be paid in the form of periodic testing, continuous updating, and constant vigilance.

For more detailed information about risk management, see Chapter 62 in this *Handbook*. For a discussion of insurance policies against information systems disasters of all kinds, see Chapter 60.

4.15 HARDWARE SECURITY CHECKLIST

Mainframes

- Are security and integrity requirements considered when selecting new equipment?
- Is a schedule of preventive maintenance enforced?
- Is a log kept of all computer malfunctions and unscheduled downtime?
- Is there an individual with responsibility for reviewing the log and initiating action?
- Are parity checks used wherever possible?
- Is there an established procedure for recording parity errors and recovering from them?
- Are forward-acting or error-correcting codes used when economically justified?
- Do operators follow prescribed procedures after a read error or other machine check halt?
- Are all operator interventions logged and explained?
- Is a job log maintained, and is it compared regularly with an authorized run list?
- Is the interval timer used to prevent excessively long runs?
- Are storage protect features such as data locks and read-only paging used?
- Are keys to software data locks adequately protected?
- Are precautions taken to prevent loss of data from volatile memory during power interruptions?
- Are standard internal and external tape and disk labeling procedures enforced?
- Are write-enable protection rings always removed from tape reels immediately after use?
- Is there a rule that new tapes and disks must be tested or certified prior to use? At regular intervals thereafter?

4 · 26 HARDWARE ELEMENTS OF SECURITY

- Are tapes and disks refinished or replaced before performance is degraded?
- Are air conditioners adequate for peak thermal loads? Are air conditioners backed up?
- Is there a schedule for frequent filter changes?
- Have all static electricity generators been disabled?
- Have all sources of water damage been eliminated?
- Is good housekeeping enforced throughout the facility?
- Is access to data terminals restricted?
- Are terminals and surrounding areas examined frequently to detect passwords carelessly left about?
- Is a log maintained of unsuccessful attempts to enter the computer from terminals?
- Is the log used to prevent further attempts?
- Is a log maintained of all successful entries to sensitive data?
- Is the log used to verify authorizations?
- Are terminals equipped with automatic identification generators?
- Are test procedures adequate to assure high-quality data transmissions?
- Is cryptography or scrambling used to protect sensitive data?
- Has a complete backup plan been formulated? Is it updated frequently?
- Does the backup plan include training, retraining, and cross-training of personnel?
- Is onsite backup available for the central processing unit? For peripherals?
- Does your backup site advise you of all changes to its hardware configuration and operating system?
- Does your backup site have enough free time available to accommodate your emergency needs?
- Do you monitor power-line voltage and frequency?
- Are the effects of brownouts, dim-outs, and blackouts known?
- Is advance warning available, and if so, is there a checklist of actions to be taken?
- Are power correctors in use? Voltage regulators? Line conditioners? Lightning spark gaps?
- Is backup power available? Dual substation supply? Motor generators? Uninterruptible power supplies?
- Does your equipment provide automatic restart and recovery after a power failure?
- Are backup plans tested realistically? At frequent intervals?

Microcomputers

In addition to the appropriate items just listed:

- Are removable disks always kept in a closed container when not actually mounted in a disk drive?
- Is it forbidden to put food or drink on or near computer equipment?
- Are personal computers securely fastened to prevent dropping or theft?
- Are air vents kept free?
- Are accurate inventories maintained?

FURTHER READING 4 · 27

- Is electrical power properly wired?
- Are uninterruptable power supplies in place?
- Has static electricity been eliminated?
- Are data communications secure?
- Is there an effective maintenance plan?

4.16 FURTHER READING

- Ayers, J. E. *Digital Integrated Circuits: Analysis and Design*. Boca Raton, FL: CRC Press, 2003. (Second edition scheduled for publication in 2009.)
- Clements, A. *Principles of Computer Hardware*, 4th ed. New York: Oxford University Press, 2006.
- Horak, R. *Telecommunications and Data Communications Handbook*. Hoboken, NJ: Wiley-Interscience, 2007.
- Kerns, D. V. *Essentials of Electrical and Computer Engineering*, 2nd ed. Upper Saddle River, NJ: Prentice-Hall, 2004.
- Pattern, D. A., and J. L. Hennessy. *Computer Organization and Design: The Hardware Software Interface*, 3rd ed. Los Angeles: Morgan Kaufmann, 2007.
- Stallings, W. *Computer Organization and Architecture: Designing for Performance*, 7th ed. Upper Saddle River, NJ: Prentice-Hall, 2005.

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 5

DATA COMMUNICATIONS AND INFORMATION SECURITY

Raymond Panko and Eric Fisher

5.1 INTRODUCTION	5·2	5.6.5 Acknowledgment Numbers	5·21
5.2 SAMPLING OF NETWORKS	5·2	5.6.6 Window Field	5·21
5.2.1 Simple Home Network	5·2	5.6.7 Options	5·21
5.2.2 Building LAN	5·4	5.6.8 Port Numbers	5·22
5.2.3 Firms' Wide Area Networks (WANs)	5·5	5.6.9 TCP Security	5·23
5.2.4 Internet	5·7		
5.2.5 Applications	5·9		
		5.7 USER DATAGRAM PROTOCOL	5·23
5.3 NETWORK PROTOCOLS AND VULNERABILITIES	5·9	5.8 TCP/IP SUPERVISORY STANDARDS	5·24
		5.8.1 Internet Control Message Protocol (ICMP)	5·24
5.4 STANDARDS	5·9	5.8.2 Domain Name System (DNS)	5·25
5.4.1 Core Layers	5·10	5.8.3 Dynamic Host Configuration Protocol (DHCP)	5·26
5.4.2 Layered Standards Architectures	5·10	5.8.4 Dynamic Routing Protocols	5·27
5.4.3 Single-Network Standards	5·11	5.8.5 Simple Network Management Protocol (SNMP)	5·27
5.4.4 Internetworking Standards	5·13		
5.5 INTERNET PROTOCOL (IP)	5·14		
5.5.1 IP Version 4 Packet	5·14	5.9 APPLICATION STANDARDS	5·28
5.5.2 IP Version 6	5·16	5.9.1 HTTP and HTML	5·28
5.5.3 IPsec	5·17	5.9.2 E-Mail	5·28
5.6 TRANSMISSION CONTROL PROTOCOL (TCP)	5·18	5.9.3 Telnet, FTP, and SSH	5·28
5.6.1 Connection-Oriented and Reliable Protocol	5·18	5.9.4 Other Application Standards	5·29
5.6.2 Reliability	5·20		
5.6.3 Flag Fields	5·20	5.10 CONCLUDING REMARKS	5·29
5.6.4 Octets and Sequence Number	5·21	5.11 FURTHER READING	5·29
		5.12 NOTES	5·29

5 · 2 DATA COMMUNICATIONS AND INFORMATION SECURITY

5.1 INTRODUCTION. Sometimes, an attacker can simply walk up to a target computer. In most cases, however, attackers must use networks to reach their targets. Some attacks even aim *at* networks, trying to bring down local area networks, wide area networks, and even the global Internet. This chapter provides an overview of networking to help readers of this *Handbook* when they come across networking concepts in other chapters or in other contexts. This chapter covers a limited number of networking concepts. Specifically, it focuses on aspects of networking that are most relevant to security.

Before beginning, readers should note three important pieces of terminology that pervade the chapter.

1. This chapter often uses the term *octet*, which is a *byte*—a collection of eight bits. Networking grew out of electrical engineering, where octet is the preferred term; it is also widely used in the international technical community.
2. The second term is *host*. Any device attached to the global Internet is called a host. This includes everything from large server hosts to client PCs, personal digital assistants, mobile telephones, and even Internet-accessible coffeepots.
3. We will distinguish between the terms *internet* and *Internet*; the latter refers to the global Internet. However, *internet* spelled in lower case is either the Internet layer in the TCP architecture (see Section 5.6) or a collection of networks that is not the global Internet.

5.2 SAMPLING OF NETWORKS. This section looks briefly at a series of increasingly complex networks, giving the reader a high-level overview of what networks look like in the real world.

5.2.1 Simple Home Network. Exhibit 5.1 shows a simple home PC network. The home has two personal computers. The network allows the two PCs to share files and the family’s single laser printer. The network also connects the two computers to the Internet.

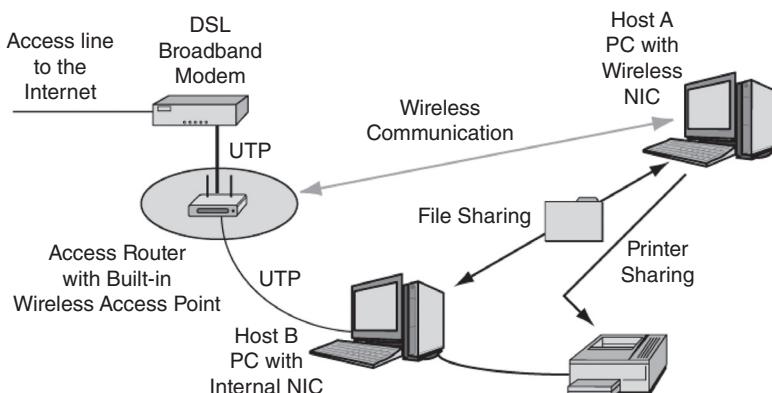


EXHIBIT 5.1 Simple Home Network

SAMPLING OF NETWORKS 5 · 3

5.2.1.1 Access Router. The heart of this network is its *access router*. This small device performs a variety of functions, most importantly these five:

1. It performs as a switch. When one PC in the home sends messages (called packets) to the other hosts, the switch transfers the packets between them.
2. The access router is a wireless access point (WAP), which permits wireless computers to connect to it. Host A connects to the access router wirelessly.
3. A router connects a network to another network—in this case, it connects the internal network to the global Internet.
4. To use the Internet, each computer needs an Internet Protocol (IP) address. We will see later that IP is the main protocol that governs communication over the Internet. The access router has a built-in Dynamic Host Configuration Protocol (DHCP) server that gives each home PC an IP address.
5. The router provides *network address translation* (NAT), which hides internal IP addresses from potential attackers. Most routers also have a firewall for added security. WAPs are easily exploited if not configured with proper authentication security. Wireless signals can be transferred up to 800 feet away or more with special equipment. Without constant monitoring to defeat intrusions, an attacker can connect to an access point without the user's knowledge and intercept all passing traffic. Using NAT is essential to keeping a home network secure. Users should always enable this feature in order to prevent their hosts from being directly accessible to the public Internet, where direct scans and attacks are prevalent.

5.2.1.2 Personal Computers. Each of the two PCs needs circuitry to communicate over the network. Traditionally, this circuitry came in the form of a printed circuit board, so the circuitry was called the computer's *network interface card* (NIC). In most computers today, the circuitry is built into the computer; there is no separate printed circuit board. However, the circuitry is still called the computer's NIC.

In this small network, the two computers share their files. Given the wireless access capability of the network, drive-by hackers could potentially read shared files as well. File sharing without strong wireless security is dangerous. It is important to set up *Wi-Fi Protected Access* (WPA or WPA2) or 802.11i security in pre-shared key (PSK) mode on both the access router/access point and each of the client PCs.

It is important to configure the PCs for security. Although NAT by itself is strong, and most routers also provide stateful-inspection firewalls (see Chapter 26 in this *Handbook*), some attacks will inevitably get through to the internal network. Hosts must have strong firewalls, antivirus programs, and antispyware programs (see Chapter 41); and they must be updated automatically when security patches are released by the operating system vendor and by application program vendors (see Chapter 40).

5.2.1.3 UTP Wiring. In Exhibit 5.1, Host B connects to the access router via copper wiring called a UTP cable, Ethernet (IEEE 802.3) cable, or commonly Cat5 (Cat5 stands for Category 5 cabling, defined in standard ANSI/TIA/EIA-568-A). It uses four-pair *unshielded twisted pair* (UTP) wiring inside the cord jacket. As Exhibit 5.2 shows, a UTP cord contains eight copper wires organized as four pairs. The two wires of each pair are twisted around each other. The *RJ-45* connectors at the ends of a UTP cord look like RJ-11 home telephone connectors but are a little wider. (*RJ* means

5 · 4 DATA COMMUNICATIONS AND INFORMATION SECURITY

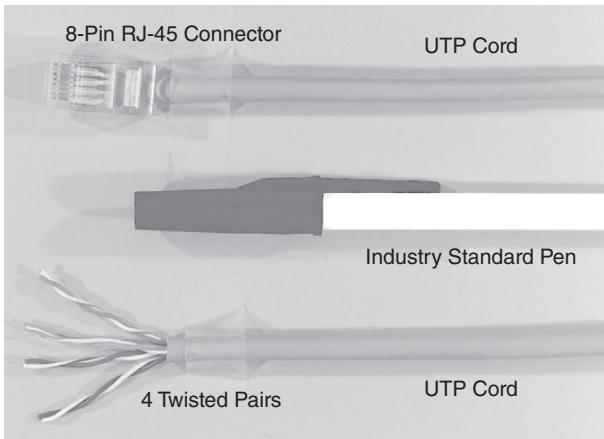


EXHIBIT 5.2 Unshielded Twisted Pair (UTP) Wiring Cord

Registered Jack and originally referred to Bell System order codes; it is now defined by the Administrative Council for Terminal Attachment, ACTA.)

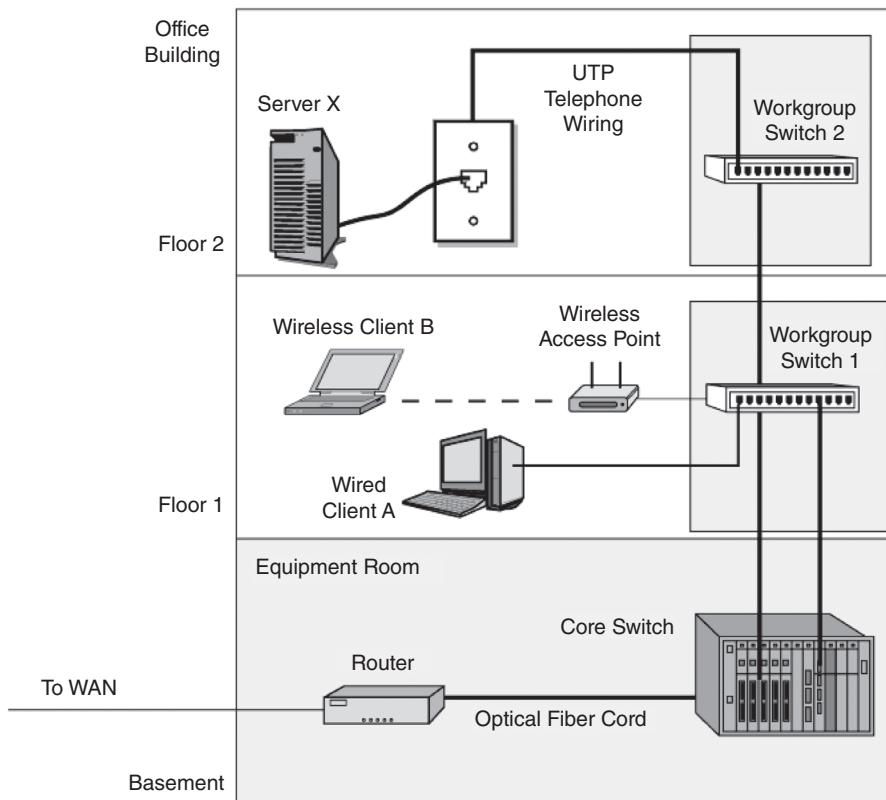
5.2.1.4 Internet Access Line. The home network needs an Internet access line to connect the home to the Internet. In Exhibit 5.1, this access line is a *digital subscriber line* (DSL) high-speed access line, and the home connects to this access line via a small box called a *DSL modem*. (The DSL modem connects to the access router via a UTP cord; it connects to the wall jack via an ordinary telephone cord.) Other Internet access technologies include slow telephone modems, fast cable modems, geosynchronous-satellite connections, and even wireless access systems. Most of these technologies are called *broadband access lines*. In general, *broadband* simply means *very fast*, although in radio transmission it describes a wide range of frequencies.

5.2.2 Building LAN. The home network shown in Exhibit 5.1 is a *local area network* (LAN). A LAN operates on a customer's premises—the property owned by the LAN user. (For historical reasons, *premises* is always spelled in the plural.) In the case of the home network, the premises consist of the user's home or apartment. Exhibit 5.3 shows a much larger LAN. Here, the premises consist of a corporate multistory office building.

On each floor, computers connect to the floor's workgroup switch via a UTP cord or a wireless access point. The workgroup switch on each floor connects to a *core switch* in the basement equipment room. The router in the basement connects the building LAN to the outside world.

Suppose that Client A on Floor 1 sends a packet to Server X on Floor 2. Client A sends the packet to Workgroup Switch 1 on the first floor. That workgroup switch sends the packet down to the core switch in the basement. The core switch then sends the packet up to Workgroup Switch 2, which passes the packet to Server X.

UTP is easy to wiretap, allowing attackers to read all packets flowing through the cord. Telecommunications closets should be kept locked at all times, and cords should be run through thick metal wiring conduits wherever possible (for more details of physical and facilities security, see Chapters 22 and 23 in this *Handbook*). UTP also generates weak radio signals when traffic flows through it. It is possible to read these signals from some distance away using highly specialized equipment. Newer

SAMPLING OF NETWORKS 5 · 5**EXHIBIT 5.3** Building LAN

specifications called Cat5e and Cat6 were developed to cut down on interference and cables can be purchased with shielding, but even then it is possible to eavesdrop.

Eavesdropping by tapping a UTP cable is not difficult once physical access is gained; however, typically there are far easier ways of gaining access to a network and far more desirable targets. Eavesdropping on a wire would reveal any passing traffic, but eavesdropping on a router or switch would reveal passing traffic on *many* wires. Physical security is an important facet of network security and must be properly addressed, but most attacks today rely on more virtual vulnerabilities.

For more extensive details of LAN security, see Chapter 25 in this *Handbook*.

5.2.3 Firms' Wide Area Networks (WANs). Although LANs operate within a company's premises, *wide area networks* (WANs) connect geographically separate sites—usually within a single corporation. Corporations do not have the regulatory rights-of-way needed to run wires through public areas. For WAN service, companies must use companies called *carriers* that do have these rights-of-way.

Exhibit 5.4 shows that most firms use multiple-carrier WANs. In the exhibit, some sites in this company are connected by point-to-point *leased lines* from a telephone company. The companies also subscribe to *switched network services* that transfer traffic between several sites. The exhibit shows that these switched network services use the *Frame Relay* technology. The company uses two separate Frame Relay networks—one to connect its own sites to one another and another to connect it to another firm.

5 · 6 DATA COMMUNICATIONS AND INFORMATION SECURITY

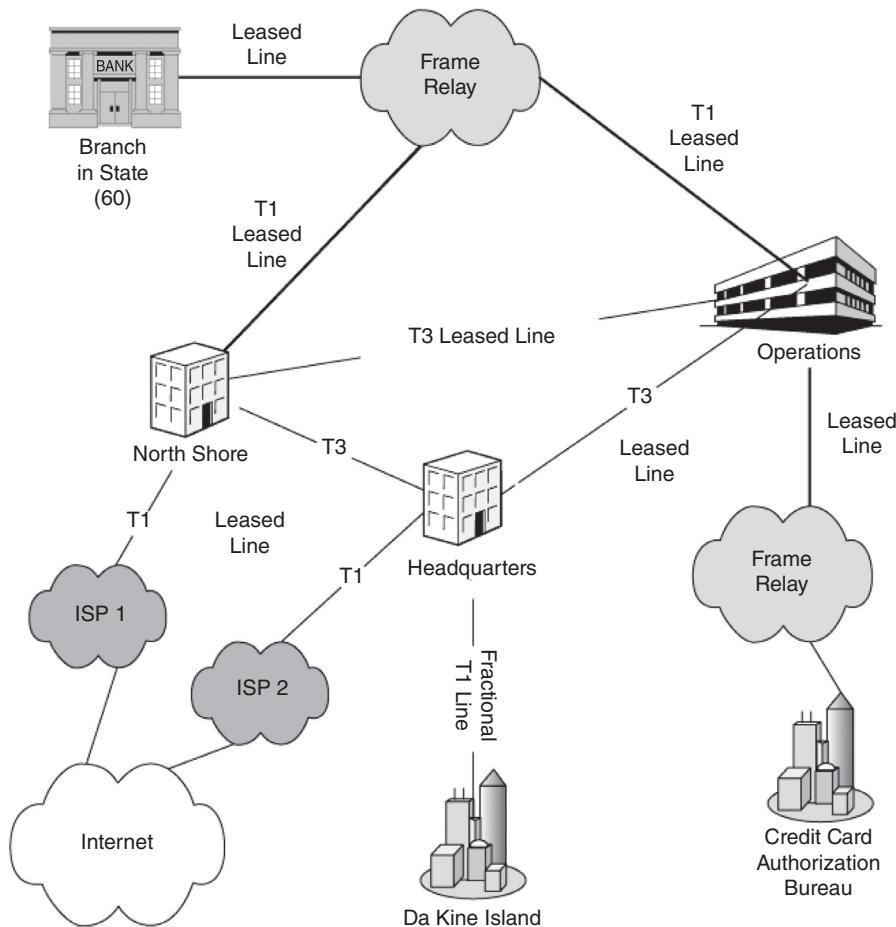
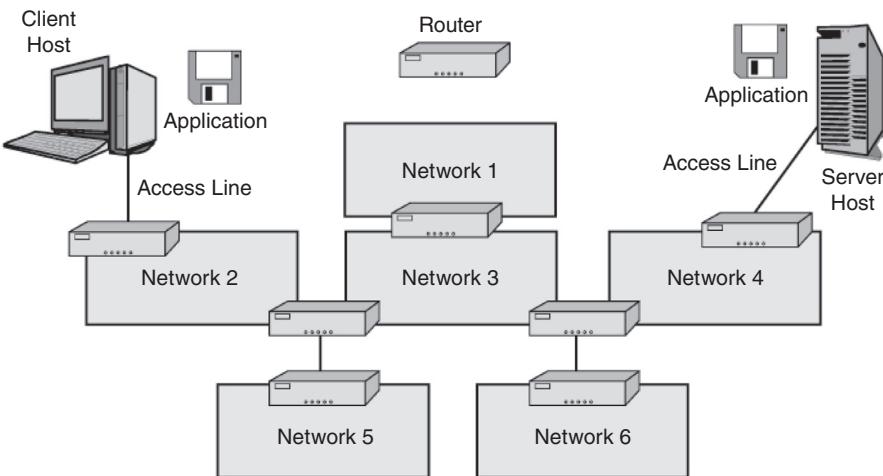


EXHIBIT 5.4 Wide Area Networks (WANs)

Carrier technology is usually considered more secure by security professionals due to its closed-access nature. Unlike the Internet, which allows anyone to connect to it, only commercial firms may connect to carrier WANs, which makes attacker access very difficult. However, attacker access is not impossible. For example, if an attacker hacks a computer owned by the carrier (or even by a customer), this breach may permit access.

In addition, the carrier alone knows how it routes traffic through its network. This should stymie attackers even if they somehow get access to the network. However, such *security through obscurity* is considered a poor practice by security professionals because it is possible for attackers who hack carrier computers to get access to routing information. (Attackers usually have much simpler attack vectors; see Chapters 15 and 19 in this *Handbook* for more details.)

Although carrier technology is more secure, it is also extremely expensive. With the development of virtual private networks (VPNs), companies can connect geographically disparate groups of computers *virtually* over the common internet. This provides much of the security benefit of WANs while cutting the implementation cost dramatically. See Chapter 32 for more information about VPN security.

SAMPLING OF NETWORKS 5 · 7**EXHIBIT 5.5** Internet

5.2.4 Internet. By the end of the 1970s, there were many LANs and WANs in the world. Many of the WANs were nonprofit networks that connected universities and research institutions. Unfortunately, computers on one network could not talk to computers on other networks. To address this problem, the Defense Advanced Research Projects Agency (DARPA) created ARPANET in 1969, the origin of today's Internet, based on the pioneering conceptual design for what J. C. R. Licklider called the *Intergalactic Computer Network* in a 1963 paper. By definition, an *internet* connects individual networks together. Later, commercial networks were allowed to join later versions of ARPANET, and it became the Internet we know today.

Exhibit 5.5 shows that devices called *routers* connect the individual networks together. Initially, these devices were called *gateways*. The term *gateway* was used instead of “router” in some early standards, but most vendors have now adopted the name “router.” There are two exceptions, the first being Microsoft, which still tends to call routers “gateways.” The second is the router directly accessible to a network, and thus the first hop when exiting a network is often called the *default gateway*.

Any computer on any network on the Internet can send messages to any computer on any other network on the Internet. The messages that travel all the way from one computer to another across the Internet are called *packets*.

Exhibit 5.6 shows that the packet travels all the way from the source host to the destination host. Along the way, it is routed through different networks until arriving at its destination.

The global Internet uses a suite of communication protocols known as *Transmission Control Protocol/Internet Protocol* (TCP/IP). In addition, many firms build separate internal TCP/IP networks for their own communication. These internal networks are called *intranets* to distinguish them from the *Internet*.

Initially, security on internal networks was comparatively light because it was assumed that external attackers would have a difficult time getting into corporate intranets. However, if a hacker takes over an internal computer connected to the intranet, light security becomes a serious problem. Consequently, most firms have been progressively hardening their intranet security.

Exhibit 5.7 shows that individual homes and corporations connect to the Internet via carriers called *Internet service providers* (ISPs). The Internet has many ISPs, but

5 · 8 DATA COMMUNICATIONS AND INFORMATION SECURITY

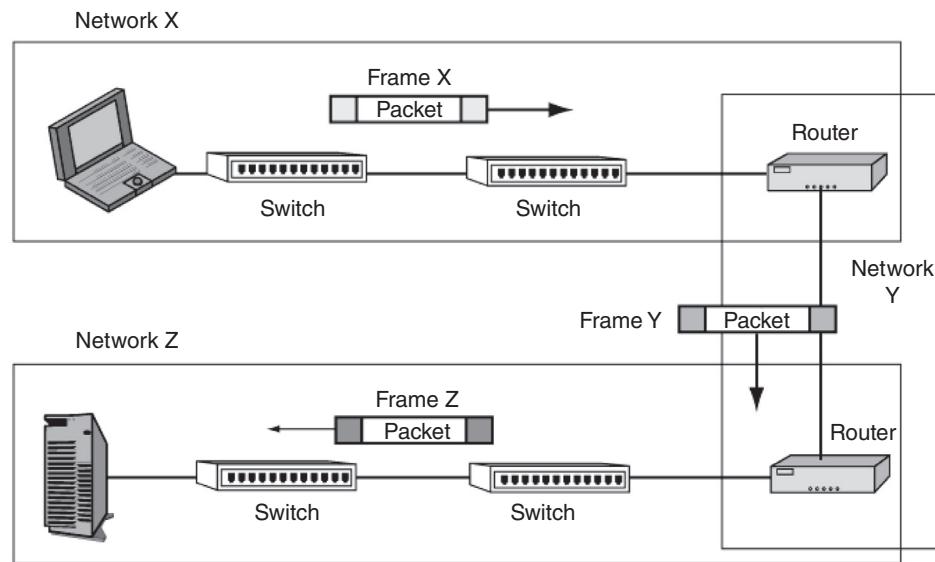


EXHIBIT 5.6 Frames and Packets

they all connect at centers that usually are called *network access points* (NAPs). These connections allow global communications for all connected hosts.

Most ISPs are commercial organizations run for profit to provide Internet access for home users. There is no central access control to the Internet; however, there are central agencies for controlling Domain Name Systems (DNSs) called *registrars*.

When the Internet was designed in the late 1970s, there was a conscious decision to promote openness and not to add the burdens of security. As a consequence of a lack of security technology and open access to almost anyone, the Internet is a security nightmare. Companies that transmit sensitive information over the Internet need to consider cryptographic protections. (See Chapters 7, 32, 33, 34, 35, and 37 in this *Handbook* for more details of cryptography and other means for achieving security on networks.)

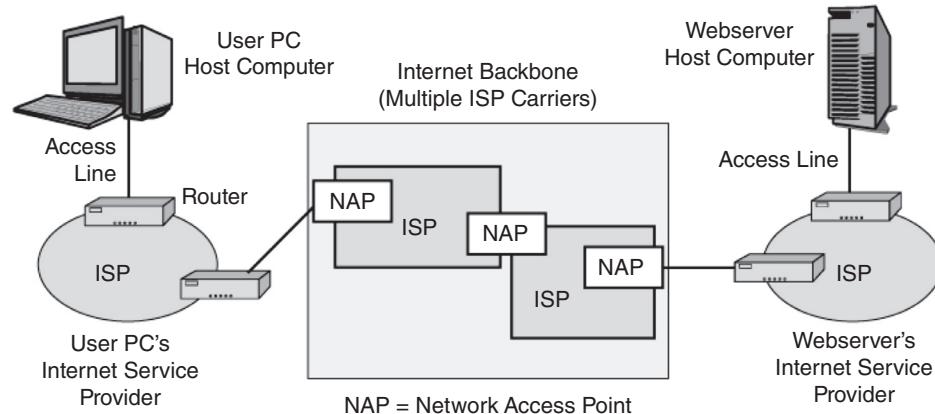


EXHIBIT 5.7 Internet Service Providers (ISPs)

STANDARDS 5 · 9

5.2.5 Applications. Although the inner workings of the Internet rely on networks, most users are only aware of the applications they commonly use that run on top of networks. Familiar personal applications include the World Wide Web, email, and instant messaging, among many others. Corporations use some of these applications, but they also use many business-specific applications, such as accounting, payroll, billing, and inventory management. Often, business applications are transaction-processing applications, which are characterized by high volumes of simple repetitive transactions. The traffic volume generated by transaction-processing and other business-oriented applications usually far outweighs the traffic of personal applications in the firm. (See Chapter 30 in this *Handbook* for details of e-commerce security.)

All programs have bugs, including security vulnerabilities. There are many applications, and keeping track of application vulnerabilities and constantly patching many applications is an enormous task that is all too easy to put off or complete only partially. (See Chapter 40 for an overview of patch management.) Also, each application must be configured with options that have high security, and security must be managed on each application (e.g., anti-virus and spam blocking in email). (See Chapter 20 for a review of spam and anti-spam measures.)

5.3 NETWORK PROTOCOLS AND VULNERABILITIES. The products of different network vendors must be able to work together (interoperate). This is possible only if there are strong communication standards to govern how hardware and software processes interact. With such standards, two or more programs can interoperate effectively.

Standards raise three security issues. One is the standard itself. For instance, the TCP standard discussed later in this chapter is difficult to attack because an attacker cannot send a false message unless he or she can guess the sequence number of the next message. This normally is very difficult to do. However, if the attacker sends an RST (reset) message, which terminates a connection, this protection is greatly reduced. In fact, it is fairly easy to send RST messages that close legitimate open connections.

A second issue is security built into the standard. Most standards were created without security, and security was added only in later versions, sometimes in an awkward way. For instance, IP, which is the main protocol for delivering packets over the Internet, originally had no security. The *IP security (IPsec, pronounced eye-pea-sek)* standards were created to address this weakness, but IPsec is burdensome and not widely used.

Another security weakness of early versions of IP, including the widely used IPv4, is the limitation on address space due to the 32 bit address field in the IPv4 packet (yielding an address space of about 4×10^9); as this edition of this *Handbook* goes to press, IPv4 address exhaustion is being addressed by the migration to IPv6, with its 128-bit addresses (an address space of about 3×10^{38}).

A further issue is the security of the implementation of standards in vendor products. Most attacks that aim at standards weaknesses attack vendor products that have security vulnerabilities unrelated to the protocols they implement.

5.4 STANDARDS. Networks and network security depend on standards. Standards have permitted global interconnectivity, unlike the early years of networking when proprietary products dominated the world of computing and interconnection was difficult or impossible.

5 · 10 DATA COMMUNICATIONS AND INFORMATION SECURITY

EXHIBIT 5.8 Three Standards Core Layers

Super Layer	Description
Application	Communication between application programs on different hosts attached to different networks on a network.
Internetworking	Transmission of packets across a routed internet. Packets contain application-layer messages.
Single network	Transmission of packets across a single-switched network.

5.4.1 Core Layers. Standards are complex, and when people deal with complex problems, they usually break these problems into smaller parts and have different specialists work on the different parts. Exhibit 5.8 shows that standards are divided into three core layers that collectively have the functionality needed to allow an application program on one network in an internet to interoperate with another program on another computer on another network.

At the *application core layer*, the two applications must be able to interact effectively. For instance, in World Wide Web access, the two application programs are the browser on the client PC and the Web server program on the Web server. The standard for Web interactions is the *Hypertext Transfer Protocol* (HTTP). Both the browser and the Web server applications have to send messages that comply with the HTTP standard.

The middle layer is the *internet core layer*. Standards at this layer govern how packets are delivered across a routed internet. One of the main standards at the internet core layer is the Internet Protocol (IP). We will see other internetworking standards later.

The lowest core layer is the *single-network core layer*. Standards at this layer govern the transmission of packets across the switches and transmission lines in a single-switched network (a LAN or WAN).

5.4.2 Layered Standards Architectures. Standards are created by standards agencies. These standards agencies first create detailed layering plans for creating standards. These specific layering plans are called *layered standards architectures*. Afterward, standards agencies create standards in the individual layers. Exhibit 5.9 shows two popular layered standards architectures and relates these standards architectures to the three core layers we saw earlier.

The *Internet Engineering Task Force* (IETF) is the standards agency for the Internet. Its standards architecture is called TCP/IP—a name taken from two of its most

EXHIBIT 5.9 Layered Standards Architectures

Super Layer	TCP/IP	OSI	Hybrid TCP/IP-OSI
Application	Application	Application Presentation Session	Application
Internet	Transport Internet	Transport Network	Transport Internet
Network	Subnet access	Data link Physical	Data link Physical

STANDARDS 5 · 11

important standards, TCP and IP. Exhibit 5.9 shows that TCP/IP has four layers. The bottom layer, the *subnet access layer*, corresponds to the single-network core layer. The top layer, in turn, is the *application layer*, which corresponds to the application core layer. The two middle layers—the *internet* and *transport* layers—correspond to the internet core layer. TCP/IP focuses primarily on internet working. Dividing this core layer into two TCP/IP layers permits greater division of labor in standards development.

The other standards architecture shown in the figure is *OSI*, which is rarely spelled out by its full name, the *Reference Model of Open Systems Interconnection*. OSI is governed by two standards agencies. One is ISO, the *International Organization for Standardization*. The other is ITU-T, the *International Telecommunications Union–Telecommunications Standards Sector*. (The official names and the official acronyms do not match because they originated in different languages.)

Exhibit 5.9 shows that OSI divides the three core layers into a total of seven layers. OSI single networks use standards at two layers—the *physical* and *data link* layers. OSI's market dominance is so strong at the physical and data link layers that the IETF rarely develops standards at these layers. The *subnet access* indication in the TCP/IP framework basically means *Use OSI standards here*.

Neither of these two standards architectures dominates. What nearly all firms use today is the hybrid TCP/IP–OSI standards architecture, which Exhibit 5.9 illustrates. This hybrid architecture uses OSI standards at the physical and data link layer and TCP/IP standards at the internet and transport layer. Corporations also use standards from some other standards architectures at the internet and transport layers, but TCP/IP standards dominate.

At the application core layer, the situation is complex. Both OSI and TCP/IP standards are used, often in combination. In fact, OSI standards often reference TCP/IP standards and vice versa. Although OSI and TCP/IP are often viewed as rivals, this is not the case at all. Several other standards agencies also create application layer standards, complicating the picture even further.

5.4.3 Single-Network Standards. As just noted, OSI standards dominate in the two single-network layers—the physical and data link layers. Exhibit 5.10 shows how the physical and data link layers are related.

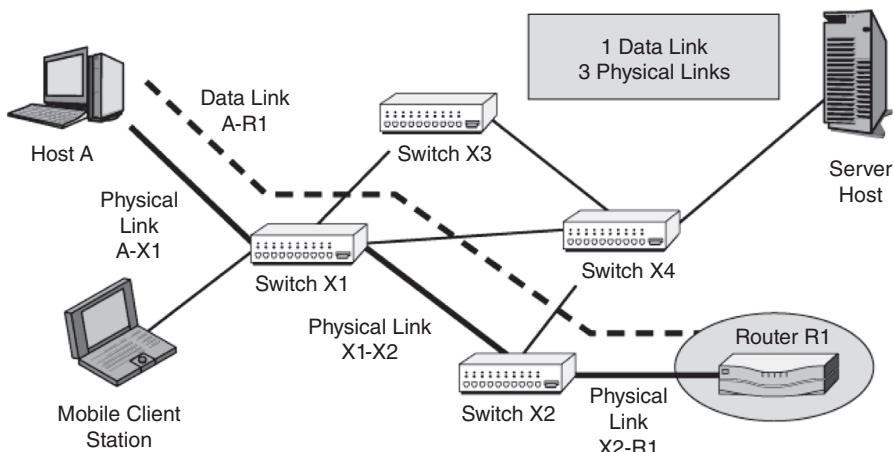


EXHIBIT 5.10 Physical and Data Link Layers

5 · 12 DATA COMMUNICATIONS AND INFORMATION SECURITY

5.4.3.1 Data Link Layer. The path that a frame takes through a single network is called the frame's data link. In Exhibit 5.10, the data link runs between Host A and Router R1. This data link passes through Switch X1 and Switch X2.

The source computer sends the frame to the first switch, which forwards the frame to the next switch along the data link, which forwards the frame further. The last switch along the data link passes the frame to the destination computer (or router, if the packet in the frame is destined for a computer on another network).

5.4.3.2 Physical Layer. Physical-layer standards govern the physical connections between consecutive devices along a data link. In Exhibit 5.10, these physical links are A–X1, X1–X2, and X2–R1. Earlier, we saw one popular transmission medium, unshielded twisted pair wire in Cat5 cables. UTP dominates in links between computers and workgroup switches (see Exhibit 5.3). UTP signals typically involve voltage changes. For instance, a high voltage may indicate a 1, while a low voltage may indicate a 0. (Actual voltage patterns usually are much more complex.)

For longer distances and very high speeds, another popular transmission medium is optical fiber, which sends light signals through thin glass tubes. Optical fiber signals actually are very simple. In a clock cycle, the light is turned on for a 1 or off for a 0.

UTP cords act like radio antennas when they carry signals. Some of the signal always radiates out, allowing people to intercept transmission signals by placing devices near (but not touching) the cord. Intercepting and interpreting electromagnetic emissions from computing devices is called van Eck phreaking (also famously codenamed "TEMPEST" by the NSA) after the Dutch scientist Wim van Eck published a paper in 1985 demonstrating how to monitor and reconstitute leaked signals from *cathode-ray terminals* (CRTs). In contrast, optical fiber requires physically tapping into the fiber cords. Physical wiretapping can also be done with UTP, but there are often far easier methods to intercept or steal traffic rather than trying to physically tap the wires.

Wireless transmission uses radio waves. This permits mobile devices to be served in ways never before possible. Wireless transmission is used for both LAN and WAN transmission.

Radio signals spread widely, even when dish antennas are used. Consequently, it is very easy for eavesdroppers to listen in on radio transmissions and do other mischief. Radio signals must be strongly encrypted, and the parties must be strongly authenticated to prevent impostors from sending radio transmission.

Radio signaling is very complex. Most radio signaling uses spread spectrum transmission, in which the information is sent over a wide range of frequencies. *Spread spectrum* transmission is used to improve propagation reliability. Radio transmission has many propagation problems, such as interference from other sources. Many propagation problems occur only at certain frequencies. By spreading the signal across a wide spectrum of frequencies and doing so redundantly, the signal will still be intelligible even if there are strong problems at some frequencies.

Prabakar Prabakaran summarized the benefits of spread-spectrum communications as follows:

Spread-spectrum systems provide some clear advantages to designers ... [H]ere are nine benefits that designers can expect when using a spread-spectrum-based wireless system.

1. Reduced crosstalk interference: In spread-spectrum systems, crosstalk interference is greatly attenuated due to the processing gain of the spread spectrum system as described earlier ...
2. Better voice quality/data integrity and less static noise ...

STANDARDS 5 · 13

3. Lowered susceptibility to multipath fading ...
4. Inherent security: In a spread spectrum system, a PN [pseudo-random number] sequence is used to either modulate the signal in the time domain (direct sequence systems) or select the carrier frequency (frequency hopping systems). Due to the pseudo-random nature of the PN sequence, the signal in the air has been “randomized.” Only the receiver having the exact same pseudo-random sequence and synchronous timing can de-spread and retrieve the original signal. Consequently, a spread spectrum system provides signal security that is not available to conventional analog wireless systems.
5. Co-existence: A spread spectrum system is less susceptible to interference than other non-spread spectrum systems. In addition, with the proper designing of pseudo-random sequences, multiple spread spectrum systems can co-exist without creating severe interference to other systems. This further increases the system capacity for spread spectrum systems or devices.
6. Longer operating distances ...
7. Hard to detect: Spread-spectrum signals are much wider than conventional narrowband transmission (of the order of 20 to 254 times the bandwidth of narrowband transmissions). Since the communication band is spread, it can be transmitted at a low power without being detrimentally by background noise ...
8. Hard to intercept or demodulate: The very foundation of the spreading technique is the code used to spread the signal ...
9. Harder to jam: The most important feature of spread spectrum is its ability to reject interference¹

The military uses frequency-hopping spread-spectrum (FHSS) transmission for security. Military spread-spectrum transmission works in such a way that makes intercepting transmissions very difficult. Civilian spread-spectrum transmission, in contrast, is designed to make connecting simple and therefore offers relatively little security.

Switches spend almost all of their time forwarding frames. However, switches spend some of their time exchanging supervisory information packets with one another to keep the network running efficiently. For example, in *Ethernet* (IEEE 802.3), which dominates LAN standards, if there are loops among the switches, the network will malfunction. If a switch detects a loop, it sends supervisory packets to other switches. The switches in the network then communicate until they determine the most appropriate path and disable other ports to prevent the internal looping. This process is governed by the *Spanning Tree Protocol* (STP, part of IEEE 802.1) or the newer *Rapid Spanning Tree Protocol* (RSTP, defined in IEEE 802.1 w and now part of IEEE 802.1D-2004).

Attackers can create denial-of-service (DOS) attacks on the switches in a network by impersonating a switch and sending a flood of false messages to the network’s real switches indicating the presence of a loop. The switches may spend so much of their time reorganizing the network that they will be unable to serve legitimate traffic. They also can attack several other supervisory protocols to make switches unavailable for processing normal packets. The 802.1AE standard is designed to limit switch-to-switch communication to authenticated switches.

5.4.4 Internetworking Standards. As noted earlier, the IETF divided the internetworking core layer into two layers—the internet and transport layers. Exhibit 5.11 shows how the two layers are related.

The internet layer forwards packets, hop by hop, among routers until the packet reaches the destination host. The main standard at the internet layer is the Internet Protocol (IP).

5 · 14 DATA COMMUNICATIONS AND INFORMATION SECURITY

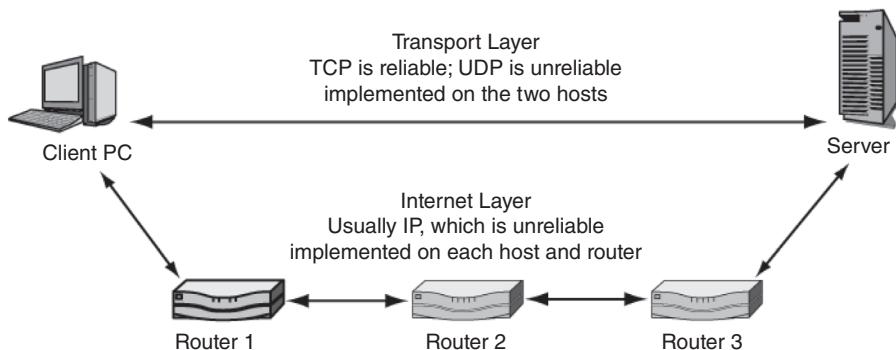


EXHIBIT 5.11 Internet- and Transport-Layer Standards

The designers of TCP/IP realized that they could not predict what services the single networks connecting routers would provide. IP was made a simple best-effort protocol, in order to assume minimal functionality in the single networks along the way. There are no guarantees that packets will arrive at all or, if they do arrive, that they will arrive in order.

To make up for the limitations of IP, a transport layer was added. The main standard designed for this layer, the *Transmission Control Protocol* (TCP), was created as a high-capability protocol that would fix any errors made along the way, ensure that packets arrived in order, slow transmission when the network became overloaded, and do several other things. For applications that did not need this level of reliability, a simpler standard was created, the *User Datagram Protocol* (UDP).

5.5 INTERNET PROTOCOL (IP). The Internet Protocol (IP) does two main things. First, it governs how packets are organized. Second, it determines how routers along the way move packets to the destination host. (Analogously, data-link-layer standards govern how frames containing packets are organized and how switches along the way move the frame across a single-switched network.)

5.5.1 IP Version 4 Packet. The main version of the Internet Protocol is Version 4 (IPv4). (There were no Versions 0, 1, 2, or 3.) This version has been in use since its definition in 1981 and will continue to be used for many years to come, although IPv6 is intended to supersede it. Exhibit 5.12 shows the IPv4 packet's organization.

A packet is a long stream of 1s and 0s. The IP *header* normally is shown on several rows, with 32 bits on each row. The first row has bits 0 through 31; the next row shows bits 32 through 63 and so on.

The header is divided into smaller units called *fields*. Fields are defined by their bit position in the packet. For example, the first four bits comprise the *version number* field. These are bits 0 through 3. In IPv4, this field holds 0100, which is 4 in binary. The *header length* field comprises the next four bits (bits 3 through 7).

5.5.1.1 First Row. As just noted, the first field (bits 0 through 3) is the version number field. In IPv4, the value is 0100 (4). In the newer version of the Internet Protocol, *IP Version 6* (IPv6), the value is 0110.

The next field is the header length field. This gives the length of the headers in 32-bit units. As Exhibit 5.12 shows, a header without options has five 32-bit lines, so this field will have the value 0101 (5 in binary).

INTERNET PROTOCOL (IP) 5 · 15

Bit 0		Bit 31	
Version (4 bits) Value is 4 (0100)	Header Length (4 bits)	Diff-Serv (8 bits)	Total Length (16 bits) length in octets
Identification (16 bits) Unique value in each original IP packet		Flags (3 bits)	Fragment Offset (13 bits) Octets from start of original IP fragment's data field
Time to Live (8 bits)	Protocol (8 bits) 1 = ICMP, 6 = TCP, 17 = UDP	Header Checksum (16 bits)	
Source IP Address (32 bits)			Destination IP Address (32 bits)
Options (if any)		Padding	
Data Field			

EXHIBIT 5.12 Internet Protocol (IP) Packet

The use of options is uncommon in practice. In fact, options tend to indicate attacks. Therefore, a value larger than 5 in the header length field indicates that the packet header has options and is therefore suspicious.

The 1-octet *dif-serv* (differential services) field was created to allow different services (priority, etc.) to be given to this packet. However, this field typically is not used.

The *total length* field gives the length of the entire IP packet in octets (bytes). Given the 16-bit length of this field, the maximum number of octets in the IP packet is 65,536 (216). Most IP packets, however, are far smaller. The length of the data field is this total length minus the length of the header in octets.

5.5.1.2 Second Row. If an IP packet is too long for a single network along the way, the router sending the packet into that network will fragment the packet, dividing its contents into a number of smaller packets. For assembly on the destination host, all fragment packets are given the same *identification field* value as in the original packet. The data octets in the original packets are numbered, and the number of the first data octet in the packet is given a *fragment offset* value (13 bits long). There are three *flag* fields (1-bit fields). One of these, *more fragments*, is set to 1 in all but the last packet, in which it is made 0. The information in these three fields allows the destination host to place the packets in order and know when there are no more packets to arrive.

IP fragmentation by routers is usually rare, and attackers can use fragmentation to hide attack information. Even if the first fragment packet is dropped by the firewall, other packets that do not have the signature information in the first header can get through. Therefore, IP fragmentation is suspicious.

5.5.1.3 Third Row. The third line begins with an ominous-sounding *time to live* (TTL) field, which has a value between 0 and 255. The sending host sets the initial value (64 or 128 in most operating systems). Each router along the way decreases the

5 · 16 DATA COMMUNICATIONS AND INFORMATION SECURITY

value by 1. If a router decreases the value to 0, it discards the packet. This process was created to prevent misaddressed packets from circulating endlessly around the Internet.

To identify hosts, attackers will use the Internet Control Message Protocol (ICMP) to *ping* many IP addresses (as discussed in Section 5.8.1). A reply tells the attacker that a host exists with that IP address. In addition, by guessing the initial TTL value and looking at the TTL value in the arriving packet, the attacker can guess how many router hops separate the attacker's host from the victim host. Sending many pings to different IP addresses can help the attacker map the routers in the target network. Often, administrators will turn off ICMP traffic outside of the internal networks in order to prevent anyone who isn't an authorized user from mapping active internal hosts.

The *data* field of the IP packet may contain a TCP message segment, a UDP datagram message, or something else, such as the ICMP messages we will discuss in Section 5.8.1. A value of 1 in this field indicates that the data field is an ICMP message. In turn, 6 indicates a TCP segment, and 17 indicates that the data field contains a UDP header.

The *header checksum* field contains a value placed there by the sender. This number is determined by a calculation based on the values of other fields. The receiving internet process redoing the calculation. If the two numbers are different, then there must have been an error along the way. If so, the router or destination host receiving the packet will simply discard the packet. There is no retransmission, so *IP* is not inherently reliable; however, one of the functions of *TCP* is to monitor the sequence numbers and initiate retransmission of missing packets. See Section 5.6.2.

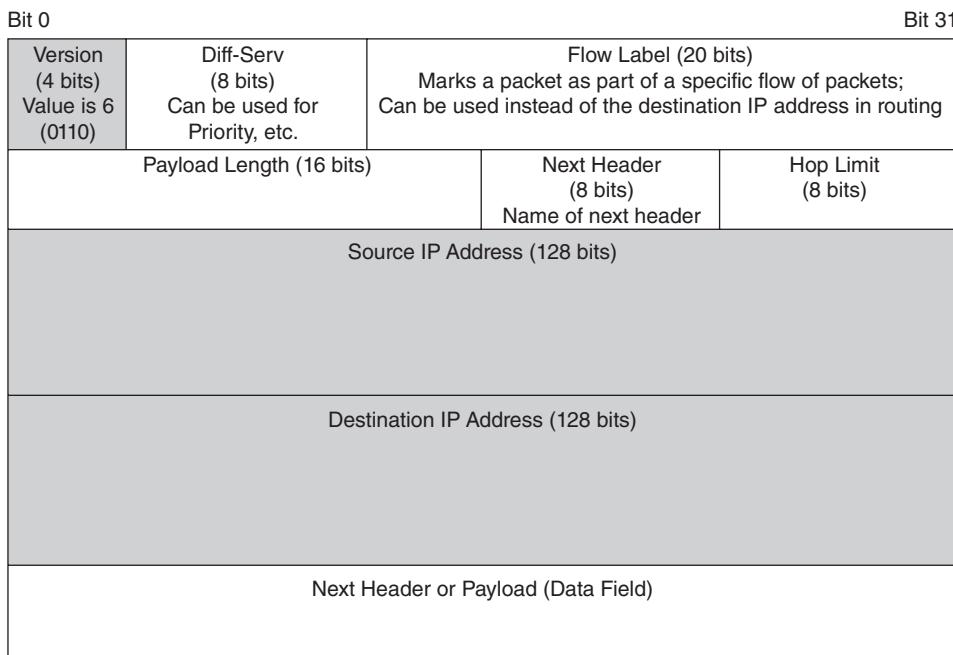
5.5.1.4 Source and Destination IP Address. When you send a letter, the envelope has an address and a return address. The analogous addresses in IP headers are the source and destination IP addresses. Note that IP addresses are 32 bits long. For human reading, these 32 bits are divided into four 8-bit *segments*, and each segment's bits are converted into a decimal number between 0 and 255. The four segment numbers are then separated by dots. An example is 128.171.17.13. Note that this dotted decimal notation is a memory and writing crutch for inferior biological entities (people). Computers and routers work with 32-bit IP addresses directly.

Many forms of firewall filtering are based on IP addresses. In addition, many attackers spoof their packet's source IP address (i.e., replace the real IP address with a false IP address).

5.5.2 IP Version 6. Although IP Version 4 is widely used, its 32-bit IP address size causes problems: It can address only 4,294,967,296 ($\sim 10^9$) devices. This relatively small size limits the number of possible IP addresses. In addition, when IP addresses were distributed, most addresses were assigned to the United States because the Internet was invented there. In fact, some U.S. universities received more IP addresses than China.

To address the limitations of the 32-bit IP address size, a new version of the Internet Protocol was created. This is *IP Version 6* (IPv6). (A Version 5 was defined, but it was never used.) Exhibit 5.13 shows the IPv6 packet organization.

One obvious change is that the IP addresses are much larger—128 bits. Each IP address, then, requires four 32-bit lines to write and is equivalent to $\sim 10^{38}$. This will provide IP addresses to allow almost every device to be a host on the Internet—including toasters and coffeepots. To give us a sense of the scale of this enormous number, it is enough to address every single molecule of water in a cube over 2 km on a side. Another popular description of the difference in size of the IPv4 and IPv6 address space is that if the address space of IPv4 were represented as a square roughly 4 cm on a side,

INTERNET PROTOCOL (IP) 5 · 17**EXHIBIT 5.13** IP Version 6 Packet

the equivalent area for IPv6 address space would cover the solar system out to Pluto's orbit.

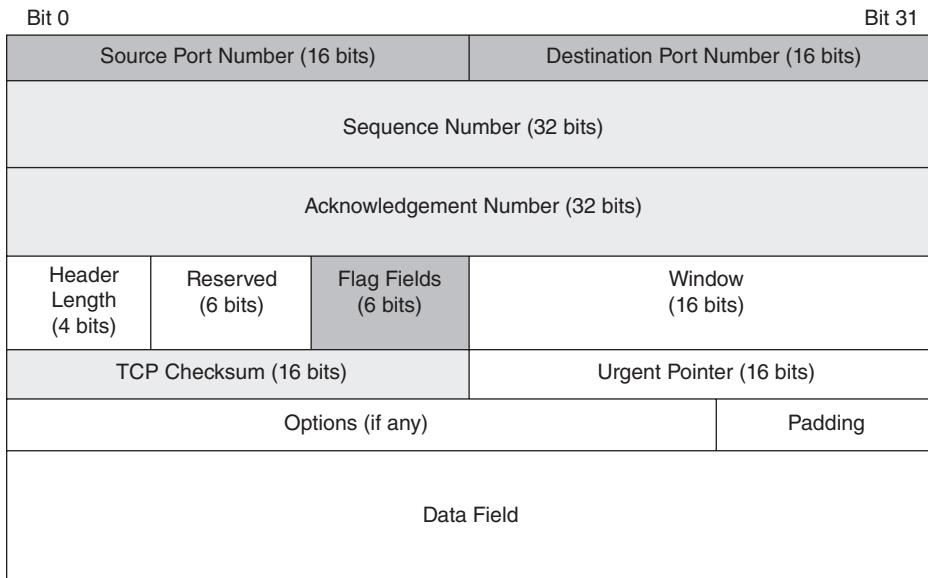
The version number field is 4 bits long, and its value is 6 (0110). There also is a *diff-serv* field and a *flow label* field that is 20 bits long. These fields allow the packet to be assigned to a category of packets with similar needs. All packets in this category would be assigned the same flow label and would be treated the same way by routers. However, this capability is not widely used.

There is a *hop limit* field that serves the same function as the time to live (TTL) field in IPv4. The *payload length*, in turn, gives the length of the data field in octets.

A major innovation in IPv6 is the *next header* field. There can be multiple headers following the first header shown in Exhibit 5.13. For instance, IPsec security is implemented with a security header. Although options are unusual in IPv4, IPv6 uses additional headers extensively. The next header field tells what the next header is. Each additional header has a *next header* field that identifies the next header or says that there is no next header.

5.5.3 IPsec. IP, which was created in the early 1980s, initially had no security at all. Finally, in the 1990s, the Internet Engineering Task Force developed a general way to secure IP transmission. This was IP security, which normally is just called IPsec. IPsec functions by protecting a packet or most of a packet and sending the protected packet inside another packet. IPsec is a general security solution because everything within the data field of the protected packet is securely encrypted, including the transport and application layer information. This includes the transport message and the application message contained in the transport message. Originally developed for IPv6, it was extended to IPv4 as well, becoming a completely general solution. See Chapter 32 in this *Handbook* for further discussion of IPsec.

5 · 18 DATA COMMUNICATIONS AND INFORMATION SECURITY



Flag fields are 1-bit fields. They include SYN, ACK, FIN, RST, PSH, and URG.

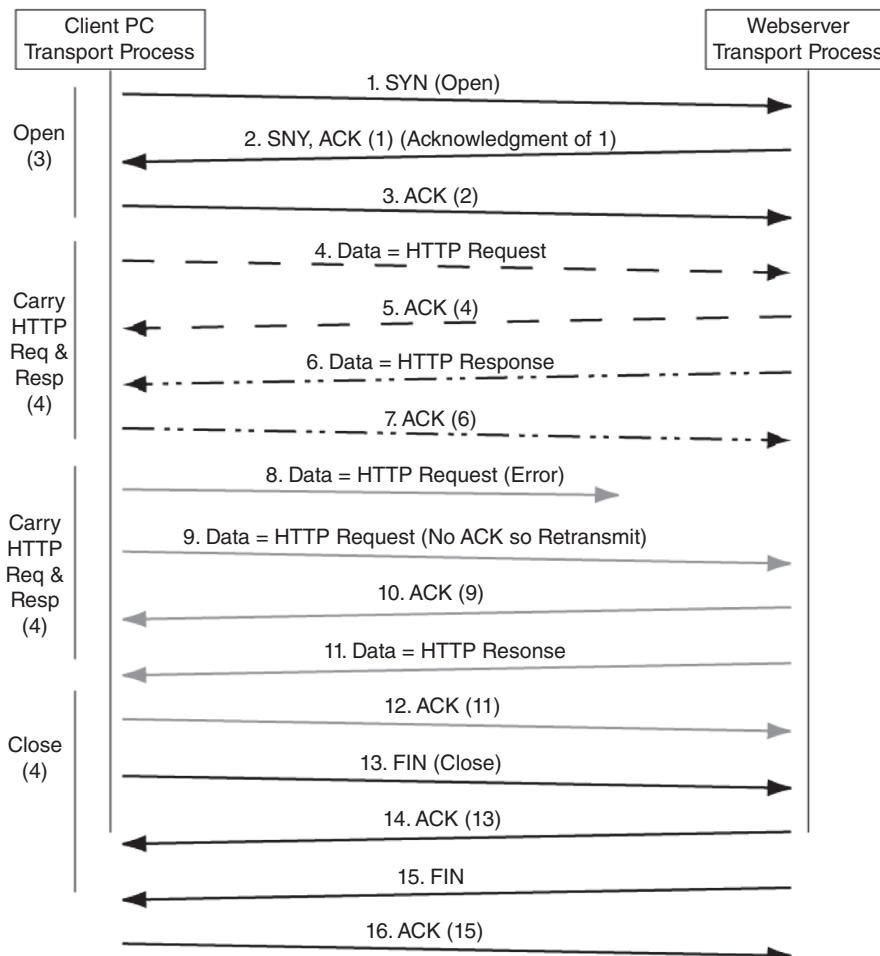
EXHIBIT 5.14 Transmission Control Protocol (TCP) Segment

5.6 TRANSMISSION CONTROL PROTOCOL (TCP). As noted earlier, the Transmission Control Protocol (TCP) is one of the two possible TCP/IP protocols at the transport layer. Exhibit 5.14 shows the TCP message, which is called a *TCP segment*.

5.6.1 Connection-Oriented and Reliable Protocol. Protocols are either connectionless or connection-oriented.

- Connection-oriented protocols are like telephone conversations. When you call someone, there is at least tacit agreement at the beginning of the conversation that you are able to speak. Explicit indicators such as “Hold, please.” and “Can I call you back?” indicate an unwillingness to proceed at the moment. Also, there is at least tacit agreement that you are done talking at the end of the conversation; simply hanging up is considered rude. “Bye” or “Talk to you later” are examples of termination signals.
- Connectionless protocols, in turn, are like email. When you send a message, there is no prior agreement, and after the message is sent, there is no built-in provision for a reply (unless you are one of those people who asks to be notified when the receiver reads the message).

Exhibit 5.15 shows a sample TCP connection. Three messages are sent to open a connection. The originator sends a TCP SYN segment to indicate that it wishes to open a TCP session. The other transport process sends back a TCP SYN/ACK segment that acknowledges the connection opening message and indicates that it is willing to open the connection. The originator then sends an ACK segment to indicate reception of the SYN/ACK segment.

TRANSMISSION CONTROL PROTOCOL (TCP) 5 · 19**EXHIBIT 5.15** Messages in a TCP Session

Attackers can use TCP connection openings to execute denial-of-service attacks that make a server unable to respond to legitimate traffic. The attacker sends a *SYN* segment to open a connection to the victim server. The victim server responds with a *SYN/ACK* message. The victim server also sets aside resources for the connection. The attacker never responds with an *ACK*, so this is called a *half-open SYN attack*. If the attacker floods a server host with *SYN* segments, the victim server will reserve so many resources that it will be overloaded and unable to serve legitimate connection opening attempts. The server may even crash. See Chapter 18 for discussion of denial-of-service attacks.

Ending a conversation, in contrast, normally takes four messages. One side sends a *FIN* segment, which the other party acknowledges. Then the other party sends a *FIN* segment, which the other side acknowledges. After the first side sends the original *FIN* segment, it will not send any new information, but it will send acknowledgments for segments sent by the other party.

5 · 20 DATA COMMUNICATIONS AND INFORMATION SECURITY

There is another way to end a session or even to reject opening one. At any point, either party can send a *RST* (reset) message. An RST message ends the conversation abruptly. There is not even an acknowledgment. It is like hanging up in a telephone conversation.

Attackers often preface an attack by attempting to identify the IP addresses of running hosts—much like thieves casing a neighborhood. One way to do this is to send TCP SYN segments to hosts. If hosts reject the SYN segment, they often send back an RST message. As noted earlier, TCP segments are carried in the data fields of IP packets. The source IP address in the packet delivering the TCP RST segment will be that of the internal host. Whenever the attacker receives an RST segment, this verifies the existence of a working host at that packet’s IP address. Firewalls often stop RST segments from leaving a site to prevent them from reaching the attacker.

5.6.2 Reliability. In addition to being connectionless or connection-oriented, protocols are either *reliable* or *unreliable*. An unreliable protocol does not detect and correct errors. Some unreliable protocols do not even check for errors. Others check for errors but simply discard a message if they find that it contains an error.

TCP is a reliable protocol. It actually corrects errors. The TCP *checksum* field is calculated using values from other fields. The sender places the result of its calculation in the checksum field. The receiver redoing the calculation and compares it with the transmitted value. If the receiving transport layer process finds that a message is correct (the values are the same), it sends an acknowledgment message. However, if the receiver detects an error in the TCP segment it receives (the values are different), it discards the segment and does nothing else.

How does a receiver know that there is an error in the message? The sender computes a value based on the other bits in the TCP segment (not just the header). The receiver redoing the calculation. If the two values match, the receiver sends an acknowledgment. If they do not match, the receiver merely drops the segment and does not send an acknowledgment.

If the segment arrives correctly, the original sender receives an acknowledgment. However, if the segment never arrives or is discarded because of damage, no reply is sent. If the original sender does not receive an acknowledgment in a specified period of time, it will resend the original segment. It will even use the original sequence number.

5.6.3 Flag Fields. *Flag* field is a general name for a 1-bit field that is logical (true or false). To say that a flag field is *set* means that its value is 1. To say that a flag field is *not set* means that its value is 0.

The TCP header contains a number of flag fields. One of these is *SYN*. To request a connection opening, the sender sets the SYN bit. The other sends a *SYN/ACK* segment, in which both the SYN and ACK bits are set. Other commonly used flags are *FIN*, *RST*, *URG*, and *PSH*.

The URG flag indicates the presence of urgent data that should be handled before earlier data octets. The urgent pointer field indicates the location of the urgent data.

If an application message is large, TCP will divide the application message into multiple TCP segments and send the segments individually. To help the receiving TCP process, the sending transport process may set the PSH (*push*) bit in the application message’s last segment. This tells the receiving transport process to push the data up to the application program immediately without buffering and delays.

TRANSMISSION CONTROL PROTOCOL (TCP) 5 · 21

5.6.4 Octets and Sequence Number. The *sequence number* field value allows the receiver to put arriving TCP segments in order even if the packets carrying them arrive out of order (including when a segment is retransmitted). Sequence numbers are also used in acknowledgments, albeit indirectly. In TCP transmission, every octet that is sent, from the very first, is counted. This octet counting is used to select each segment's sequence number.

- For the first segment, a random initial sequence number (*ISN*) is placed in the sequence number field.
- If the segment contains data, the number of the first octet contained in the data field is used as the segment's sequence number.
- For a purely supervisory message that carries no data, such as an ACK, SYN, SYN/ACK, FIN, or RST segment, the sequence number is increased by 1 over the previous message.

One dangerous attack is TCP *session hijacking*, in which an attacker takes over the role of one side. This allows the hijacker to read messages and send false messages to the other side. To accomplish session hijacking, the attacker must be able to predict sequence numbers because if a segment arrives with an inappropriate sequence number, the receiver will reject it. TCP session hijacking is likely to be successful only if the initial sequence number is predictable. Few operating systems today pick initial sequence numbers in a predictable way, but predictable sequence numbers were common in earlier operating systems, some of which are still in use.

5.6.5 Acknowledgment Numbers. When a receiver sends an acknowledgment, it sets the ACK bit. It also puts a value in the *acknowledgment number* field to indicate which segment is being acknowledged. This process is needed because the sender sends many segments and because acknowledgments may be delayed.

You might think that the acknowledgment number would be the sequence number of the segment being acknowledged. Instead, it is the number of the last octet in the data field plus 1. In other words, the acknowledgment number gives the octet number of the first octet in the next segment to be sent. This seems a bit odd, but it makes certain calculations easier for the receiver.

5.6.6 Window Field. Flow control limits the rate at which a side sends TCP segments. The TCP *window* field allows one to limit how many more octets the other side may send before getting another acknowledgment. The process is somewhat complex and has no known security implications at the time of this writing (June 2013). In acknowledgments, the ACK bit is set, and both the acknowledgment and window size fields are filled in.

5.6.7 Options. Like the IPv4 header, the TCP header can have options. However, while IP options are rare and cause for suspicion, TCP uses options extensively. One common option, often sent with the initial SYN or SYN/ACK segment, is the *maximum segment size* (MSS) option. This gives the other side a limit on the maximum size of TCP segment data fields (not on segment sizes as a whole). The presence of TCP options, then, is not suspicious by itself.

5 · 22 DATA COMMUNICATIONS AND INFORMATION SECURITY

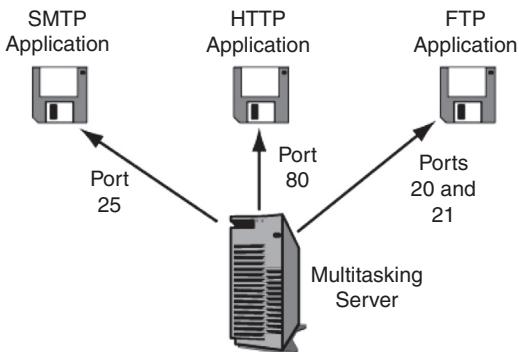


EXHIBIT 5.16 Multitasking Server Host and Port Numbers

5.6.8 Port Numbers. We have now looked at most fields in the TCP header. The first two fields warrant special mention.

5.6.8.1 Port Numbers on Servers. *Port number* fields mean different things for clients and servers. For a server, it represents a specific application running on that server, as Exhibit 5.16 shows. Servers are multitasking computers, which means that they can run multiple applications at the same time. Each application is specified by a different port number.

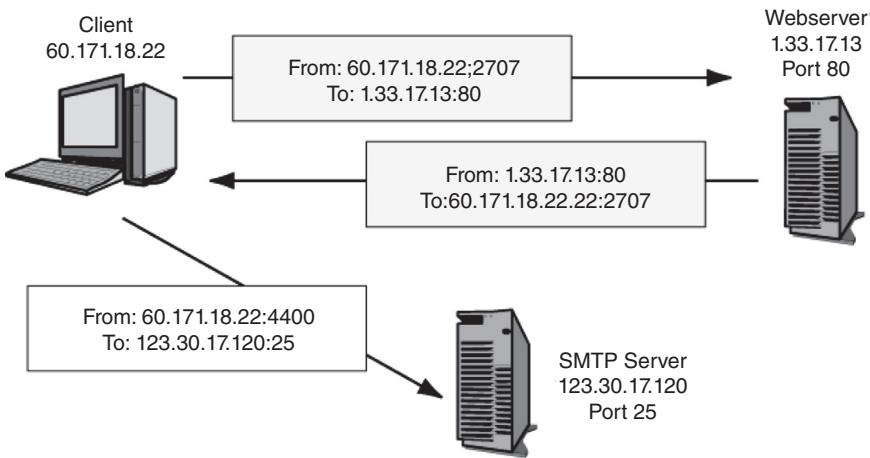
For instance, on a server, a Web server application program may run on TCP Port 80. Incoming TCP segments that have 80 as their destination port number are passed to the Web server application. Actually, TCP Port 80 is the well-known port for Web server programs, meaning that it is the usual port number for the application. Although Web servers can be given other TCP port numbers, this makes it impossible for users to establish connections unless they know or can guess the nonstandard TCP port number.

The TCP port range from 0 to 1023 is reserved for the well-known port numbers of major applications, such as HTTP and email. For instance, *Simple Mail Transfer Protocol* (SMTP) mail server programs usually are run on TCP Port 25, while *File Transfer Protocol* (FTP) requires two well-known port numbers—TCP Port 21 for supervisory control and TCP Port 20 for the actual transfer of files.

5.6.8.2 Port Numbers on Clients. Client hosts use TCP port numbers differently. Whenever a client connects to an application program on a server, it generates a random *ephemeral port number* that it uses only for that connection. On Windows machines, the ephemeral TCP port numbers range from 1024 to 4999.

The Microsoft port number range for ephemeral port numbers may differ from the official IETF range, with values of 5000–65534. The use of nonstandard ephemeral port numbers by Windows and some other operating systems causes problems for firewall filtering.

5.6.8.3 Sockets. Exhibit 5.17 shows that the goal of internetworking is to deliver application messages from one application on one machine to another application on another machine. On each machine, there is a TCP port number that specifies the application (or connection) and an IP address to specify a computer. A *socket* is a combination of an IP address and a TCP port number. It is written as the IP address,

USER DATAGRAM PROTOCOL 5 · 23**EXHIBIT 5.17** Sockets

a colon, and the TCP port number. A typical socket, then, would be something like *128.171.17.13:80*.

Attackers often do *socket spoofing*—both IP address spoofing and port spoofing. For instance, in TCP session hijacking, if the attacker wishes to take over the identity of a client, it must know both the client’s IP address and ephemeral port number. Of course, these fields are transmitted in the clear (without encryption) in TCP, so an attacker with a sniffer that captures and reads traffic flowing between the client and server can easily obtain this information.

5.6.9 TCP Security. Like IP, TCP was created without security. However, although IPsec has made IP secure, the IETF has not created a comparable way to secure TCP. One reason for this is IPsec’s ability to secure all transport layer traffic transparently, without modification to transport layer protocols. The IETF has made IPsec the centerpiece of its security protections and a single method to handle upper-layer security. Communicating partners that want TCP security should implement IPsec.

However, few TCP sessions are protected by IPsec. Consequently, some pairs of users employ an option in TCP, which adds an electronic signature to each TCP session. This signature proves the identity of the sender. This option, described in RFC 2385, requires the two parties to share a secret value. This option is awkward because it provides no way to share keys automatically, and it does not provide encryption or other protections. The option is used primarily in the Border Gateway Protocol (BGP). BGP is used to exchange routing information between administrative systems—say a corporate system and an internet service provider. BGP always uses one-to-one connections, the communicating parties usually know each other quite well, and the two parties have long-term relationships, which makes key exchange less burdensome and risky. Outside of BGP, however, the RFC 2385 electronic signature option does not appear to be used significantly. Even in BGP, it is widely seen as very weak security.

5.7 USER DATAGRAM PROTOCOL. As noted earlier, TCP is a protocol that makes up for the limitations of IP. TCP adds error correction, the sequencing of IP packets, flow control, and other functionality that we have not discussed.

5 · 24 DATA COMMUNICATIONS AND INFORMATION SECURITY

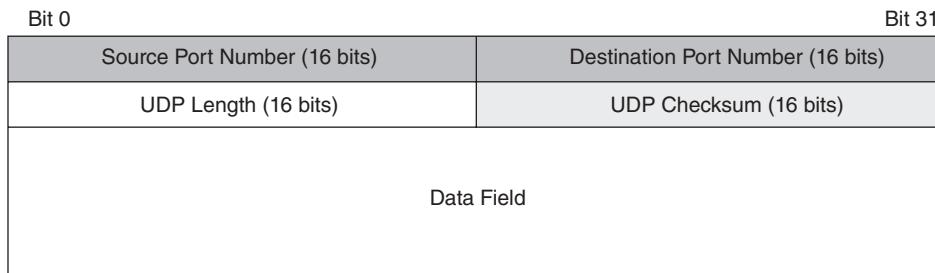


EXHIBIT 5.18 User Datagram Protocol (UDP)

Not all applications need the reliable service offered by TCP. For instance, in voice over IP (VOIP), there is no time to wait for the retransmission of lost or damaged packets carrying voice. In turn, the *Simple Network Management Protocol* (SNMP), which is used for network management communications, sends so many messages back and forth that the added traffic of connection-opening packets, acknowledgments, and other TCP supervisory segments could overload the network. Consequently, voice over IP, SNMP, and many other applications do not use TCP at the transport layer.

Instead, they use the *User Datagram Protocol* (UDP). This protocol is connectionless and unreliable. Each UDP message (called a UDP *datagram*) is sent on its own. There are no openings, closings, or acknowledgments.

As a consequence of the simplicity of UDP's operation, the UDP datagram's organization is also very simple, as Exhibit 5.18 illustrates. There are no sequence numbers, acknowledgment numbers, flag fields, or most of the other fields found in TCP.

There are source and destination port numbers, a UDP header length to allow variable-length UDP datagrams, and a UDP checksum. If the receiver detects an error using the checksum, it simply discards the message. There is no retransmission.

The fact that both TCP and UDP use port numbers means that whenever you refer to port numbers for well-known applications, you also need to refer to whether the port numbers are TCP or UDP port numbers. This is why the well-known port number for Web servers is TCP port 80.

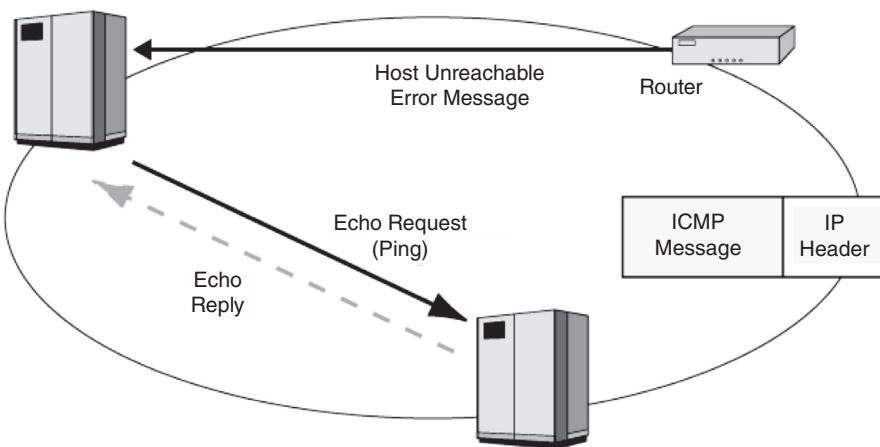
TCP's sequence numbers make TCP session hijacking very difficult. The receiver will discard messages with the wrong sequence numbers even if the source and destination sockets are correct. UDP lacks this protection, making UDP a somewhat more dangerous protocol than TCP.

Like TCP, UDP has no inherent security. Companies that wish to secure their UDP communication must use IPsec.

5.8 TCP/IP SUPERVISORY STANDARDS. So far, we have looked at standards that deliver a stream of packets across an internet and that perhaps check for errors and provide other assurances. However, the TCP/IP architecture also includes a number of supervisory protocols that keep the Internet functioning.

5.8.1 Internet Control Message Protocol (ICMP). The first supervisory protocol on the Internet was the Internet Control Message Protocol (ICMP). As Exhibit 5.19 shows, ICMP messages are delivered in the data fields of IP packets.

The best-known pair of ICMP message types is the ICMP *echo* message and the *echo reply* message. Suppose that a host sends an ICMP echo message to an IP address. If a host is active at that address, it may send back an ICMP echo reply message. This

TCP/IP SUPERVISORY STANDARDS 5 · 25**EXHIBIT 5.19** Internet Control Message Protocol (ICMP)

process is often called *pinging* because the most popular program for sending ICMP echo message is called *Ping*. The echo message is a very important tool for network management. If the network manager suspects a problem, he or she will ping a wide range of host addresses to see which of them are reachable. The pattern of responses can reveal where problems exist within a network.

Attackers also love to ping a wide range of host IP addresses. This can give them a list of hosts that are reachable for attacks. Another popular network management and attack tool is *traceroute* (or *tracert* on Windows PCs). *Traceroute* is similar to ping, but traceroute also lists the routers that lie between the sending host and the host that is the target of the traceroute command. This allows an attacker to map the network. Border firewalls often drop echo reply messages leaving the firm to the outside.

Many ICMP messages are error messages. For instance, if a router cannot deliver the packet, it may send back an ICMP error message to the source host. This error message will provide as much information as possible about the type of error that occurred.

If an attacker cannot ping destination hosts because a firewall stops them, attackers often send IP packets that are malformed and so will be rejected. The ICMP error message is delivered in an IP packet, and the source IP address in this packet will reveal the IP address of the sending router. By analyzing error messages, the attacker can learn how routers are organized in a network. This information can be very useful to attackers.

5.8.2 Domain Name System (DNS). To send a packet to another host, a source host must place the destination host's IP address in the destination address field of the packets. Often, however, the user merely types the host name of the destination host, for instance, *cnn.com*.

Unfortunately, host names are only nicknames. If the user types a host name, the computer must learn the corresponding IP address. As Exhibit 5.20 shows, the host wishing to send a packet to a target host sends a *Domain Name System* (DNS) request message to the DNS server. This message contains the host name of the target host. The DNS response message sends back the target host's IP address. To give an analogy, if you know someone's name, you must look up their telephone number in a telephone directory if you want to call them. In DNS, the human name corresponds to the

5 · 26 DATA COMMUNICATIONS AND INFORMATION SECURITY

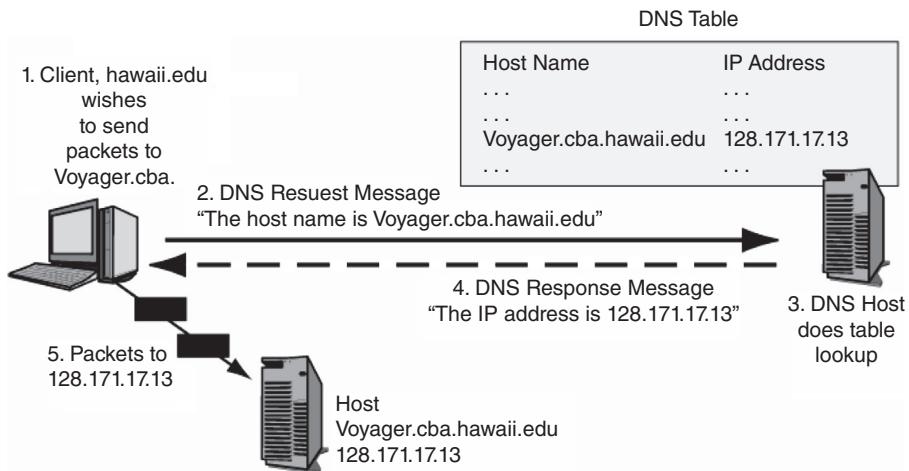


EXHIBIT 5.20 Domain Name System (DNS) Server

host name, the telephone number corresponds to the IP address, and the DNS server corresponds to the telephone directory.

DNS is critical to the Internet's operation. Unfortunately, DNS is vulnerable to several attacks. For example, in *DNS cache poisoning*, an attacker replaces the IP address of a host name with another IP address. After cache poisoning, a legitimate user who contacts a DNS server to look up the host name will be given the false IP address, sending the user to the attacker's chosen site. Denial-of-service attacks are also too easy to accomplish. RFC 3833 lists a number of DNS security issues.²

Several attempts to strengthen DNS security have been developed, under the general banner of Domain Name System Security Extensions (DNSSEC), especially RFC 2535.³ However, both the original DNSSEC specifications and the newer DNSSEC *bis* specifications (RFCs 4033-4035⁴) have proven to be insufficient. Developing a security standard that is sufficiently backwardly compatible for Internet-scale implementation has proven to be extremely difficult.

If the DNS server does not know the host name, it contacts another DNS server. The DNS system contains many DNS servers organized in a hierarchy. At the top of the hierarchy are 13 *DNS root servers*. Below these are DNS servers for *top-level domains*, such as *.com*, *.edu*, *.ie*, *.uk*, *.nl*, and *.ca*. Each top-level domain has two or more top-level DNS servers for their domain. Second-level domain names are given to organizations (e.g., *Hawaii.edu* and *Microsoft.com*). Organizations are required to maintain DNS servers for computers within their domain.

If attackers could bring down the 13 root servers, they could paralyze the Internet. Widespread paralysis would not occur immediately, but in a few days, the Internet would begin experiencing serious outages.

5.8.3 Dynamic Host Configuration Protocol (DHCP). Server hosts are given static (permanent) IP addresses. Client PCs, however, are given dynamic (temporary) IP addresses whenever they use the Internet. The *Dynamic Host Configuration Protocol* (DHCP) standard that we saw earlier in the chapter makes this possible. A DHCP server has a database of available IP addresses. When a client requests an IP address, the DHCP server picks one from the database and sends it to the client. The next time the client uses the Internet, the DHCP server may give it a different IP address.

TCP/IP SUPERVISORY STANDARDS 5 · 27

The fact that clients may receive different IP addresses each time they get on the Internet causes problems for *peer-to-peer* (P2P) applications. A presence server or some other mechanism must be used to find the other party's IP address. A lack of accepted standards for presence (including presence security) is a serious issue now that P2P applications are widespread. In fact, most security considerations in P2P presence servers have been used in P2P piracy applications, with an eye toward avoiding discovery by legitimate authorities.

5.8.4 Dynamic Routing Protocols. How do routers on the Internet learn what to do with packets addressed to various IP addresses? They frequently talk to one another, exchanging information about the organization of the Internet. These exchanges must occur frequently because the structure of the Internet changes frequently as routers are added or dropped. Protocols for exchanging organization information are called *dynamic routing protocols*. There are many dynamic routing protocols, including the *Routing Information Protocol* (RIP), *Open Shortest Path First* (OSPF), the *Border Gateway Protocol* (BGP), and Cisco Systems' proprietary *Enhanced Interior Gateway Routing Protocol* (EIGRP). Each is used under different circumstances. These protocols have widely different security features, and different versions of each protocol have different levels of functionality.

An attacker who can impersonate a router can send false dynamic routing protocol messages to other routers. These false messages could cause the routers to fail to deliver their packets. The attacker could even cause packets to pass through the attacker's computer (called a *man-in-the-middle attack* or MIMA) in order to read their contents.

The protocols just listed have widely different security features, and different versions of each protocol have different levels of security functionality.

5.8.5 Simple Network Management Protocol (SNMP). Networks often have many elements—routers, switches, and host computers. Managing dozens, hundreds, or thousands of devices can be nearly impossible. To make management easier, the IETF developed the *Simple Network Management Protocol* (SNMP). As Exhibit 5.21 shows, the manager program can send SNMP messages to managed devices to determine their conditions. The manager program can even send configuration messages that can change the ways in which remote devices operate. This allows the manager to fix many problems remotely.

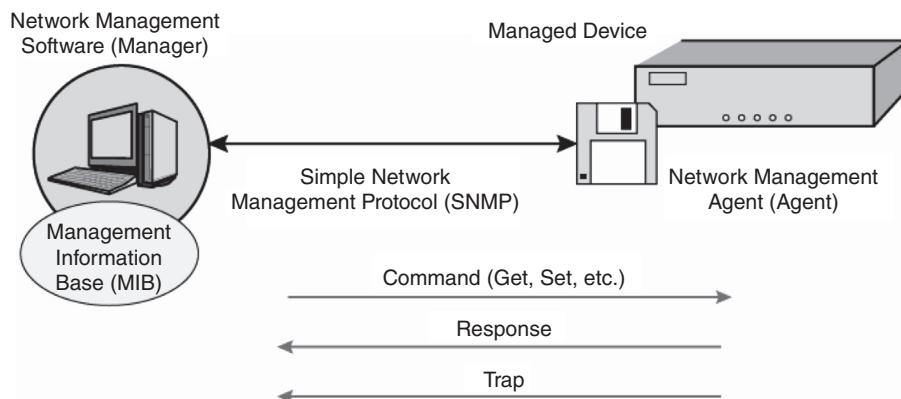


EXHIBIT 5.21 Simple Network Management Protocol (SNMP)

5 · 28 DATA COMMUNICATIONS AND INFORMATION SECURITY

Many firms disable remote configuration because of the damage that attackers could do with it. They could simply turn off all ports on switches and routers, or they could do more subtle damage.

5.9 APPLICATION STANDARDS. Most applications have their own application-layer standards. In fact, given the large number of applications in the world, there are literally hundreds of application-layer standards.

As corporations get better at defending against attacks at lower layers, attackers have begun to focus their attention on application vulnerabilities. If an attacker can take over an application running with high privileges, he or she obtains these privileges. Many applications run at the highest privileges, and attackers that compromise them own the box.

5.9.1 HTTP and HTML. Many applications have two types of standards. The *transport standard* transfers application-layer messages between applications on different machines; for the World Wide Web, this is the *Hypertext Transfer Protocol* (HTTP). The other is a standard for *document structure*. The main document-structure standard for the WWW is the *Hypertext Markup Language* (HTML).

Netscape, which created the first widely used browser, also created a security standard to protect HTTP communication. This was *Secure Sockets Layer* (SSL). Later, the Internet Engineering Task Force took over SSL and changed the name of the standard to *Transport Layer Security* (TLS).

5.9.2 E-Mail. Popular transfer standards for email are the *Simple Mail Transfer Protocol* (SMTP), *Post Office Protocol* (POP), and *Internet Message Access Protocol* (IMAP) for downloading email to a client from a mailbox on a server. Popular document-body standards include RFC 2822 (for all-text messages), HTML, and Multipurpose Internet Mail Extensions (MIME). S/MIME (Secure MIME) adds public-key encryption (see Chapter 7) to MIME and is defined in RFCs 2634, 3850, and 3851.

An obvious security issue in email is content filtering. Viruses, spam, phishing messages, and other undesirable content should be filtered out before they reach users and can do damage. (For more information on spam and other low-technology attacks, see Chapter 20 in this *Handbook*; for malware and spam countermeasures see Chapters 26, 27, 31, and 41.)

Another security issue in email is securing messages flowing from the sending client to the sender's mail server, to the receiver's mail server, and to the receiving client. Fortunately, there are security standards for part or all of the message flows, including SSL/TLS and S/MIME among others. Unfortunately, the IETF has been unable to agree on a security standard.

When Web mail, which uses HTTP and HTML for email communication, is used, then SSL/TLS can work between the sender and the sender's mail server and between the receiver's mail server and the receiver. Transmission between the email servers is another issue. Of course, senders can send encrypted message bodies directly to receivers. However, this prevents filtering at firewalls. Users should be particularly careful about using Web mail via wireless connections. (See Chapters 32 and 33.)

5.9.3 Telnet, FTP, and SSH. The two earliest applications on the Internet were the *File Transfer Protocol* (FTP) and *Telnet*. FTP provides bulk file transfers between hosts. Telnet allows a user to launch a command shell (user interface) on another computer. Neither of these standards has any security. Of particular concern is that

NOTES 5 · 29

both send passwords in the clear (without encryption) during login. The newer *Secure SHell* (SSH) standard can be used in place of both FTP and Telnet while providing high security by encrypting all transferred traffic between the hosts.

5.9.4 Other Application Standards. There are many other applications and therefore application standards. These include *Voice over IP* (VoIP; see Chapter 34 in this *Handbook*), *peer-to-peer applications* (P2P; see Chapter 35), and *service-oriented architecture* (SOA) and Web service applications (see Chapters 21, 30, and 31), among many others. Most applications have serious security issues. Application security has become perhaps the most complex aspect of network security (see Chapters 38, 39, and 40).

5.10 CONCLUDING REMARKS. It is impossible to understand information security without a strong knowledge of networking. This chapter is designed to give you a working overview of networking. It is likely to be sufficient if you run into basic networking questions while reading other chapters in this *Handbook*. However, to work in security, you will need a much stronger knowledge of networking. The books and other resources cited in Section 5.11 are a good start in that direction.

5.11 FURTHER READING

- Comer, D. E. *Internetworking with TCP/IP Vol. 1: Principles, Protocols, and Architecture*, 6th ed. Addison-Wesley, 2013.
- Ferrero, A. *The Eternal Ethernet*, 2nd ed. Boston: Addison-Wesley, 1999.
- FitzGerald, J. *Business Data Communications and Networking*, 11th ed. Wiley, 2011.
- Freedman, A. *Computer Desktop Encyclopedia*. Point Pleasant, PA: Computer Language Company, 2013. Available online from www.computerlanguage.com
- Gibson, D. *Microsoft Windows Networking Essentials*. Sybex, 2011
- Hallberg, B. *Networking, A Beginner's Guide*, 5th ed. McGraw-Hill Osborne Media, 2009.
- Hummel, S. L. *Network Design Fundamentals*. CreateSpace Independent Publishing Platform, 2013.
- Kurose, J. F., and K. W. Ross. *Computer Networking: A Top-Down Approach*, 6th ed. Pearson, 2012.
- Panko, R. R. *Business Data Networks and Telecommunications*, 8th ed. Prentice-Hall, 2010.
- Panko, R. R. *Corporate Computer and Network Security*, 2nd ed. Prentice-Hall, 2009.
- Panko, R. R., and J. Panko. *Business Data Networks and Security*, 9th ed. Prentice-Hall, 2012.
- Palmer, M. *Hands-On Networking Fundamentals*, 2nd ed. Cengage Learning, 2012.
- Roberts, R. M. *Networking Fundamentals*, 2nd ed. Goodheart-Willcox, 2011.
- Spurgeon, C. *Ethernet: The Definitive Guide*. O'Reilly Media, 2000.
- Stevens, W. R. *TCP/IP Illustrated, Volume 1: The Protocols*, 2nd ed. Addison-Wesley Professional, 2011.

5.12 NOTES

1. P. Prabakaran, "Tutorial on Spread Spectrum Technology," *EE Times | Design*, May 6, 2003, www.eetimes.com/design/communications-design/4008962/Tutorial-on-Spread-Spectrum-Technology
2. www.faqs.org/rfcs/rfc3833.html
3. www.faqs.org/rfcs/rfc2535.html
4. www.dnssec.net/

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 6

LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

Gary C. Kessler

6.1 OVERVIEW	6·2	6.5.1 OSI Model versus LAN Model Architectures	6·15
6.1.1 LAN Characteristics	6·2	6.5.2 IEEE 802 Standards	6·17
6.1.2 LAN Components	6·2	6.5.3 IEEE 802.3 CSMA/CD Standard	6·19
6.1.3 LAN Technology Parameters	6·3	6.5.4 Ethernet II	6·21
6.1.4 Summary	6·3	6.5.5 IEEE 802.5 Token-Ring Standard	6·22
6.2 LAN TOPOLOGY	6·3	6.5.6 IEEE 802.2 LLC Standard	6·23
6.2.1 Network Control	6·4	6.5.7 Summary	6·24
6.2.2 Star Topology	6·4		
6.2.3 Ring Topology	6·4		
6.2.4 Bus Topology	6·6		
6.2.5 Physical versus Logical Topology	6·6	6.6 INTERCONNECTION DEVICES	6·24
		6.6.1 Hubs	6·25
		6.6.2 Switches	6·25
6.3 MEDIA	6·8	6.6.3 Bridges	6·25
6.3.1 Coaxial Cable	6·8	6.6.4 Routers	6·26
6.3.2 Twisted Pair	6·9	6.6.5 Summary	6·27
6.3.3 Optical Fiber	6·10		
6.3.4 Wireless Media	6·11	6.7 NETWORK OPERATING SYSTEMS	6·27
6.3.5 Summary	6·13		
6.4 MEDIA ACCESS CONTROL	6·13		
6.4.1 Contention	6·13	6.8 SUMMARY	6·28
6.4.2 Distributed Polling	6·14	6.9 FURTHER READING	6·30
6.5 LAN PROTOCOLS AND STANDARDS	6·15	6.10 NOTES	6·30

This chapter provides a broad overview of local area network (LAN) concepts, basic terms, standards, and technologies. These topics are important to give the information security professional a better understanding of the terms that might be used to describe a particular network implementation and its products. The chapter also is written with an eye to what information security professionals need to know; for a more complete

6 · 2 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

overview of the topic, the reader is referred to general LAN texts, such as those listed in Section 6.10.

6.1 OVERVIEW. There are a number of ways to describe a LAN, and each will provide a glimpse as to implementation and product differences as well as points of security exposures. This section introduces various terms and perspectives as a basis for the discussion in the following sections.¹

6.1.1 LAN Characteristics. One way of describing LANs is to describe the characteristics that distinguish a *local* network from other types of networks. The most common characteristics are:

- Small geographic scope (the two most distant stations may be up to 5 kilometers [km] or so apart)
- Fast speed (data rates well in excess of 1 million [mega] bits per second [Mbps] and up to 1 billion [giga] bits per second [Gbps])
- Special media (common use of coaxial cable and optical fiber, as well as twisted pair)
- Private ownership

This type of network, then, has a very different look and feel than the Internet or some other public or private wide-area networks (WANs). More people have access to the LAN infrastructure than to the infrastructure of just about any WAN. LAN users can easily “spy” on each other by sniffing packets, something that is generally very difficult on the Internet. A single user can bring the LAN to a standstill.

The corporate LAN is generally the users’ primary access to the Internet. The users on the LAN are behind the corporate firewall and router; some studies suggest that they are responsible for 80 percent of security incidents.

Other LANs include hotel networks, which are often used by criminals as handy sources of confidential information available on unprotected systems temporarily hooked up to these Internet-access services.

Often the success of these attacks is due to users’ lack of education and awareness, such as choosing poor passwords, not maintaining up-to-date virus signature files, or computers attached to the LAN without firewalls. Sometimes the attacks are more sophisticated, such as using a packet sniffer to learn another user’s password or taking steps to degrade network performance.

6.1.2 LAN Components. In general, there are four basic components required to build a LAN, providing their own vulnerabilities and exposures from a security perspective:

1. **Computers.** These are the basic devices that are connected on the network. Read “computer” very broadly; the term can include personal computers (PCs), minicomputers, mainframes, file servers, printers, plotters, mobile devices (e.g., smartphones and tablets), communications servers, and network interconnection devices. It can also include protocol analyzers.
2. **Media.** These are the physical means by which the computers are interconnected. LAN media include unshielded twisted pair (UTP), coaxial cable (coax), optical fiber, and wireless (radio) devices. The wireless media have connection points

LAN TOPOLOGY 6 · 3

throughout an area where devices can attach to the network, and every place is a potential connection point in a wireless environment.

3. **Network interface card (NIC).** This is the physical attachment from the computer to the LAN medium. Older NICs are internal cards; the only item that is actually seen is the physical attachment to the LAN, often an RJ-45 jack. An increasing number of adapters use the universal serial bus (USB) slots on modern computers. Although NICs range widely in price depending on their capabilities, intended use, and vendor, an internal 1 Gbps Ethernet NIC for a desktop personal computer (PC) could be purchased for less than \$7, a USB Ethernet adapter for about \$14, and a wireless USB adapter for less than \$13 at the time of writing in January 2013.
4. **Software.** The three components above provide physical connectivity. Software—often called a *network operating system* (NOS)—is necessary for the devices to actually take advantage of the resource sharing that the LAN can provide. The NOS can support many types of services such as file sharing, print sharing, client/server operation, communications services, and more.

While the LAN needs to be examined in a holistic fashion, each of these components at each attached node also may require examination.

6.1.3 LAN Technology Parameters. One final way of discussing the specific operation of the LAN is to describe the technology:

- **Physical topology.** The physical layout of the medium.
- **Logical topology.** The logical relationship of the LAN nodes to each other.
- **Media Access Control (MAC) Standard.** The specification describing the rules that each node follows to determine when it is its turn to transmit on the medium.
- **Use of the Logical Link Control (LLC) protocol.** Defines the frame format employed above the MAC layer, and additional services.
- **Use of higher-layer protocols.** Defines the node-to-node communicating protocols and additional higher-layer applications.

6.1.4 Summary. It does not matter how a LAN is classified or described. It is essential, however, that the LAN be understood from a variety of perspectives to be able to apply a network-security examination.

6.2 LAN TOPOLOGY. WANs typically use some sort of switched technology, such as traditional circuit switching, packet switching (e.g., X.25), or fast packet switching (e.g., frame relay or asynchronous transfer mode [ATM]). Indeed, point-to-point lines typically connect the network switches so that there is a single data transmission on the line at one time.

Historically, LANs have been broadcast networks, meaning that every LAN station hears every transmission on the medium. LAN topologies, then, have to support the broadcast nature of the network and provide full connectivity between all stations.

The *topology* of a network is used to describe two issues. The *physical topology* describes how the LAN stations are physically connected so that they can communicate with each other. The *logical topology* describes how the broadcast nature of the LAN is actually affected, and, therefore, how stations participate in the process of obtaining

6 · 4 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

permission to transmit on the medium. There are three common topologies found in LANs: star, ring, and bus.

6.2.1 Network Control. Since LANs are broadcast networks, it is imperative that only a single node be allowed to transmit at any one time. All LANs use a *distributed access-control* scheme, meaning that all nodes follow the same rules to access the network medium and no one LAN node controls the other nodes' access. In this way, LAN nodes can come online and offline without bringing the network down.

This description is *not* meant to suggest that there are no critical elements in a LAN. Indeed, if a central hub, switch, or transmitter fails, the LAN will crash. Distributed control *does* suggest, however, that all nodes (user stations) follow the same access rules, and failure of a single *node* will not bring the LAN down. The access-control scheme is defined by the MAC protocol.

6.2.2 Star Topology. In a *star topology* (see Exhibit 6.1), all devices on the LAN are interconnected through some central device. Since LANs use distributed access-control schemes, all communication is from one node to another, and the central device merely provides a pathway between pairs of devices.

Physical star topologies have a tremendous advantage over other topologies in that they greatly ease network administration, maintenance, reconfiguration, and error recovery. Disadvantages include the single point of failure.

6.2.3 Ring Topology. In a *ring topology*, the nodes are connected by a set of point-to-point links that are organized in a circle (see Exhibit 6.2). Stations connect to

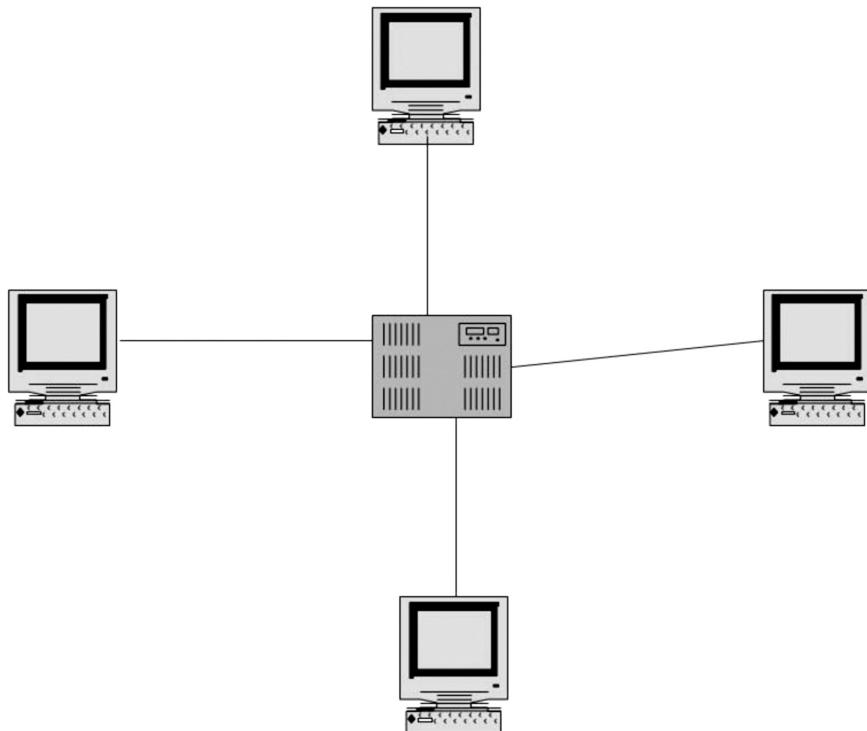
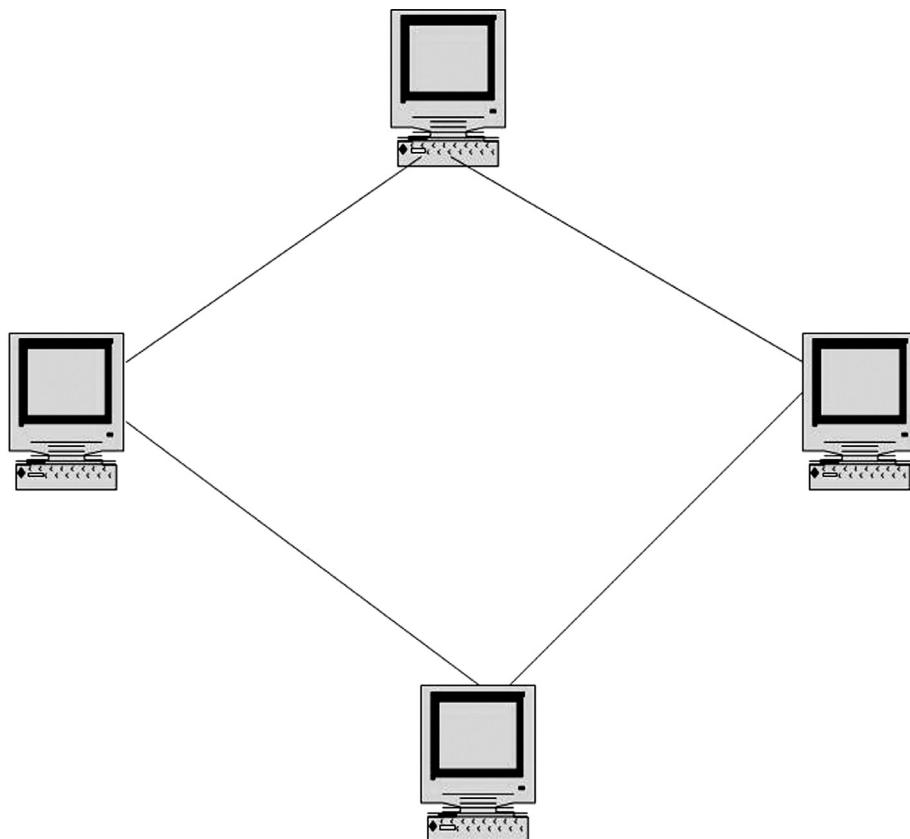


EXHIBIT 6.1 Star Topology

LAN TOPOLOGY 6 · 5**EXHIBIT 6.2** Ring Topology

the medium using *active taps* that are actually bit repeaters; a bit is read from the input line, held for a single bit time, then transmitted out to the output line.

A station transmits a message on the network by sending out a bit stream on its outgoing link; thus, rings are unidirectional in nature. Since all of the other stations see the bits one at a time, the intended receiver has no prior warning about an incoming message. For this reason, the transmitter is responsible for removing the message from the ring when the bits come back around. The MAC scheme ensures that multiple stations do not transmit at the same time.

In addition, a ring is a *serial broadcast* network. Because a station sends a message one bit at a time, every other station will see the message as it passes through but each will be receiving a different part of the message at any point in time.

Rings are a common physical LAN topology. However, unlike stars, they have multiple points of failure: if one link or one active tap fails, the integrity of the ring is destroyed. If the probability of failure of a single element is p and there are n elements, then the probability of failure $P\{F\}$ of the LAN is

$$P\{F\} = 1 - (1 - p)^n$$

As the number of elements rises, the probability of network failure rises exponentially. This problem is of such a critical nature that nearly all ring products use a star-wiring scheme or have some sort of redundancy built in for just this eventuality.

6 · 6 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

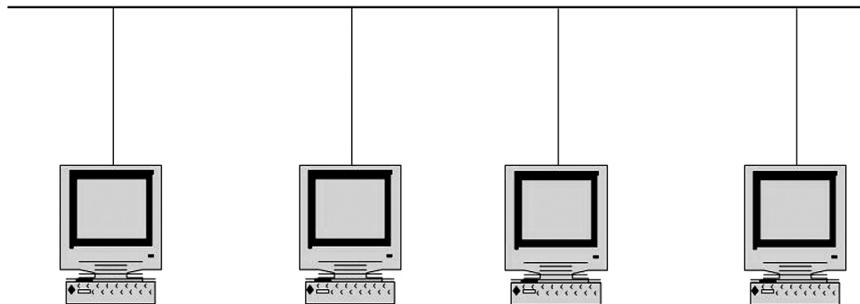


EXHIBIT 6.3 Bus Topology

6.2.4 Bus Topology. In a *bus topology* (see Exhibit 6.3), all devices are connected to a single electrically continuous medium; for this reason, this topology is also called a common cable or shared medium network. Nodes attach to the medium using a *passive tap*, one that monitors the bit flow without altering it. This is similar to the operation of a voltmeter; it measures the voltage on a power line without changing the available voltage.

Bus networks are analogous to the way appliances are connected to an alternating current (AC) power line. All of the devices draw power from the same source, even if they are on different physical segments of the power distribution network within the building. In addition, the operation of the devices is independent of each other; if the coffeepot breaks, the toaster will still work.

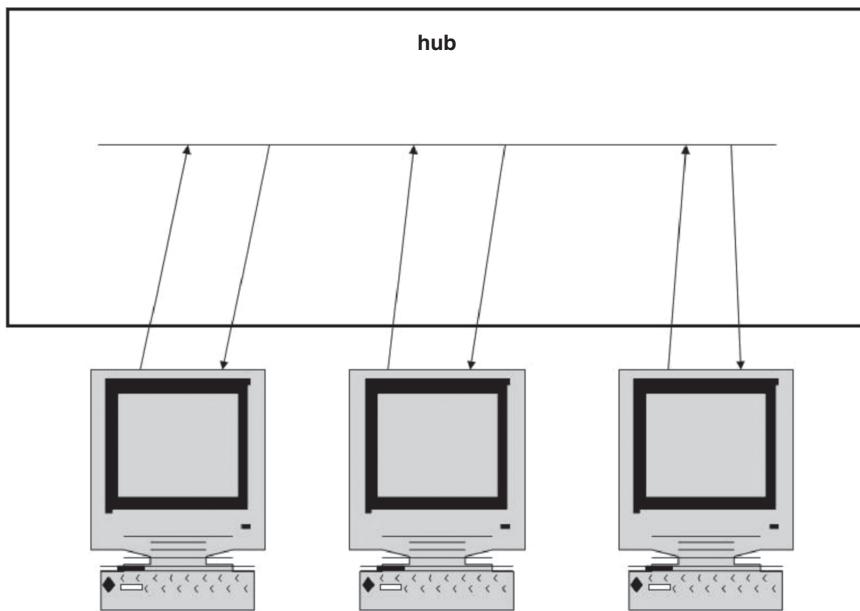
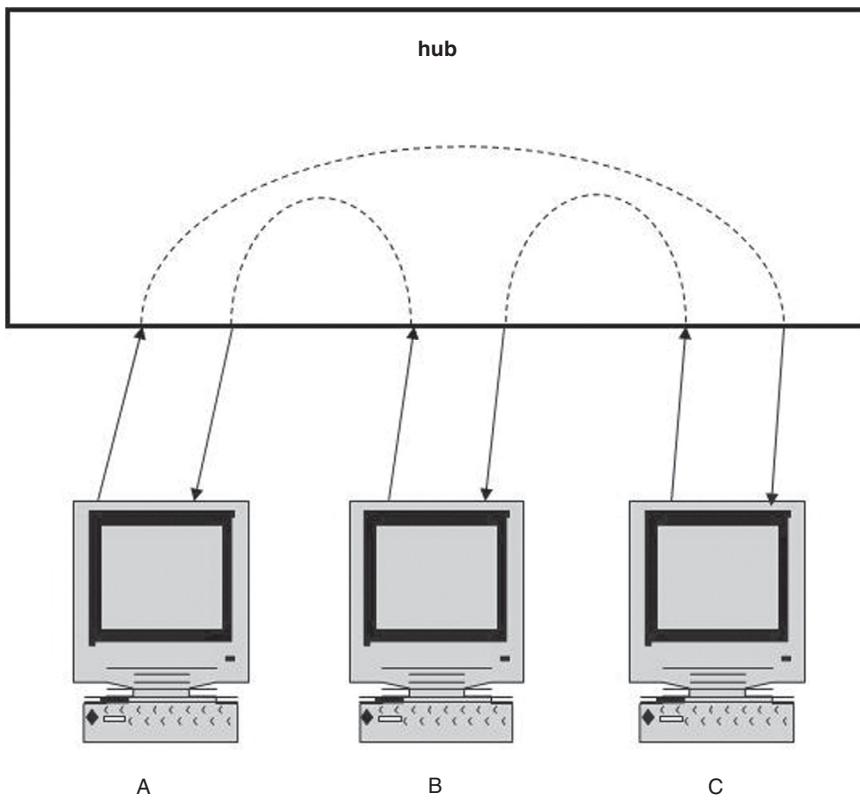
A bus is a *simultaneous broadcast* network, meaning that all stations receive a transmitted message at essentially the same time (ignoring propagation delay through the medium). Most home and business LANs employ a *baseband bus* where direct current (DC) signals are applied directly to the bus by the transmitter without any modification. In addition, transmissions on a baseband bus are broadcast bidirectionally and cannot be altered by the receivers. Bus LAN technologies are employed on cable television systems. For example, they employ a *broadband bus* where the signals are modulated (i.e., frequency shifted) to certain frequencies for transmission in one direction or another.

Buses are the oldest LAN topology and are generally limited in the type of medium that they can use. They do not usually suffer from single-point-of-failure problems.

6.2.5 Physical versus Logical Topology. A distinction was made above between the *physical* and *logical* topology of a LAN. *Physical* topology describes how the stations are physically positioned and attached to each other whereas the *logical* topology describes how the signals propagate and the logical operation of the network.

In all of today's commonly used LANs, the logical topology differs from the physical topology. The most common LAN configuration today is a *star-wired bus* (see Exhibit 6.4). This type of network has a star topology where all stations are physically attached with point-to-point links to a central device. This central device contains a bus that interconnects all of the I/O ports in such a way that when one station transmits a message, all stations will receive it. Since this acts exactly like a simultaneous broadcast, or bus, network, we categorize this configuration as a physical star, logical bus.

Another common configuration is a *star-wired ring* (see Exhibit 6.5). In this configuration, the bits will travel in logical order from station A to B, C, A, and so forth,

LAN TOPOLOGY 6 · 7**EXHIBIT 6.4** Star-Wired Bus**EXHIBIT 6.5** Star-Wired Ring

6 · 8 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

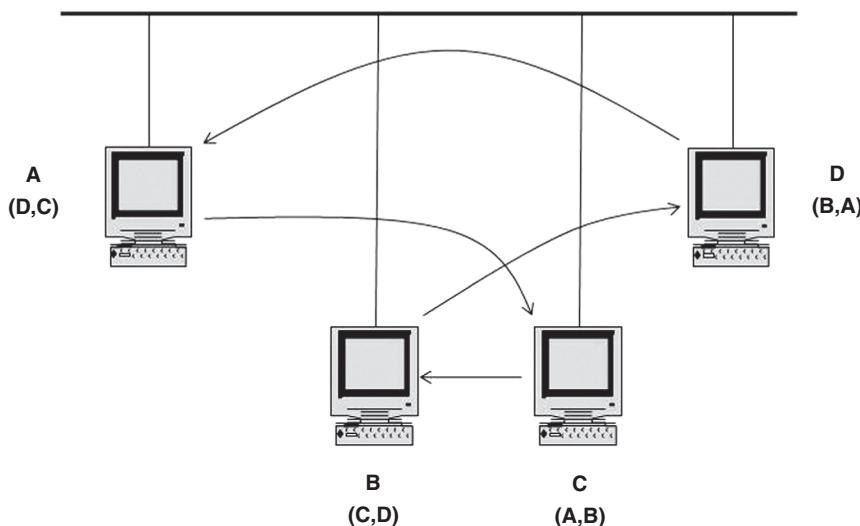


EXHIBIT 6.6 Bus-Wired Ring (The station identifier is shown above the station ID of the predecessor and successor stations in the logical ring.)

which matches the serial broadcast operation of a ring. We call this a physical star, logical ring.

Although uncommon today, another hybrid technology is the *bus-wired ring* (see Exhibit 6.6). In this configuration, nodes are passively attached to a single cable, forming a physical bus. Each station maintains a table specifying the address of predecessor and successor stations, thus forming a logical ring.

6.3 MEDIA. The next paragraphs discuss the three primary types of LAN media currently in use. Due to their relatively high speed, small geographic size, and protected environments, a number of media types can be employed with LANs.

6.3.1 Coaxial Cable. Coaxial cable (coax) is the original LAN medium. It gets its name from the physical composition of the cable itself (see Exhibit 6.7). At the center of the cable is a conductor, usually made of copper, which is surrounded by an insulator that, in turn, is surrounded by another conductor that acts as an electrical shield. Since the shield completely surrounds the central conductor and the two have a common axis, the shield prevents external electrical noise from affecting signals on



EXHIBIT 6.7 Coaxial Cable

MEDIA 6·9**EXHIBIT 6.8** Unshielded Twisted Pair

the conductor and prevents signals on the conductor from generating noise that affects other cables.

Coaxial cables vary in size from $\frac{1}{4}$ –1 inch (6.35–25.4 mm), depending on the thickness of the conductor, shield, and insulation. Applications for coax range from cable television to LANs. Speeds in excess of several hundred Mbps at distances of several hundred to several thousand meters can be achieved. Coaxial cable also has a high immunity from electromagnetic and radio frequency interference. However, it is easy to *tap*.

Coaxial cable is only seen in physical bus LANs such as Ethernet. The original Ethernet specification, in fact, called for a thin coaxial cable; a later version that employed thin (CATV) coax was dubbed Cheapernet. Coax is not typically found in star or ring networks.

Coaxial cable is easy to wiretap.²

6.3.2 Twisted Pair. The medium enjoying the largest popularity for LAN applications today is twisted pair. Twisted pair cable consists of two insulated copper conductors that are twisted around each other (see Exhibit 6.8). This is typically 22- to 26-gauge (i.e., 0.025"/0.644 mm to 0.016"/0.405 mm) wire, the same as is used for telephone wiring. Twisting the conductors around each other minimizes the effect of external electrical radiation on the signal carried on the wire; if external voltage is applied to one wire of the pair, it will be applied equally to the other wire. The twisting, then, effectively eliminates the effect of the external noise. As the number of twists per inch increases, the noise reduction characteristics improve; unfortunately, so does the overall amount of cable and the cost. Most twisted pair for telephony applications has 10 to 15 twists per foot.

The type of twisted pair cable shown in Exhibit 6.8 is called *unshielded twisted pair* (UTP) because the wire pair itself is not shielded. The data-carrying capacity of UTP is generally indicated by its *category*:

Level 1 (sometimes called Category 1, or Cat 1) cable is older 0.4 MHz cable used for some telephone and modem applications, but is generally unsuited for data applications.

Level 2 (sometimes called Category 2, or Cat 2) cable is 4 MHz cable used for legacy data terminal systems, such as IBM 3270 BISYNC.

Category 3 (Cat 3) cable has a maximum bandwidth of 16 MHz and is rated for 10 Mbps over a wire segment of 100 m (although speeds of 100 Mbps can often be achieved). Cat 3 cable is rated up to 16 Mbps and is primarily used today for telephones.

Category 4 (Cat 4) cable, with a maximum bandwidth of 20 MHz, was used for IBM's 16 Mbps Token Ring networks. It is not commonly seen today.

6 · 10 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

Category 5 (Cat 5) cable has a maximum bandwidth of 100 MHz and is rated for voice or data at speeds up to 100 Mbps over a wire segment of 100 meters. Cat 5e is rated for full-duplex and 1 gigabit (Gbps) Ethernet. These are the most common LAN cables in use today.

Category 6 (Cat 6) cable is rated to 250 MHz over a wire segment between 15 and 100 meters in length. Cat 6 is intended for use for very-high-speed broadband applications at data rates up to 10 Gbps. Cat 6a cable is a variant of Cat 6, rated to 500 MHz and 10 Gbps.

Category 7 (Cat 7) cable is rated up to 600 MHz and uses four pair. Each pair of wires in the cable sheath, and the sheath itself, are shielded to prevent electromagnetic interference at data rates up to 10 Gbps. Cat 7a is rated up to 1,000 MHz and 10 Gbps data rates.

UTP is commonly found in physical star-wired bus and ring LANs; it is never used in a physical bus and rarely in a physical ring. As with coaxial cable, it is easy to wiretap. In addition, many LANs using UTP also make connections through *patch panels*, which are frequently unprotected because the technicians installing the connections are unaware of the security issues of providing centralized access to dozens or hundreds of connections.

Another twisted pair variant is *shielded twisted pair* (STP), where each cable pair is surrounded by a metallic shield that provides the same function as the outer conductor in coaxial cable. STP was only used in the IBM Token Ring, a star-wired ring.

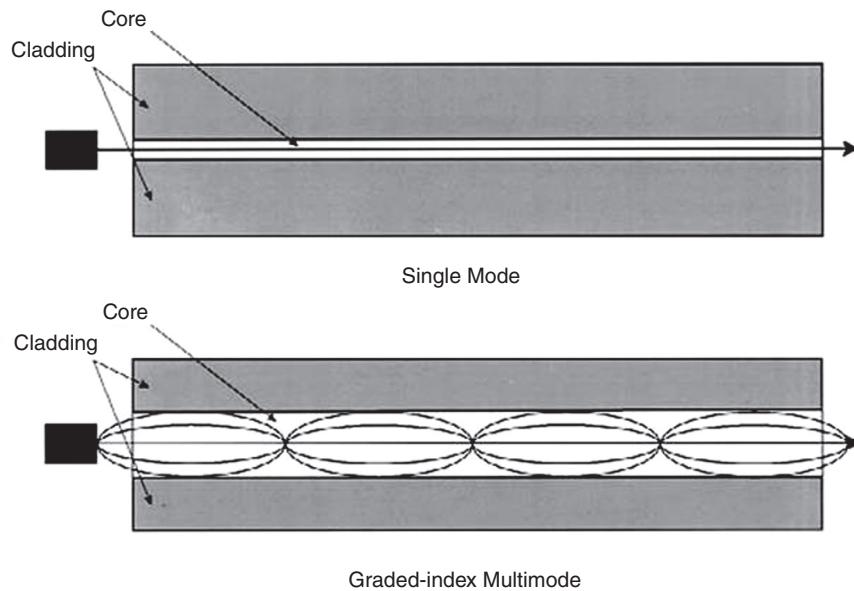
6.3.3 Optical Fiber. Optical fiber is a thin flexible medium that acts as a wave-guide for signals in the 10^{14} - to 10^{15} -Hz range, which includes the visible light spectrum and part of the infrared spectrum. Optical fiber is a great medium for digital communications; it is essentially immune to any type of radio or magnetic interference and very difficult (using highly specialized equipment) to tap surreptitiously. Theoretically able to achieve data rates on the order of trillions of bits per second, optical fiber has been shown to reach data rates of 100 Gbps over a 4,350-mile (7,000-km) fiber; the practical limit is usually due to the electronics performing optical-electrical conversion.

In WAN applications, this speed limit is exceeded in one of two ways.

1. An optical switch can terminate optical fiber without any electrical-optical conversion.
2. Dense wave division multiplexing (DWDM) allows 100 or more 10 Gbps bit streams to be carried on a single-fiber strand simultaneously. These technologies may well eventually find their way to the LAN.

The electronics are a critical part of any optical fiber system. The incoming electrical signal to be transmitted on the fiber is converted to an optical signal by the transmitter. Common optical sources are a *light-emitting diode* (LED) or *injection laser diode* (ILD). LEDs are less expensive than ILDs but are limited to lower speeds. The optical signal is received by a device called a *photodiode*, which essentially counts photons and converts the count to an electrical signal. Common photodiodes include the *positive-intrinsic-negative* (PIN) *photodiode* and *avalanche photodiode* (APD). The PIN is less expensive than the APD but is limited to lower speeds.

The physical and transmission characteristics of optical fiber are shown in Exhibit 6.9. At the center of an optical fiber cable is the *core*, a thin, flexible medium capable of carrying a light signal. The core is typically between 2 and 125 micrometers

MEDIA 6 · 11**EXHIBIT 6.9** Optical Fiber Cable

(μm), or microns, in diameter and may be made from a variety of glass or plastic compounds. Surrounding the core is a layer called the *cladding*. The optical characteristics of the cladding are always different from the core's characteristics so that light signals traveling through the core at an angle will reflect back and stay in the core. The cladding may vary in thickness from a few to several hundred microns. The outermost layer is the *jacket*. Composed of plastic or rubber, the jacket's function is to provide the cable with physical protection from moisture, handling, and other environmental factors.

Two types of optical fiber cable are used for voice and data communications, differentiated by their transmission characteristics (see Exhibit 6.9). *Multimode fiber* (MMF) has a core diameter between 50 and 125 μm . Because this diameter is relatively large, light rays at different angles will be traveling through the core. This phenomenon, known as *modal dispersion*, has the effect of limiting the bit rate and/or distance of the cable. MMF cable is generally limited to a maximum cable length of 2 km. *Single-mode fiber* (SMF) eliminates the multiple path problem of MMF by using a thin core with a diameter of 2 to 8 μm . This thin-core cable results in a single propagation path so that very high bandwidths over large distances (up to 10 km) can be achieved.

SMF is the most expensive type of fiber and is usually used for long-haul data and telecommunications networks. MMF is commonly used on LANs; it is less expensive but can still handle the required data rates and distances.

Optical cable is extremely difficult to wiretap, but it's easy to cut. Such cables should be protected by shielded conduits, not placed in accessible locations such as next to a baseboard on the floor of a public corridor—in a hospital!

6.3.4 Wireless Media. Wireless LANs use radio signals to interconnect LAN nodes. Wireless LANs are increasingly common in environments where:

- It is difficult to install new wiring (e.g., in a building with asbestos in the walls).
- There are mobile users (e.g., in a hospital or car rental agency).

6 · 12 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

- Right-of-ways for wiring are hard to obtain (e.g., campus environments that span roadways).
- A temporary network is necessary (e.g., at a conference or meeting).
- Residential areas have no other networking facilities.
- Conference centers, hotels, and colleges and universities need wide and easy network access.

Wireless LANs generally employ infrared, spread spectrum, or microwave communications technology. *Infrared* (IR) is used for a variety of communications, monitoring, and control applications. It is also used for such non-LAN applications as home entertainment remote control, building security intrusion and motion detectors, medical diagnostic equipment, and missile guidance systems. For wireless LANs, the most common IR communications band uses signals with a wavelength in the range 800 to 1,000 nanometers (nm, or 10^{-9} m). *Diffused IR* operates at data rates between 1 to 4 Mbps at distances up to 200 feet, and can be used for stationary or mobile LAN nodes. *Directed Beam IR*, which requires line-of-sight, operates at data rates from 1 to 10 Mbps at distances up to 80 feet. IR systems are limited to a single room because the signals cannot pass through walls.

Spread spectrum is a wireless communications technology in the region of 2.4 or 5 gigahertz (GHz, or billions of cycles per second), where the actual frequency of the transmitted signal is deliberately varied during transmission. Originally, the frequency shifting was for security purposes to prevent monitoring of the communications channels. Two types of spread spectrum technology are used in LANs:

1. In *frequency hopping spread spectrum* (FHSS),³ the transmitter sends the signal over a set of radio frequencies, hopping from frequency to frequency at split-second intervals in what appears to be a random sequence. The sequence is not random, however, and the receiver changes frequencies in synchronization with the transmitter. FHSS can support data rates from 1 to 3 Mbps up to a distance of 330 feet (100 m).
2. In *direct sequence spread spectrum* (DSSS), each bit in the original data stream is represented by multiple bits in the transmitted signal, spreading the signal across a wide frequency range. One result of DSSS is that the system can achieve a greater bandwidth than the original signal. DSSS can support data rates in excess of 20 Mbps up to a distance of 1,000 feet (300 m).

Another modulation scheme used in wireless LANs is *orthogonal frequency-division multiplexing* (OFDM). OFDM is a variant of frequency division multiplexing and uses a forward error-correction scheme and orthogonal subcarriers in order to minimize frequency crosstalk and bit errors.

Wireless access points can be purchased for less than \$100, making this an attractive alternative to even UTP-based LANs in many scenarios.

Microwave LANs refers to communications in the area of 1, 5, and 19 GHz. Electromagnetic energy with a frequency higher than 1 GHz and data rates up to 20 Mbps can be maintained for distances up to 130 feet. One major disadvantage of microwave is that Federal Communications Commission (FCC) licensing is required for many of these frequencies.

MEDIA ACCESS CONTROL 6 · 13

6.3.5 Summary. In the early 1980s, coaxial cable was the most commonly used LAN medium. Twisted pair, used for telephony applications, was not used in LANs because high speeds could not be achieved. Optical fiber technology was still in its infancy and was very expensive. All of this changed by the early 1990s, when the electronics to drive twisted pair had dramatically improved, and optical fiber technology had greatly matured. It is rare to see coaxial cable used in a LAN today; instead, UTP (less costly than coax) or optical fiber (higher speeds than coax) are more often employed. Wireless LANs are a viable alternative to wire-based networks, yet more difficult to secure.⁴ However, growth in wireless network access to the Internet outstripped fixed broadband subscriptions by the late 2000s; in 2011, rates of growth for wireless subscriptions were 200 to 300 percent of the rates of growth for fixed subscriptions.⁵

6.4 MEDIA ACCESS CONTROL. As mentioned, LANs are broadcast networks connecting peer devices, all having equal access to the medium. These characteristics place two requirements on the protocol that controls access to the network:

1. There can be only one station transmitting at any given time since multiple transmitters would result in garbled messages.
2. All stations must follow the same rules for accessing the network since there is no master station.

The schemes controlling access to the network medium are called *media access control* (MAC) protocols. Although many different LAN MAC schemes have been introduced in working products, the most common ones are essentially variants of two approaches: *contention* and *distributed polling*. These schemes will be discussed below, along with reference to appropriate Institute for Electronics and Electrical Engineers (IEEE) LAN standards.

6.4.1 Contention. A contention network can be compared to a group of people sitting around a conference table without a chairperson. When someone wants to speak, it is necessary first to determine whether anyone else is already speaking; if someone else is speaking, no one else can begin until that person has stopped. When a person detects silence at the table, he or she starts to talk. If two people start to talk at the same time, a *collision* has occurred and must be resolved. In the human analogy, collisions are resolved in one of two ways: Either both speakers stop and defer to each other (“polite backoff”) or both continue speaking louder and louder until one gives up (a “rudeness algorithm”).

The contention scheme used in LANs is similar to the polite backoff situation, and is called *carrier sense multiple access with collision detection* (CSMA/CD). CSMA/CD is one of the oldest LAN MAC schemes in use today, used originally in Ethernet and becoming the basis of the IEEE 802.3 standard (to be described). Although there have been other contention schemes used on LANs, CSMA/CD is the one that has survived and thrived in the marketplace.

CSMA/CD works on logical bus networks. When a station is ready to transmit, it first listens to the network medium (“carrier sense”). If the station detects a transmission on the line, it will continue to monitor the channel until it is idle. Once silence is detected, the station with a message to send will start to transmit. Stations continue to monitor the channel during transmission so that if a collision is detected, all transmitters stop transmitting.

6 · 14 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

CSMA/CD networks employ a *backoff scheme* so that the first collision does not bring the network down. Without a backoff scheme, all transmitters would detect a collision and stop transmitting; after again hearing silence on the line, however, all stations would once again start transmitting and would again collide with each other. The backoff scheme causes stations to make a random decision whether to transmit or not after silence is detected on the channel after a collision has occurred.

CSMA/CD uses a backoff scheme called *truncated binary exponential backoff*. Although this name is a mouthful, it actually describes the process very precisely. When a station is ready to transmit and detects silence on the line, it will attempt to send a message with a probability of 1 (i.e., 100 percent likelihood that it will transmit); this probability is called the *persistency* of the MAC scheme.⁶ If a collision occurs, the station will stop transmitting and again wait for silence on the line. When silence is again detected, the station will transmit with a probability of $\frac{1}{2}$ (i.e., there is a 50 percent chance that it will transmit and a 50 percent chance that it will not). If two stations were involved in the collision and they both back off to a $\frac{1}{2}$ -persistent condition, then there is a 50 percent chance that one will transmit and one will defer at the next transmission opportunity, a 25 percent chance that both will defer at the next opportunity, and a 25 percent chance that both will collide again.

If a station collides again, its persistency is again cut in half, now to $\frac{1}{4}$. All stations involved in the collision(s) drop their persistency and each station independently determines whether it will transmit at the next occurrence of silence or not.

As long as collisions occur, the persistency is continually cut in half until the station either successfully transmits or has 16 unsuccessful attempts to transmit the message. After 16 failed attempts, the station gives up.⁷ After the station successfully transmits or has 16 unsuccessful attempts, the station's persistency returns to 1 and the operation continues as before.

Wireless LANs also use a form of contention, but it is generally not CSMA/CD because collision detection is not practical in a wireless environment. Instead, the stations still employ CSMA—they listen for an idle channel—but they do not necessarily transmit when the channel is idle. Instead, they wait to see if the channel remains idle for some period of time in an attempt to stave off a collision. This is a form of *CSMA with collision avoidance* (CSMA/CA).

6.4.2 Distributed Polling. Imagine that the same group of people is sitting around the same conference table, still without a chairperson. One person at the table has a microphone and can say anything to anyone in the room. Everyone in the room, of course, will hear the message. The rule here is that the only person who is allowed to speak is the one with the microphone; furthermore, the person will hold on to the microphone only while he or she has something to say and can hold on to it only for some maximum amount of time. When the first person is done talking or the time limit is reached, the microphone is passed to the next person at the table. Person 2 can now speak or immediately pass the microphone on to person 3. Eventually, the first person at the table will get the microphone back and get another opportunity to talk.

The scheme just described is implemented in LANs with a scheme called *token passing*. This is the basis for the IBM Token Ring and represents the second most commonly used LAN MAC algorithm. Token passing, in one variant or another, is the basis for the IEEE 802.4 and 802.5 standards, as well as for the Fiber Distributed Data Interface (FDDI).

Token passing requires a logical ring topology. When a station has data to send to another station, it must wait to receive a bit pattern representing the *token*. Tokens are

LAN PROTOCOLS AND STANDARDS 6 · 15

sent in such a way that only one station will see it at any given time; in this way, if a station sees the token, it has temporary, exclusive ownership of the network.

If a station receives the token and has no data to send, it passes the token on. If it does have data to send, it generates a *frame* containing the data. After sending the frame, the station will generate and send another token.

A *token ring* network is a logical ring implemented on a physical topology that supports a serial broadcast operation (i.e., a star or a ring). Each station receives transmissions one bit at a time and regenerates the bits for the next station. A station transmitting a frame will send the bits on its output link and receive them back on its input link. The transmitter, then, is responsible for removing its message from the network. When finished transmitting, the station transfers control to another station by sending the bits comprising a token on its output link. The next station on the ring that wants to transmit *and* sees the token can then send its data frame. Token rings (standardized in 802.5 and FDDI) are the most common implementation of token passing.

A *token bus* network (as specified in 802.4) is conceptually similar to the token ring, except that it is implemented using a simultaneous broadcast topology (i.e., a bus). In this physical topology, all stations hear all transmissions. A station that wants to send data to another will address a frame to the intended receiver on the network, as in a CSMA/CD bus. When done transmitting, the station will address a token to the next station logically in the ring; while all stations will hear the token transmission, only the one station to which it is addressed will pick it up. After receiving a token, a station may or may not transmit data, but it is, in any case, responsible for passing the token to the next station in the logical ring. Eventually, the token will return to the first station.

6.5 LAN PROTOCOLS AND STANDARDS. The Open Systems Interconnection (OSI) Reference Model continues to be the standard framework with which to describe data communications architectures, including those for LANs. The basic LAN protocol architecture maps easily to the OSI model, as discussed in this section.

6.5.1 OSI Model versus LAN Model Architectures. Although the LAN protocol architecture can be related to the OSI model, there is not a perfect one-to-one mapping of the protocol layers (see Exhibit 6.10). The OSI Physical Layer is analogous to a LAN Physical Layer (PHY). Both specify such things as:

- Electrical characteristics of the interface
- Mechanical characteristics of the connector and medium
- Interface circuits and their functions
- Properties of the medium
- Signaling speed
- Signaling method

Most LAN physical layer specifications actually comprise two sublayers. The lower sublayer describes physical layer aspects that are specific to a given medium; the higher sublayer describes those aspects that are media-independent.

The OSI Data Link Layer, responsible for error-free communication between any two communicating devices, is represented by two sublayers in a LAN. The lower sublayer is the MAC, which deals with issues of how the station should access the

6 · 16 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

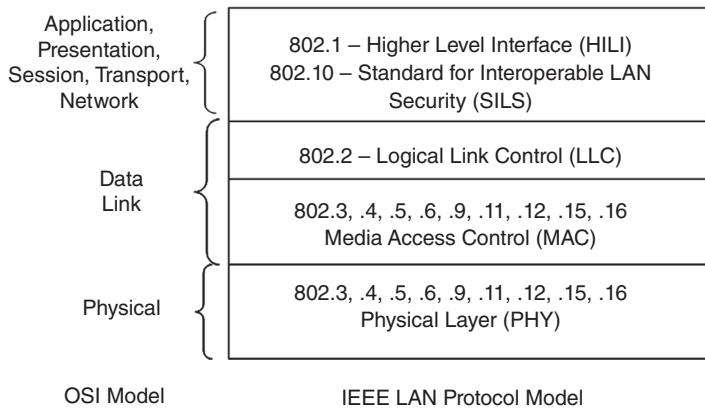


EXHIBIT 6.10 IEEE versus LAN Protocol Models

network medium. The MAC is responsible for error-free communication over the PHY and specifies such things as:

- Framing
- Addressing
- Bit-error detection
- Control and maintenance of the MAC protocol
- Rules governing medium access

The upper sublayer is called the Logical Link Control (LLC). The LLC protocol is responsible for maintaining a logical connection between two communicating LAN stations. The LLC specifies such rules as:

- Frame sequencing
- Error control
- Establishment and termination of a logical connection
- Addressing of higher layer services

Recalling that the main functions of the network layer are routing and congestion control, there are two reasons that no LAN protocol layer acts strictly like the OSI Network Layer:

1. There is no need for a routing algorithm in a broadcast network because all stations receive all transmissions; the address of the intended receiver is included in the transmission itself.
2. Congestion control is also not an issue in a broadcast network; a broadcast network must be limited to a single transmitter at a time, and this is accomplished by the MAC layer.

There are no standards for LANs corresponding to the upper four layers of the OSI model. Even in the less organized 1980s, end-to-end protocols as such were not required in a LAN environment because the end-to-end communication was limited to nodes on the LAN, and for that the MAC guaranteed error-free communication.

LAN PROTOCOLS AND STANDARDS 6 · 17

Only when LAN interconnection, via WAN and LAN access to the Internet, gained popularity did other end-to-end protocols become necessary. IP (and other network layer protocols) grew in demand as well. Those protocols are associated with the communications software as part of a network operating system (NOS), and these will be discussed later.

6.5.2 IEEE 802 Standards. Although they are not directly related to security, it is useful to be familiar with the standards describing LANs, the most common of which are the IEEE 802 standards. The IEEE Computer Society formed the Project 802 Committee in February 1980 to create standards for LANs as part of its more general work on standards for microprocessors; no other organization was making any similar standardization efforts. Originally, there was to be a single LAN standard, operating at a speed between 1 and 20 Mbps. The standard was divided into three parts: PHY, MAC, and a high-level interface (HILI) to allow other protocol suites to have a common protocol boundary with the LAN. The original MAC was based on the Ethernet standard, but other MAC schemes were quickly added and, over the years, the 802 committee has addressed many LAN schemes. They all have in common an interface to a single LLC protocol that provides a common interface between the HILI and any MAC.

A description of the Project 802 working groups (WG) and their status as of October 2012 follows.⁸

802.1—High-Layer LAN Protocols Working Group. Provides the framework for higher-layer issues, including protocol architecture, security, end-to-end protocols, bridging, internetworking, network management, and performance measurement.

802.2—Logical Link Control Working Group. Provides a consistent interface between any LAN MAC and higher-layer protocols. Depending on the options employed, the LLC can provide error detection and correction, sequential delivery, and multiprotocol encapsulation. The 802.2 standard is described in more detail in Section 6.5.6. This WG has been *disbanded*.⁹

802.3—Ethernet Working Group. Defines the MAC and PHY specifications for a CSMA/CD bus network. This specification is discussed in more detail in Section 6.5.3. (The 802.3 CSMA/CD standard is based on Ethernet, described in Section 6.5.4.)

802.4—Token Bus Working Group. Defines the MAC and PHY specifications for a token-passing bus based on work originally done at General Motors as part of the Manufacturing Automation Protocol (MAP). Well suited for factory floors and assembly lines, MAP never achieved widespread use. This WG has been disbanded.

802.5—Token Ring Working Group. Defines the MAC and PHY specifications for a token-passing ring. Although this WG is currently disbanded, the specification is discussed in more detail in Section 6.5.5 for historical purposes.

802.6—Metropolitan Area Network (MAN) Working Group. Defines the MAC and PHY specifications for a MAN. In particular, the 802.6 standard defines a MAC and PHY called Distributed Queue Dual Bus (DQDB), which was one of the MACs employed with the Switched Multimegabit Data Service (SMDS) and Connectionless Broadband Data Service (CBDS). Introduced in the early 1990s, neither service remains in use today. This WG has been disbanded.

6 · 18 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

802.7—Broadband Technology Advisory Group (BBTAG). Advises other 802 subcommittees about changes in broadband technology and their effect on the 802 standards. This WG has been disbanded.

802.8—Fiber Optics Technology Advisory Group (FOTAG). Advises other 802 subcommittees about changes in optical fiber technology and their effect on the 802 standards. This WG has been disbanded.

802.9—Integrated Services LAN (ISLAN) Working Group. Defines the MAC and PHY specifications for integrated voice/data terminal access to integrated services networks, including ISLANs and MANs, and Integrated Services Digital Networks (ISDN). The only practical implementation was deployed in IsoEthernet products, described in the IEEE 802.9a standard. This WG has been disbanded.

802.10—Security Working Group. Defines procedures for providing security mechanisms on interconnected LANs, including cryptography and certificates. This WG has been disbanded.

802.11—Wireless LAN (WLAN) Working Group. Defines MAC and PHY specifications for “through the air” media. The original 802.11 standard defined operation at 1 or 2 Mbps using the 2.4-GHz range and DSSS or FHSS spread spectrum technology; nominal maximum distances were 330 feet (100 m). The most common variants today are:

- 802.11b—Data rates up to 11 Mbps at a nominal maximum distance up to 460 feet (140 m) on a frequency of 2.4 GHz using DSSS.
- 802.11g—Data rates up to 54 Mbps at a nominal maximum distance up to 460 feet (140 m) on a frequency of 2.4 GHz using DSSS or OFDM.
- 802.11n—Data rates up to 150 Mbps at a nominal maximum distance up to 820 feet (250 m) on a frequency of 2.4 or 5 GHz using OFDM.

Future 802.11 standards are expected that will provide data rates up to 866.7 Mbps on a frequency of 5 GHz using OFDM.

802.12—Demand Priority Working Group. Describes one of the MAC and PHY specifications originally proposed for 100 Mbps LAN speeds and dubbed *100BASE-VG/AnyLAN*. Largely unused, and the WG has been disbanded.

802.13. (This number was never assigned to a WG because it was felt that the *13* would hamper products in the marketplace.)

802.14—Cable Modem Working Group. Originally intended to describe LANs for cable TV systems. This WG has been disbanded.

802.15—Wireless Personal Area Network (WPAN) Working Group. Defines a MAC and PHY for a short distance wireless network between portable and mobile devices such as PCs, personal digital assistants (PDAs), cell phones, pagers, and other communications equipment.

802.16—Broadband Wireless Access (BBWA) Working Group. Defines the MAC and PHY for high-speed wireless network access over relatively short distances. BBWA standards address the “first-mile/last-mile” connection in wireless metropolitan area networks, extending the reach of residential broadband services such as cable modem or digital subscriber line (DSL).

802.17—Resilient Packet Ring (RPR) Working Group. Defines standards to support the development and deployment of RPR local, metropolitan, and wide

LAN PROTOCOLS AND STANDARDS 6 · 19

area networks for resilient and efficient transfer of data packets at rates scalable to many gigabits per second. This WG is currently in *hibernation*.¹⁰

802.18—Radio Regulatory Technical Advisory Group (RR-TAG). On behalf of other 802 WGs using radio-based communication, this TAG monitors, and actively participates in, ongoing national and international radio regulatory activities.

802.19—Wireless Coexistence Working Group. Develops and maintains policies defining the responsibilities of 802 standards developers to address issues of coexistence with existing standards and other standards under development.

802.20—Mobile Broadband Wireless Access (MBWA) Working Group. Defines the specification for a packet-based wireless interface that is optimized for IP-based services. The goal is to enable worldwide deployment of affordable, ubiquitous, always on, and interoperable multivendor mobile broadband wireless access networks that meet the needs of business and residential end user markets. This WG is currently in hibernation.

802.21—Media Independent Handover Services Working Group. Developing standards to enable handover and interoperability between heterogeneous network types including both 802 and non-802 networks.

802.22—Wireless Regional Area Networks (WRAN) Working Group. Developing a standard for a radio-based PHY, MAC, and air interface for use by license-exempt devices on a noninterfering basis in the spectrum allocated to broadcast television.

802.23—Emergency Services Working Group. This working group was created to define an IEEE 802 framework for LANs that would comply with applicable civil authority requirements for communications systems. This working group has been disbanded.

802.24—Smart Grid Technology Advisory Group. This TAG was created to provide liaison between the 802 committee and the smart grid industry and regulatory bodies, and to provide coordination and collaboration amongst 802 working groups related to smart grids.

6.5.3 IEEE 802.3 CSMA/CD Standard. The original IEEE 802.3 standard, first published in 1985, describes the PHY and MAC for a CSMA/CD bus network operating over thick coaxial cable. Today, an 802.3 network implementation can employ any of a number of media types, including UTP and optical fiber. Without question, star-wired UTP implementations are the most popular.

The 802.3 committee anticipated the different media types that might be used, and they developed a nomenclature to identify the actual physical implementation, using the format:

[speed (Mbps)][signaling type][segment length (m) or media type]

The original 802.3 specification, for example, operated at 10 Mbps, used baseband (digital) signaling and limited a single coaxial cable segment to a length of 500 m (1,640 feet); the cable was designated 10BASE5. In fact, the largest distance between two 802.3 stations could be 2.8 km (9,200 feet), so repeaters might be used to interconnect several 500-m coaxial cable segments.

6 · 20 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

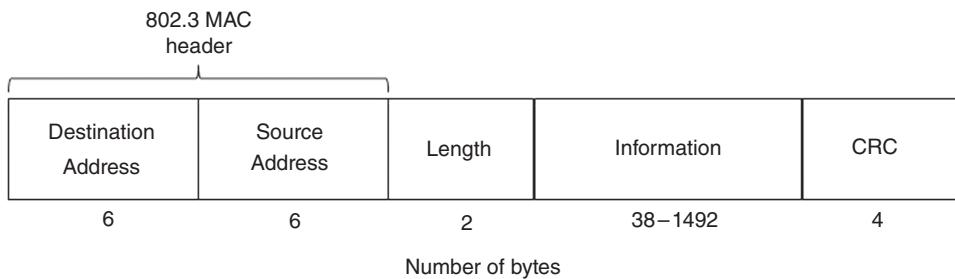


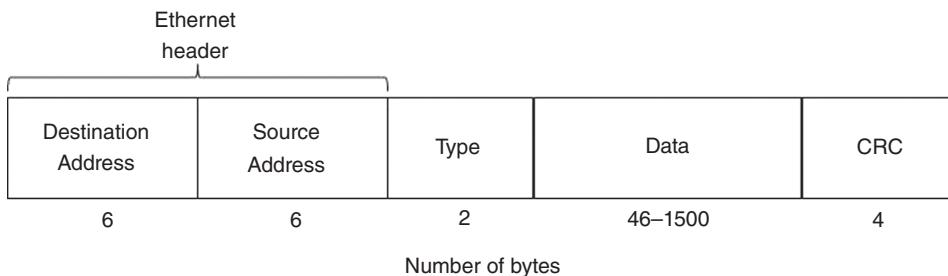
EXHIBIT 6.11 IEEE 802.3 Frame Format

A less expensive version, called *CheaperNet*, was later introduced that operated over thin coaxial cable segments limited to 185 m (610 feet); this PHY is denoted 10BASE2.

In the mid-1980s, AT&T introduced a product called StarLAN, which operated at 1 Mbps over UTP. Although this product has long been relegated to obscurity, it was the first to break the 1-Mbps barrier on UTP. Subsequent versions of 802.3 that employ UTP all use a star topology where each network node connects directly back to a central hub. The first 10-Mbps version of 802.3 was denoted 10BASE-T, the *T* indicating use of the UTP medium (which structured wiring standards say is limited to a distance of 100 m, or 330 feet). The 10-Mbps optical fiber version of 802.3 is 10BASE-F. Today, of course, 100-Mbps and 1-Gbps versions (i.e., 100BASE-T and 1000BASE-T) are available. Full-duplex Ethernet takes advantage of the point-to-point links in a star configuration and effectively doubles the line speed by allowing both stations to transmit at the same time.

Exhibit 6.11 shows the format of an IEEE 802.3 MAC frame, primarily for reference purposes. The fields and their functions are:

- **Preamble.** Used for clock synchronization; employs 7 repetitions of the 8-bit pattern 10101010. (8 binary bits = 1 byte = 1 octet)
- **Start frame delimiter (SFD).** The bit pattern 10101011 denotes the actual beginning of the frame. 1 octet.
- **Destination address (DA).** 48-bit MAC address of the station that should receive this frame. An all-1s address in 48 binary bits (ff-ff-ff-ff-ff-ff in hexadecimal) is the *broadcast address*, indicating that all stations should receive this message.
- **Source address (SA).** 48-bit MAC address of the station sending this frame.
- **Length.** Number of octets in the LLC data field, a value between 0 and 1500. 2 octets.
- **LLC Data.** Data from LLC (and higher layers). This field contains a 3-octet LLC Header, 5-octet 802.2 Subnetwork Access Protocol (SNAP) header, and 38 to 1492 octets of higher layer data.
- **PAD.** Additional octets to ensure that the frame is at least 64 octets in length; this minimum is required by CSMA/CD networks as part of the collision detection mechanism.
- **Frame check sequence (FCS).** Remainder from CRC-32 calculation used for bit error detection. 4 octets.

LAN PROTOCOLS AND STANDARDS 6 · 21**EXHIBIT 6.12** Ethernet II Frame Format

6.5.4 Ethernet II. The IEEE’s CSMA/CD standard is based on the Ethernet specification developed at Xerox’s Palo Alto Research Center (PARC) in the mid-1970s. When Xerox first decided to market Ethernet, there was no OSI model or any LAN standards or products. Given that environment, Xerox sought industry support for this new specification. The Ethernet specification has been jointly distributed (and marketed) by Digital Equipment Corporation (DEC, now Compaq), Intel, and Xerox (hence sometimes known as *DIX Ethernet*). While the 802.3 standard is based on Ethernet II, the two are not exactly the same.

Exhibit 6.12 shows the format of an Ethernet MAC frame, primarily for purposes of comparison to the IEEE frame. The fields and their functions are:

- **Preamble.** Used for clock synchronization; employs the bit pattern 10101010 ... 10101011. 8 octets.
- **Destination address (DA).** 48-bit MAC address of the station that should receive this frame. An all-1s address (ff-ff-ff-ff-ff-ff) is the *broadcast address*, indicating that all stations should receive this message.
- **Source address (SA).** 48-bit MAC address of the station sending this frame.
- **Protocol identifier (PID).** Indicator of the protocol information transported in the Information field. Sample values include 2048 and 2054 to indicate the Internet Protocol (IP) and Address Resolution Protocol (ARP), respectively. 2 octets.
- **Information.** Protocol data unit from the protocol identified in the PID field. 46 to 1,500 octets. (It is the responsibility of the higher layer to ensure that there are at least 46 octets of data in the frame.)
- **Frame check sequence (FCS).** Remainder from CRC-32 calculation used for bit error detection. 4 octets.

The point in comparing the frame formats of Ethernet and 802.3 is primarily of historical purposes because today’s implementations are 802.3 and not Ethernet. That said, it is interesting to note that the two specifications are, in fact, different. It is a minor thing, perhaps, and was a common misnomer in the industry to refer to *IEEE 802.3 Ethernet* (even the IEEE 802.3 committee is now known as the *Ethernet Working Group*), but it was an important difference to both a network administrator and a security professional.

In particular, in years past, if one LAN device only understood Ethernet encapsulation, it would not be able to communicate successfully with another LAN device that only understood IEEE 802.3 encapsulation. Both devices, however, can share the same

6 · 22 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

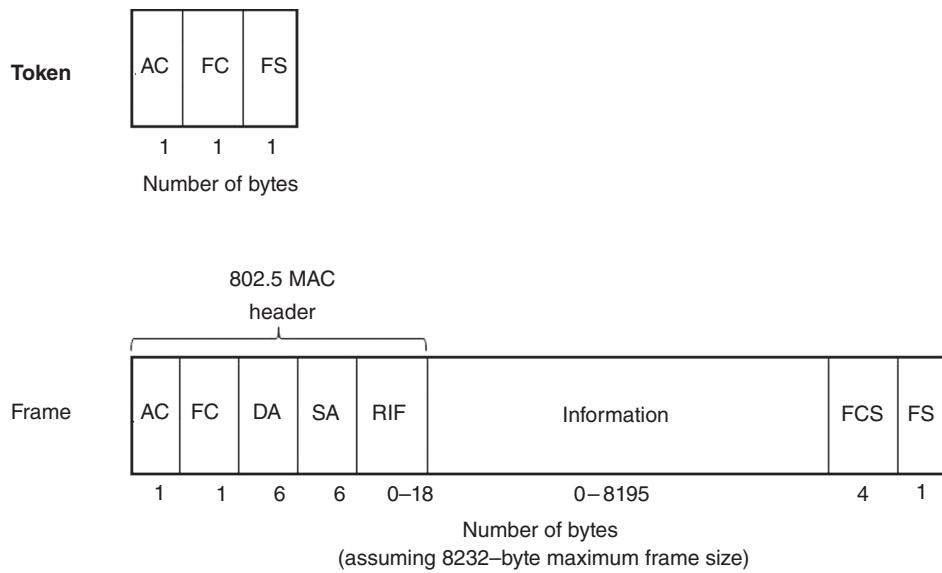


EXHIBIT 6.13 IEEE 802.5 Token and Frame Formats

medium backbone because the electronics are the same. A NetWare server running the Internetwork Packet Exchange (IPX) network layer protocol over IEEE 802.3 frames, for example, could easily share the network with a UNIX host running IP over Ethernet, but it maintained some immunity from attack by an IP host for the NetWare server because the two networks could not cross-communicate.

6.5.5 IEEE 802.5 Token-Ring Standard. The IEEE 802.5 token-ring standard was based on the IBM product of the same name. Both the standard and the product date back to about 1985. It is described here for historical purposes.

The token ring has a logical ring topology, although it was built as a physical star. Designed to operate with STP or UTP cable, most implementations operated at speeds of 16 Mbps or higher. The 802.5 MAC was essentially the same as the token passing scheme described in Section 6.4.2. The fields of the MAC frame (see Exhibit 6.13) are:

- **Start delimiter (SD).** Marks the actual beginning of the transmission. Bit pattern JK0JK000, where J and K represent special symbols on the line.¹¹ 1 octet.
- **Access control (AC).** Indicates whether this transmission is a *token* (i.e., no data) or a *frame* (i.e., contains data). This field also contains information about the priority of this transmission. 1 octet.
- **Frame control (FC).** Indicates if this frame carries LLC (and higher-layer) data or MAC management information; if it is MAC-specific information, this field also indicates the MAC frame type. 1 octet.
- **Destination address (DA).** 48-bit MAC address of the station to which this frame is intended.
- **Source address (SA).** 48-bit MAC address of the station sending this frame.
- **Routing information field (RIF).** An optional field, used only in multiple-ring networks utilizing source routing *and* in which the intended receiver is on a different ring than the transmitter. In source routing, the transmitter can specify

LAN PROTOCOLS AND STANDARDS 6 · 23

the intended path of this frame, designating up to eight intermediate networks.¹² 0 to 18 octets.

- **Information (INFO).** Contains an LLC frame or MAC management information. No maximum length is specified by the standard, but the length of this field will be limited by the time required to transmit the entire frame, controlled by the *token holding time* parameter.
- **Frame check sequence (FCS).** Remainder from a CRC-32 calculation to detect bit errors in the frame. 4 octets.
- **End delimiter (ED).** Demarks the end of the transmission, with the bit pattern JK1JK1IE, where J and K are as described in the SD field. The I-bit indicates whether this frame is the last frame of a multiple-frame sequence and the E-bit indicates whether a bit error was detected by the receiver (E); these bits are cleared by the original sender when the frame returns to that station. 1 octet.
- **Frame status (FS).** The bit pattern AC00AC00; these bits indicate whether the frame's destination address was recognized by any station on the network (A) and whether this frame was successfully copied by the intended receiver (C). 1 octet.

As shown, a token comprises just three octets, the SD, AC, and ED fields. A station sends a frame whenever there is user data or MAC information to send. The station must wait until it receives a *token* before it can generate a *frame*.

The transmitting station is responsible for generating a new token after it transmits a single frame. Recall that the transmitted bits come back to the sender, and it is this station that removes the bits from the network. According to the original standard, the transmitter will send a token after sending all of the bits of the frame and must wait until it has seen at least the returning SA field to verify that it is, in fact, removing its own frame from the network. Optionally, early token release allows the transmitter to generate a new token immediately after finishing sending the bits from its frame, even if the SA field has not yet returned. This latter option was developed to improve performance in very large token ring environments, such as the American National Standards Institute (ANSI) FDDI standard.

Today, 802.5 token rings are primarily limited to IBM environments, and there is a lot to be found there. FDDI is more commonly found in multibuilding campus environments, used as a backbone to interconnect Ethernet/802.3 networks. FDDI is being phased out; the last FDDI product vendor dropped out of the marketplace in 1999.

6.5.6 IEEE 802.2 LLC Standard. The IEEE 802.2 LLC protocol was intended to provide a common interface between 802 LAN MACs and higher-layer applications. With the LLC, the underlying MAC scheme is transparent to the application just as the application is transparent to the MAC.

The LLC was designed to support any number of services, the most common being an unacknowledged connectionless service (primarily used in contention networks) and an acknowledged connection-oriented service (primarily used in token ring environments).

The LLC is loosely based on the Higher-layer Data Link Control (HDLC) bit-oriented protocol in both operation and frame format (see Exhibit 6.14). The LLC frame appears in the Information field of a MAC frame. The first two fields of the LLC header are the *Destination Service Access Point* (DSAP) and the *Source Service Access Point* (SSAP) fields, originally intended to identify the higher-layer services at the source

6 · 24 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

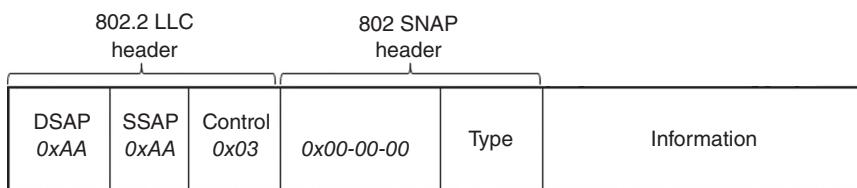


EXHIBIT 6.14 IEEE 802.2 LLC Frame Transporting SNAP Header (which in turn indicates IEEE organization and EtherType protocol identifiers)

and destination node. This is similar in concept to *ports* in the Transmission Control Protocol (TCP) and User Datagram Protocol (UDP) but was never well implemented, and the DSAP and SSAP values are typically the same. The third field is the Control field, identifying the type of frame.

The Subnetwork Access Protocol (SNAP) is an IEEE 802 protocol that can be used to identify *any* protocol created by *any* agency, and is commonly used above the LLC layer. In this case, the SNAP header immediately follows the LLC header. Use of SNAP is indicated by the LLC fields when both DSAP and SSAP fields are set to a value of 170 (0xAA) and the Control field is set to a value of 3 (octal 03) to indicate that it is an Unnumbered Information frame.

The SNAP header has two fields. The 3-byte *Organizationally Unique Identifier* (OUI) field refers to the organization that developed either the higher-layer protocol or a way to refer to the protocol. The 2-byte *Type* field identifies the protocol using the Organization-defined number.

The Internet Protocol (IP) and Address Resolution Protocol (ARP) provide example uses of SNAP. The common format of a SNAP header encapsulating these protocols would be to set the OUI value to 0 (0x00-00-00) to identify IEEE/ISO as the organization. The Type field would then use the *EtherType* values of 2048 (0x08-00) and 2054 (0x08-06) to indicate use of IP and ARP, respectively.

6.5.7 Summary. This section has covered the important LAN standards governing what is most likely to be seen in the industry today. Table 6.1 summarizes some of the discussion about the most common LAN topologies, media, MAC schemes, and standards.

6.6 INTERCONNECTION DEVICES. LAN interconnection devices are used to attach individual LANs to each other in order to build a large enterprise network. They can also interconnect LAN components across a WAN and provide LAN access to the

TABLE 6.1 LAN Characteristics

Physical Topology	Logical Topology	Media	MAC	Speed (Mbps)	Standard
Bus	Bus	Coax	CSMA/CD	10	802.3, Ethernet
Star	Bus	UTP, Fiber	CSMA/CD	1–1000	802.3
Star	Bus	Wireless	CSMA/CA	1–150	802.11
Star	Ring	UTP	Token passing	16	802.5
Ring	Ring	Fiber	Token passing	100	FDDI

INTERCONNECTION DEVICES 6 · 25

Internet. Several types of such devices are used for LAN interconnections, including hubs, switches, bridges, and routers. The major distinction between these devices is the OSI layer at which they operate, and all are discussed in the next sections.

6.6.1 Hubs. *Hubs* are used to build physically star-wired LANs, using media that are basically point-to-point in nature (such as UTP and optical fiber). Note that it is the internal wiring of the hub that determines its *logical* nature, so that a logical bus or ring LAN can be physically star-wired.

So-called *Ethernet hubs* support 10, 100, and/or 1,000 Mbps Ethernet or 802.3 networks. Different hubs will have a different number of ports, generally ranging from 4 to 32. Hubs provide physical connectivity only; when a frame arrives on one port, the hub will broadcast the frame back out to all other ports, which simulates the broadcast bus environment. Multiple hubs can be interconnected to form reasonably large networks.

Token-ring hubs, generally called *multistation access units* (MAUs), look similar to Ethernet hubs but have different internal wiring. When an MAU receives a transmission on one port, it merely forwards that transmission, a bit at a time, to the next port sequentially on the MAU. In this way, it simulates the ring environment.

6.6.2 Switches. *Switches* are generally employed in the CSMA/CD environment and extend the capabilities of a hub. A switch operates at a combination of PHY and MAC layers. In addition to providing physical connectivity like a hub, a switch learns the MAC address of all stations attached to it. When a frame arrives on a switch port, the switch looks at the destination MAC address and places the frame on the port associated with that address (which might be the port leading to another switch).

Switches are used primarily to improve performance. Given the scenario described earlier, multiple stations can transmit simultaneously without collision. Furthermore, switches can operate in full-duplex mode, meaning that a single station can both transmit and receive at the same time. A 10 Mbps switched Ethernet LAN, for example, can achieve performance similar to that of a 100-Mbps hubbed Ethernet LAN. (This was a real boon in those environments where it is not viable to upgrade 10-Mbps NICs and wiring.)

There is a subtle security ramification to the use of switches versus hubs. In particular, if a user places a packet sniffer on a hubbed LAN, the sniffer will see every frame because the hub simulates the broadcast environment. A packet sniffer on a switched network will not be as effective; it will only pick up those frames that are specifically addressed to the LAN broadcast address. That said, many switches come with an administrative port that can be set to monitor all ports for troubleshooting purposes.

6.6.3 Bridges. A *bridge* provides a point-to-point link between the two LANs, usually those employing similar MAC schemes. Bridges operate at the MAC layer, and their operation is controlled by the MAC address.

Ethernet environments commonly employ *learning bridges*. In a very simple case, consider a bridge interconnecting two LANs, #1 and #2 (see Exhibit 6.15). When any LAN station sends a frame, both destination and source MAC addresses are included in the transmission. As frames appear on the networks, the bridge sees all of the source addresses and builds a table associating the MAC addresses with one LAN or the other, eventually learning the location of all of the network's stations. This process is sometimes called *backward learning* because the bridge learns the location of stations that transmit.

6 · 26 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

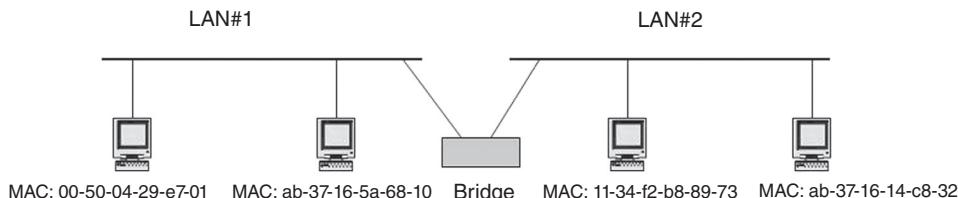


EXHIBIT 6.15 Two LANs Interconnected via a Bridge

A bridge is a simple frame store-and-forward device. Like all stations on the LAN, a bridge examines the destination address of any transmitted frames. If a transmission on LAN #1 contains a destination address of a station on LAN #2, the bridge will forward the frame. If a transmission contains an unknown destination address, the bridge will also forward the frame.

Although a bridge bases its decisions on the MAC address, it is not an intelligent device; that is, it knows that a station with a particular MAC address is in one direction or another, but it does not know precisely where that station is. Because bridges have to build tables containing all of the stations' addresses that they learn, bridges do not scale particularly well to large networks. Bridges also extend the *broadcast domain* (i.e., if a frame transmitted on LAN #1 is sent to the broadcast address, it will be forwarded to LAN #2).

6.6.4 Routers. A *router* is conceptually similar to a bridge in that it is also a store-and-forward device. A router, however, works at the Network Layer and is therefore a much more powerful device than a bridge. As Exhibit 6.16 shows, every LAN device has both a MAC (hardware) and Network Layer (software) address (in this case, IP is the sample Network Layer address). Because Network Layer addresses are hierarchical, the networks themselves have a network identifier number (the NET_ID in Exhibit 6.16). Network Layer addresses are well suited to environments where intermediate devices have to find a best route between networks.

Like a bridge, a router is considered to be just another station on a LAN to which it is attached. If the router sees a transmission on LAN #1 (with a NET_ID of 192.168.16.0) containing a destination address of a station on another network, it will route the packet to the correct destination network, even if that means going through another router to get there.

This example also demonstrates another major difference between bridges and routers. In a bridged environment, a station on LAN #1 sends a frame to some MAC address and has no knowledge of whether the intended destination is on the same LAN or not; the bridge will forward the frame if necessary, but this is all transparent to sender

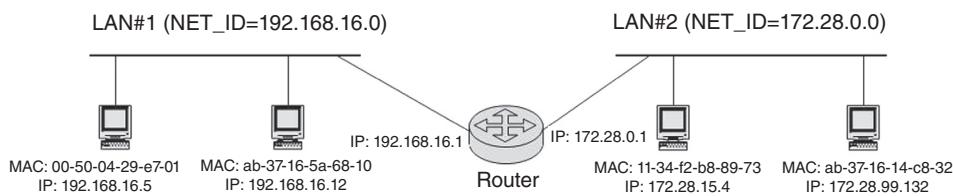


EXHIBIT 6.16 Two LANs Interconnected via a Router

NETWORK OPERATING SYSTEMS 6 · 27

and receiver. In a routed environment, however, the sender can tell if the receiver is on the same or different network merely by examining the destination network address. In fact, the router only gets involved if the packet has to leave the local network; that is why in an IP environment, for example, an address of a default gateway (router) has to be provided.

Routers also limit the broadcast domain. If a station on LAN #1 transmits a frame using the broadcast MAC address, the frame goes no further than the router.

Routers build their routing tables very differently than bridges. Whereas bridges learn the relative location of a station by observing a frame's source address, packets learn the Network Layer address by the use of routing protocols that allow groups of routers to exchange routing information.¹³

6.6.5 Summary. Hubs, switches, bridges, and routers are all commonly employed LAN interconnection devices. These are tools in the kit of everyone who works with LANs, as the building blocks of everything from small and intermediate-size local networks to large enterprise networks and the global Internet.

6.7 NETWORK OPERATING SYSTEMS. Just as an operating system manages computer resources, a *network operating system* (NOS) provides the software that controls the resources of a LAN. NOSs generally comprise software that provides at least these functions:

- *Hardware drivers* are the software that allows the NOS to communicate with the NIC.
- *Communications software* allows applications running on different LAN nodes to communicate.
- *Services* are the functional aspects of the NOS and the reason that people use a LAN in the first place. Sample services include file services (file sharing), print services (commonly shared printers), message services (email), communication services (LAN access to the Internet), and fax services (commonly shared facsimile).

NOSs are typically classified as being peer to peer or client/server. A *peer-to-peer* LAN allows any LAN node to communicate with any other LAN node, and any LAN node can provide services to other nodes. In a *client/server* (or *server-based*) environment, every node is either a client or a server. In this scenario, *servers* are special nodes that offer services to other servers or to clients, while *clients* are the ordinary end-user workstations. Clients can only communicate with a server.

When evaluating or investigating the security of a LAN, the software is the most common point of exposure, vulnerability, and exploitation, particularly for remote attacks.

Some sample NOSs that have had historical significance include:

AppleTalk. Apple Macs have come with integrated LAN capabilities since their inception in 1985. Originally using a scheme called LocalTalk, AppleTalk was a peer-to-peer network running over a 10-Mbps CSMA/CD LAN. The Network Layer protocol historically associated with AppleTalk was called the

6 · 28 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

Datagram Delivery Protocol (DDP). Apple dropped support for AppleTalk in 2009 and supports TCP/IP-based networking.

Microsoft Networking. Microsoft operating systems have come with LAN capabilities since Windows for Workgroups (WfW or Windows 3.11). Employing a nonroutable protocol called the Network Basic Input/Output System (NetBIOS) Extended User Interface (NetBEUI), Windows client systems (Windows 3.11 and later) can be easily used to build an inexpensive, simple peer-to-peer LAN for file and print sharing. NetBEUI is nonroutable because it does not provide an addressing mechanism to allow interconnected yet distinct NetBEUI subnetworks; if two NetBEUI networks are attached in any way, they will appear to be one large network. (This is why the hard drive of improperly configured Windows systems can be viewed from across the Internet.) Microsoft defined NetBIOS over TCP/IP (NBT) in the 1980s to support encapsulation of NetBIOS messages in TCP and UDP messages. Most Microsoft networking today relies more on TCP/IP than NetBEUI.

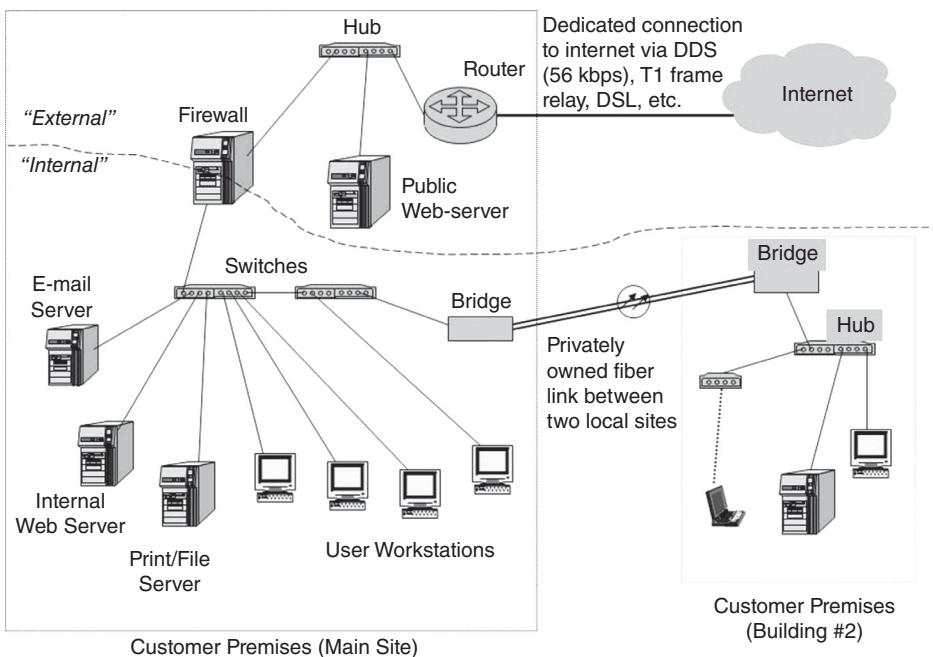
Microsoft Windows Server (including Windows NT Server, Windows Server 2003, Windows Server 2008, and Windows Server 2012). The NOS of choice for Windows network environments, Windows Server is Microsoft's client/server operating system. Client systems can run nearly any Windows operating system, from Windows XP or NT Workstation to Windows VISTA, Windows 7, or Windows 8. Clients on a Windows Server network themselves can form a peer-to-peer network. Microsoft networking is reliant on Active Directory and TCP/IP.

Novell NetWare. Novell offered one of the first PC-class networks in the early 1980s using a proprietary star-based LAN. By the 1990s, NetWare was the best-known client/server NOS and accounted for more than 70 percent of the NOS market. The Network Layer protocol associated with classical NetWare is the Internetwork Packet Exchange (IPX) protocol. By the turn of the century, Microsoft was eating dramatically into NetWare's popularity and TCP/IP was dominating as the protocol-of-choice. At around this time, Novell embraced Linux and TCP/IP, and NetWare has been replaced with the Open Enterprise Server.

UNIX and Linux. TCP/IP has been the network communications protocol for UNIX systems since 1984, allowing UNIX-based hosts to build client/server (and, ultimately, peer-to-peer) networks. TCP/IP has also been integral to Linux since its inception in 1991. With TCP/IP, any system can run server (daemon) software to provide services to other systems, so that any system can act as a client or server, depending on application.

6.8 SUMMARY. Exhibit 6.17 shows a possible network design that includes many of the elements that have been described in this chapter (and a few that have not). This network's router provides the interface to the Internet and is attached via some sort of dedicated connection, such as a point-to-point 56 Kbps or T1 (1.544 Mbps) leased line, frame relay, or digital subscriber line.

In this scenario, the router is physically located at the main site. From a security perspective, the organization may segment its network into an *external* and *internal* side, the internal being protected by a firewall.¹⁴ The external network includes the router, public Web server, and firewall. Those three systems are interconnected through

SUMMARY 6 · 29**EXHIBIT 6.17** LAN Scenario

a hub to which they each attach via a Cat 5 UTP cable. In this scenario, the hub could actually implement 10BASE-T or 100BASE-T Ethernet, or even a token ring.

The external and internal networks are connected through the firewall, which, in this case, will have two NICs. The two networks are separate and distinct; the firewall does not extend the broadcast domain of either network, and, in fact, these two networks would have different IP network identifiers.

The internal network at the main site is a collection of servers and user workstations that are interconnected via a set of switches. In this example, these are 8-port 100-Mbps Ethernet switches. Since there are more than 8 devices, the switches themselves need to be interconnected. There are several options for that:

- *Stackable* switches physically attach to each other, extending the switch's backplane to create a larger switch (in this case, a 16-port switch).
- An *optical fiber* link can be used to interconnect the switch, usually at backplane speeds in the 1+-Gbps range.
- A *UTP* link might be used to interconnect the switches via two of the 100-Mbps ports.

To connect the LAN in Building #2 with the LAN at the main site, a point-to-point connection between a pair of bridges would suffice. In this case, the buildings are several kilometers apart, necessitating use of optical fiber.

In Building #2, there is another hub-based LAN, with a laptop using wireless technology, communicating with an access node that is also attached to the hub.

This chapter has only skimmed the surface of LAN concepts, standards, and technologies. Their study is important to the security professional, however, because LANs are the basis of all networking. As a *network of networks*, the Internet comprises

6 · 30 LOCAL AREA NETWORK TOPOLOGIES, PROTOCOLS, AND DESIGN

millions of local networks. This chapter indicates many of the points of potential vulnerability or compromise in a system.

6.9 FURTHER READING

- Cisco Systems. *Internetworking Technology Handbook*. Cisco DocWiki Web-site. June 20, 2012. http://docwiki.cisco.com/wiki/Internetworking_Technology_Handbook
- Gast, M. “Wireless LAN Security: A Short History.” O’Reilly | Wireless Devcenter Website. April 19, 2002. www.oreillynet.com/pub/a/wireless/2002/04/19/security.html
- Mikalsen, A., and P. Borgesen. *Local Area Network Management, Design and Security: A Practical Approach*. Hoboken, NJ: John Wiley & Sons, 2002.
- Riley, S., and R. A. Breyer. *Switched, Fast, and Gigabit Ethernet*, 3rd ed. Indianapolis, IN: New Riders Publishing, 1998.
- Stallings, W., and T. Case. *Business Data Communications—Infrastructure, Networking and Security*, 7th ed. Upper Saddle River, NJ: Prentice-Hall, 2012.

6.10 NOTES

1. See Chapter 5 in this *Handbook* for additional background information on LANs and WANs.
2. See Sections 22.4.5 and 23.9.1 in this *Handbook* for further discussion of wiretaps.
3. It is an interesting point of trivia to note that the original frequency-hopping spread spectrum technique was co-invented by Hollywood star Hedy Lamarr in 1940 and given, free of charge, to the U.S. Navy.
4. See Chapter 33 in this *Handbook* for more details of wireless LAN security.
5. M. Kende, “Internet Global Growth: Lessons for the Future,” Analysys Mason Knowledge Centre (Website), 2012, www.analysysmason.com/internet-global-growth-lessons-for-the-future (p. 14).
6. Since CSMA/CD transmits with a probability of 1, it is sometimes referred to as being 1-persistent.
7. As an aside, although the station can experience 16 collisions, the probability of transmission will never fall below 1/1024, or 2-10, since Ethernet and IEEE 802.3 do not allow more than 1,024 devices on the network. This is the source of the word “truncated” in the name of the scheme.
8. Up-to-date status information about the 802 committee can be found at the LAN/MAN Standards Committee Web site at <http://grouper.ieee.org/groups/802/>.
9. A WG is disbanded when it is considered that there is no more work for the IEEE to undertake in this topic area.
10. A WG will go into hibernation when there are no new projects to undertake. This status indicates a WG that has reached status quo.
11. The term “special symbol” requires explanation. The signaling scheme used in the token ring PHY standard is called Differential Manchester. In this signaling scheme, the signal is at a positive voltage for half of the bit time and at a negative voltage for the other half of the bit time, meaning that each bit has a sum total of 0 volts (resulting in what is sometimes called DC balancing). The J and K symbols are Differential Manchester code violations, where one symbol is at negative voltage for an entire bit time and the other at positive voltage for an entire bit

NOTES 6 · 31

time. These code violations have the benefit of being able to indicate special events and can be used for synchronization. J and K symbols are always used in pairs to maintain DC balancing.

12. Source routing is a very rarely used option in IP and is, in fact, a security problem; firewall administrators routinely set up filters to block IP packets with source routing. Source routing in an 802.5 network, however, is a normal feature and is not considered to be a security threat because this information has no impact on the WAN.
13. In the IP environment, common routing protocols include the Border Gateway Protocol (BGP), Open Shortest Path First (OSPF), and Routing Information Protocol (RIP).
14. This is a very simplistic firewall design with the internal and external network. The focus of this diagram is on the LAN components, however, rather than the specific security architecture.

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 7

ENCRYPTION

Stephen Cobb and Corinne LeFrançois

7.1 INTRODUCTION TO CRYPTOGRAPHY	7·1	7.4 PUBLIC KEY ENCRYPTION	7·21
7.1.1 Terminology	7·2	7.4.1 Key-Exchange Problem	7·23
7.1.2 Role of Cryptography	7·3	7.4.2 Public Key Systems	7·23
7.1.3 Limitations	7·6	7.4.3 Authenticity and Trust	7·26
7.4.4 Limitations and Combinations		7.4.4 Limitations and Combinations	7·27
7.2 BASIC CRYPTOGRAPHY	7·6	7.5 PRACTICAL ENCRYPTION	7·27
7.2.1 Early Ciphers	7·6	7.5.1 Communications and Storage	7·28
7.2.2 More Cryptic Terminology	7·8	7.5.2 Securing the Transport Layer	7·28
7.2.3 Basic Cryptanalysis	7·8	7.5.3 X.509v3 Certificate Format	7·31
7.2.4 Brute Force Cryptanalysis	7·9		
7.2.5 Monoalphabetical Substitution Ciphers	7·11	7.6 BEYOND RSA AND DES	7·36
7.2.6 Polyalphabetical Substitution Ciphers	7·12	7.6.1 Elliptic Curve Cryptography	7·36
7.2.7 The Vigenère Cipher	7·13	7.6.2 RSA Patent Expires	7·37
7.2.8 Early-Twentieth-Century Cryptanalysis	7·14	7.6.3 DES Superseded	7·38
7.2.9 Adding up XOR	7·15	7.6.4 Quantum Cryptography	7·39
7.6.5 Snake Oil Factor		7.6.5 Snake Oil Factor	7·44
7.3 DES AND MODERN ENCRYPTION	7·17	7.7 STEGANOGRAPHY	7·44
7.3.1 Real Constraints	7·17	7.8 FURTHER READING	7·45
7.3.2 One-Time Pad	7·17	7.9 NOTES	7·46
7.3.3 Transposition, Rotors, Products, and Blocks	7·18		
7.3.4 Data Encryption Standard	7·20		
7.3.5 DES Strength	7·20		
7.3.6 DES Weakness	7·21		

7.1 INTRODUCTION TO CRYPTOGRAPHY. The ability to transform data so that they are accessible only to authorized persons is just one of the many valuable services performed by the technology commonly referred to as encryption. This technology has appeared in other chapters, but some readers may not be familiar with its principles and origins. The purpose of this chapter is to explain encryption technology in basic terms and to describe its application in areas such as file encryption, message scrambling, authentication, and secure Internet transactions. This is not a theoretical

7 · 2 ENCRYPTION

or scientific treatise on encryption, but a practical guide for those who need to employ encryption in a computer security context.

Organizations around the world increasingly rely on cryptography to communicate securely and to store information safely. Typically, the algorithms used by the Department of Defense (DoD) organizations are employed and maintained for many years. For example, the Data Encryption Standard (DES) has been used in some form for over 20 years.¹

This chapter is a brief overview of cryptography and its practical applications to the needs of normal business users, as distinct from the needs of high-security government agencies. A thorough examination of the mathematics that are the foundation of these topics is beyond the scope of this chapter, but we provide suggested readings for further study.

7.1.1 Terminology. This list of basic terms will be helpful for readers as they continue through this chapter:

Algorithm—a finite list of well-defined instructions for accomplishing some task that, given an initial state, will terminate in a defined end state.

Cipher—the core algorithm used to encrypt data. A cipher transforms plaintext into ciphertext that is not reversible without a key.

Ciphertext—text in encrypted form, as opposed to the plain text. We show ciphertext in **UPPERCASE** throughout this chapter.

Codes—a list of equivalences (*a codebook*) allows the substitution of meaningful text for words, phrases, or sentences in an innocuous message; for example, “I will buy flowers for Mama tomorrow for her party at 7 pm” might be decoded to mean “Launch the attack on the mother ship next week on Sunday.”

Decrypt/Decipher—the process of retrieving the plaintext from the ciphertext.

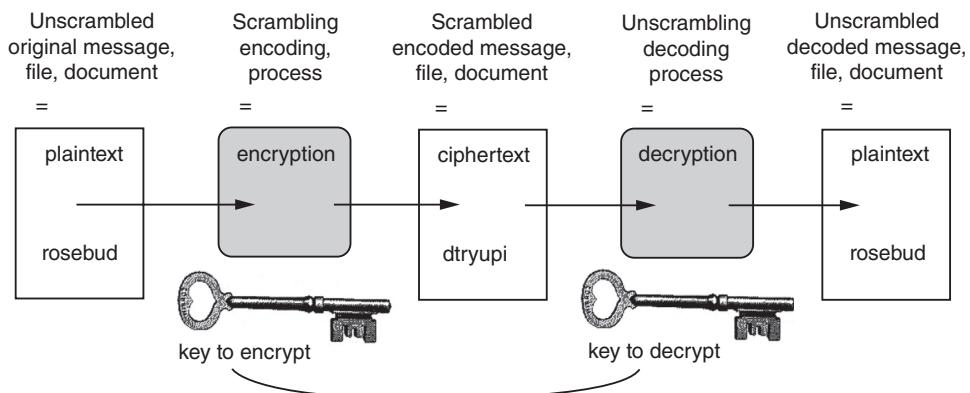
Encrypt/Encipher—to alter plaintext using a secret code so as to be unintelligible to unauthorized parties.

Key—a word or system for solving a cipher or code.

Plaintext—the original message to be encoded or enciphered. We show plaintext in lowercase throughout this chapter.

The science of cryptology (sometimes abbreviated as *crypto*) is the study of secure communications, formed from the Greek words *κρυπτοσ* (*kryptos*), meaning “hidden,” and *λογοσ* (*logos*), “word.” More specifically, it is the study of two distinct, yet highly intertwined, fields of study: cryptography and cryptanalysis. Cryptography is “the science of coding and decoding messages so as to keep these messages secure.”² Cryptanalysis is the art and science of “cracking codes, decoding secrets, violating authentication schemes, and in general, breaking cryptographic protocols,”³ all without knowing the secret key. Systems for encrypting information are referred to as *cryptosystems*.

Systems for encrypting information may also be referred to as *ciphersystems*, from *cipher*, meaning “zero,” or “empty” (a word rooted in the Arabic *sifir*). Terms using cipher and crypto are interchangeable, with some authors preferring cipher to avoid the religious and cultural connotations of *crypt*, a word with the same root as “encryption.” Thus, encryption may be referred to as encipherment, decryption referred to as decipherment, and so on.

INTRODUCTION TO CRYPTOGRAPHY 7 · 3**Exhibit 7.1** Diagram of Cryptographic Terms

The most obvious use of encryption is to scramble the contents of a file or message, using some form of shared secret as a *key*. Without the key, the scrambled data remain hidden and cannot be unscrambled or *decrypted*. The total number of possible keys for an encryption algorithm is called the *keyspace*. The keyspace is a function of the length of the key and the number of possible values in each position of the key. For a *keylength* of n positions, with each position having v possible values, then the keyspace for that key would be v^n . For example, with three positions and two values per position (e.g., 0 or 1), the possible keys would be 000, 001, 010, 011, 100, 101, 110, and 111 for a total keyspace of 8.

In cryptographic terms, the contents of a file before encryption are *plaintext*, while the scrambled or encoded file is known as *ciphertext* (see Exhibit 7.1). As a field of intellectual activity, cryptology goes back many millennia. Used in ancient Egyptian, China, and India, it was discussed by the Greeks and regularly employed by the Romans. The first European treatise on the subject appeared in the fourteenth century. The subject assumed immense historic importance during both world wars. The British success in breaking codes that the Germans used to protect military communications in World War II was a major factor in both the outcome of the war and in the development of the first electronic computer systems.

Since then, cryptography and computer science have developed hand in hand. In 1956, the United States National Security Agency (NSA), the U.S. Government department in charge of monitoring the worldwide flow of information, began funding improvements in computer hardware, pumping some \$25 million into Project Lightning. This five-year development effort, intended to produce a thousand-fold increase in computing power, resulted in over 150 technical articles. It also gave rise to more than 300 patent applications and succeeded in advancing the frontiers of hardware design. The NSA, based in Fort Meade, Maryland, was also involved in the creation of DES as the commercial encryption standard for much of the last 20 years. Today, the NSA is widely believed to have the world's largest collection of supercomputers and the largest staff of cryptanalysts.

7.1.2 Role of Cryptography. The central role of cryptography in computer security is ensuring the confidentiality of data. But cryptography can support other pillars of computer security, such as integrity and authenticity. This section looks at the different roles of cryptography.

7.4 ENCRYPTION

7.1.2.1 Confidentiality. The role of encryption in protecting confidentiality can be seen in a classic definition of encryption: “Encryption is a special computation that operates on messages, converting them into a representation that is meaningless for all parties other than the intended receiver.”⁴

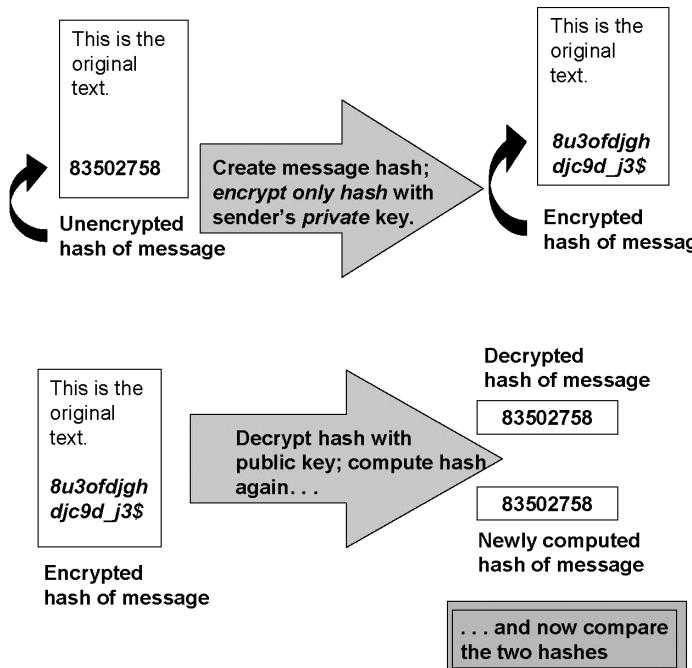
Much of the literature on cryptography discusses the technology in terms of ensuring the confidentiality of messages, but this is functionally equivalent to protecting the confidentiality of data. The use of the term “message” reflects the traditional use to which cryptography has been put, both before and after the advent of computers. For example, Julius Caesar encrypted messages to Cicero 2,000 years ago, while today messages between a Web browser and a Web server are encrypted when performing a “secure” transaction.

When applying cryptography to computer security, it is sometimes appropriate to substitute the term “files” for “messages.” For example, hard drive encryption programs protect data files stored on a hard drive. However, data files take the form of messages when they are transferred from one computer to another, across a network, the Internet, or via phone lines. Practically speaking, data being transferred in this manner are exposed to a different set of dangers from those that threaten data residing on a computer in an office. Thus, the use of encryption to render files useless to anyone other than an authorized user is relevant both to files in transit and to those that reside on a server or a stand-alone computer, particularly when the latter is a laptop, notebook, or PDA.

7.1.2.2 Integrity. In the second half of the last century, following the advent of programmable computer systems, the ability of cryptography to transform data was applied in many new and interesting ways. As will be seen in a moment, many cryptographic techniques use a lot of mathematical calculation. The ability of computers to perform many calculations in a short period of time greatly expanded the usefulness of cryptography, and also inspired the development of ever-stronger ciphersystems.

Maintaining the integrity of data is often as important as keeping them confidential. When writing checks, people take pains to thwart alteration of the payee or the amount. In some cases, integrity is more important than confidentiality. Changing the contents of a company press release as it passes from the company to the press could have serious consequences. It is not only human actions that threaten data integrity; mechanical failures and logical errors can also change data. It is vital that such changes be detected, as was discussed in Chapter 4 of this *Handbook*, where it was observed that “[a]ll data movements and translations increase the likelihood of internal error, and for this reason parity checks and validity tests have become indispensable.”

That chapter covered the role of parity bits for error detection, the function of redundancy checks, and the use of checksums to provide a modification-detection capability. A type of cryptographic hash or checksum called a Message Authentication Code (MAC) can protect against intentional, but unauthorized, data modification as well as against accidental modification. A MAC is calculated by applying a cryptographic algorithm and a secret value, called the *key*, to the data. The data are later verified by applying the cryptographic algorithm and the same secret key to the data to produce another MAC; this MAC then is compared to the initial MAC. If the two MACs are equal, then the data are considered authentic (see diagram in Exhibit 7.2, which uses the public key cryptosystem, discussed later). Otherwise, an unauthorized modification is assumed (any party trying to modify the data without knowing the key would not know how to calculate the appropriate MAC corresponding to the altered data).

INTRODUCTION TO CRYPTOGRAPHY 7 · 5**EXHIBIT 7.2** Message Authentication Code Using Public Key

Cryptosystem

Source: Copyright © 2008 M. E. Kabay. Used with permission.

7.1.2.3 Authentication. In the context of computer security, authentication is the ability to confirm the identity of users. For example, many computers now ask users to log on before they can access data. By requesting a user name and password, systems attempt to assure themselves that only authentic users can gain access. However, this form of authentication is limited—it merely assures that the person logging on is someone who knows a valid user name and password pair. Cryptography plays a very important role in efforts to ensure stronger authentication, from encrypting the password data to the creation and verification of electronic identifiers such as digital signatures. These will be described in more detail later in this chapter, along with the differences between public key and private key cryptography, both of which may be used in these schemes.

Using a public key system, documents in a computer system can be electronically signed by applying the originator's private key to the document. The resulting digital signature and document then can be stored or transmitted. The signature can be verified using the public key of the originator. If the signature verifies properly, the receiver has confidence that the document was signed using the private key of the originator and that the message had not been altered after it was signed. Because private keys are known only to their owner, it is also possible to verify the originator of the information to a third party.

7.1.2.4 Nonrepudiation. An aspect of computer security that has increased greatly in significance, due to the growth in internetwork transactions, is nonrepudiation. For example, if someone places an electronic order to sell stocks that later increase

7 · 6 ENCRYPTION

in value, it is important to prove that the order definitely originated with the individual who placed it. Made possible by public key cryptography, nonrepudiation helps ensure that the parties to a communication cannot deny having participated in all or part of the communication.

7.1.3 Limitations. One role that cryptography cannot fill is defense against data destruction. Although encryption does not assure availability, it does represent a very valuable extra line of defense for computer information when added to physical security, system access controls, and secure channels of communication. Indeed, when computers are mobile, or data are being communicated over insecure channels, encryption may be the main line of defense. However, even though applied cryptography can provide computer users with levels of security that cannot be overcome without specialized knowledge and powerful computers, encryption of data should not be thought of as an alternative to, or substitute for, system access control. According to Seberry and Pieprzyk⁵, the role of cryptography is to protect “information to which illegal access is possible and where other protective measures are inefficient.”

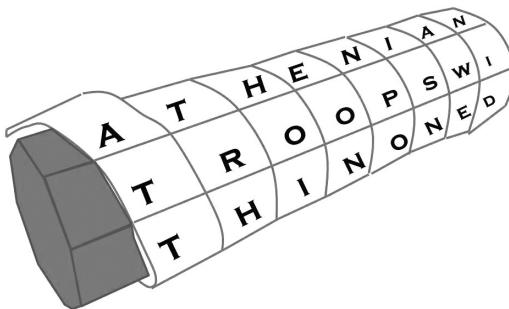
Encryption-based file access controls should be a third barrier after site and system access controls, if for no other reason than that encryption systems alone do little to prevent people deleting files.

7.2 BASIC CRYPTOGRAPHY. The aim of cryptography is to develop systems that can encrypt plaintext into ciphertext that is indistinguishable from a purely random collection of data. This implies that all of the possible decrypted versions of the data except one will be hopelessly ambiguous, with none more likely to be correct than any of the others. One of the simplest ways to create ciphertext is to represent each character or word in the plaintext by a different character or word in the ciphertext, such that there is no immediately apparent relationship between the two versions of the same text.

7.2.1 Early Ciphers. It is believed that the earliest text to exhibit the baseline attribute of cryptography, having a slight modification of the text, occurred in Egypt nearly 4,000 years ago. A scribe used a number of unusual symbols to confuse or obscure the meaning of the hieroglyphic inscriptions on the tomb of a nobleman named Khnumhotep II.⁶

It is also believed that the first effective military use of cryptography was a simple transposition cipher (see Section 7.3.3) by the Spartans, who “as early as 400 BCE employed a cipher device called the scytale for secret communication between military commanders.”⁷ The scytale was a cylindrical or tapered stick with a thin strip of leather or parchment wrapped around it spirally.⁸ The message to be hidden was written lengthwise with no blank spaces. When unraveled, the parchment appeared to hold nothing but random letters. To read the parchment, the recipient had to have a stick with exactly the same dimensions as the sender. The distribution of appropriate decoding scytales took place before the military commanders departed for the field.⁹ For example, a particular combination of stick and strip could allow the cleartext (shown in lowercase):

atheniantroopswithinonedaysmarchofromebereadynow

BASIC CRYPTOGRAPHY 7 · 7**EXHIBIT 7.3** Scytale in Use

Source: Copyright © 2008 M. E. Kabay. Used with permission.

to be broken into up to six rows of eight letters to be written across the rolled-up strip, in this way:

athenian
troopswi
thinoned
aysmarch
ofrombe
readynow

The message might appear on the scytale as shown schematically in Exhibit 7.3.

Reading the unwrapped strip without the stick would produce this ciphertext (shown in uppercase):

ATTAORTRHYFEHOISREEONMODNPOAMYISNRENAWCBCONIDHEW

“The first attested use of [a substitution cipher] in military affairs comes from the Romans.”¹⁰ During that time, Julius Caesar encoded all his messages by simply replacing every letter with the letter three positions away. For example, the letter *a* would become the letter *d*, the letter *b* would become the letter *e*, and so on. Now called the Caesar cipher, this scheme is best-known of all the monoalphabetic algorithms (see Section 7.2.5).¹¹ Consider the Caesar cipher illustrated in the next comparison using the modern English alphabet, with the letters of the alphabet simply shifted three places.

Plaintext: abcdefghijklmnopqrstuvwxyz
Ciphertext: DEFGHIJKLMNOPQRSTUVWXYZABC

To encrypt a message, the sender finds each letter of the message in the plaintext alphabet and uses the letter below it in the ciphertext alphabet. Thus, the clear message:

Plaintext: beware the ides of march

is transformed into the encrypted message:

Ciphertext: EHZDUH WKH LGHV RI PDUFK

7 · 8 ENCRYPTION

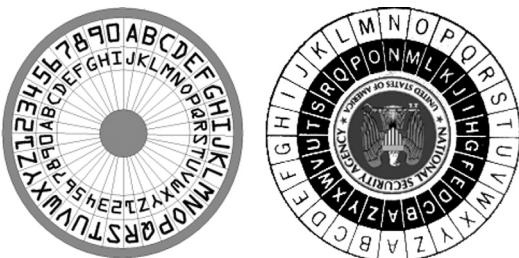


EXHIBIT 7.4 Code Wheels and the NSA Seal

This type of cipher is known as a *substitution cipher*. Although the Caesar cipher is relatively simple, substitution ciphers can be very powerful. Most examples of the Caesar cipher shift the alphabet three places, as shown, so that the ciphertext line begins with *d*, but some authors suggest Caesar might have used other numbers, so the term “Caesar cipher” is used for all ciphers that conform to this algorithm (an *algorithm* being a formula or recipe for solving a problem).

This level of encryption might seem rudimentary, but it is an important starting point for much that follows. For example, one way to visualize the Caesar cipher is as a pair of rings, one inside the other, as shown in Exhibit 7.4. Both circles contain the letters of the alphabet. If one is rotated relative to the other, the result is a cipher wheel, something well suited to automation. Eventually this happened, at first mechanically, then electrically, and today digitally. Automation facilitates repetition, and messages encrypted with a substitution cipher can be more difficult to decipher if multiple different substitutions are used. Thus, the code wheel earned a place in the seal of the NSA, the U.S. government agency most influential in the development of encryption.

7.2.2 More Cryptic Terminology. The key or password for the Caesar cipher presented in the last section is the number of places the alphabet has been shifted, in this case three. Because this key must be kept private in order for the message to remain protected, it must be delivered to the recipient for the message to be decoded, or decrypted, back to plaintext. That is why the Caesar cipher is described as a *private key* algorithm and also a *symmetrical encryption* algorithm, the same private key being used to encrypt and decrypt the message. Algorithms of this type can be defeated by someone who has the key, an encrypted message, and knowledge of the algorithm used. This might sound like a statement of the obvious; however, as will be seen later in this chapter, there are encryption algorithms that use keys that can be openly exchanged without rendering the encrypted data accessible. Knowledge of the algorithm used can often be derived, or reverse-engineered, by analysis of its output.

Another seemingly obvious fact is that when a private key cipher is used in an effort to achieve confidentiality, one problem is swapped for another. The problem of exchanging messages while keeping the contents from unintended recipients is replaced by the problem of exchanging keys between sender and receiver without disclosing the keys. This new problem is known as the *key-exchange problem*. The key-exchange problem will be examined in more detail later.

7.2.3 Basic Cryptanalysis. “The first people to understand clearly the principles of cryptography and to elucidate the beginnings of cryptanalysis were the Arabs.”¹² By the fifteenth century, they had discovered the technique of letter frequency distribution analysis and had successfully decrypted a Greek message on its way to the

BASIC CRYPTOGRAPHY 7 · 9

Byzantine Emperor.¹³ In 1492, a man known as al-Kalka-shandi described this technique in an encyclopedia. He also described several cryptographic techniques, including substitution and transposition ciphers.¹⁴

Returning to the Caesar cipher, consider how this code could be broken using the science of cryptanalysis. When examined for a length of time, this particular code is fairly transparent. As soon as several letters are identified correctly, the rest fall into place. For example, because “the” is the most common three-letter word in the English language, testing “XLI” against “the” reveals that each letter of plaintext has a fixed relationship to the ciphertext: a shift of three to the right.

If that difference is applied to the rest of the message, the result is a piece of plaintext that is intelligible and thus assumed to be the correct solution to the problem. However, even in this simple example several sophisticated processes and assumptions are at work; they deserve closer attention before looking at more complex codes. First, the test of “the” against “XLI” assumes that the plaintext is English and that the attacker has some detailed knowledge of that language, such as the frequency of certain words. Second, it is assumed that the ciphertext follows the plaintext in terms of word breaks. Typically, this is not the case. Ciphertext usually is written in blocks of letters of equal length to further disguise it, as in:

Ciphertext: EHZDU HWKHL GHVRI PDUFK

When the recipient of the message decrypts it, the result, while not exactly easy reading, is nevertheless entirely intelligible:

Plaintext: bewar ethei desof march

Also note the convention of ignoring the case of individual letters and placing all plaintext in lowercase while all ciphertext is in capitals.

7.2.4 Brute Force Cryptanalysis. The next thing to note about the Caesar cipher is that, using the English alphabet, there are 26 possible keys. This means that someone intercepting the encrypted message could mount a standard form of attack known as *brute force cryptanalysis*. This method runs possible keys through the decryption algorithm until a solution is discovered. Statistically speaking, the correct key is reached after testing only half of all possible keys. In Exhibit 7.5, a spreadsheet table details a brute force attack on the Caesar ciphertext. In the example, the plaintext appears in line 6, Key #3.

Note that three items of information are required for this attack, and all three of them are relevant to encryption on personal computers:

1. A knowledge of the encryption algorithm used
2. The number of possible keys
3. The language of the plaintext

Using a computer in an office is somewhat different from sending messages on the field of battle (at least on a good day). Unlike an enemy spy, someone who is attempting to gain unauthorized access to data already has a fairly good idea of which algorithm is being used. (There are relatively few in use, and they often are directly associated with particular applications). This takes care of the first item. The primary obstacle to a brute force attack is the second item, number of keys. In the case of the Caesar

7 · 10 ENCRYPTION

	B6	$=IF(CODE(B$2)-$A6<65,CHAR(CODE(B$2)-$A6+26+32),CHAR(CODE(B$2)-$A6+32))$																									
		CRYPTO.XLS:2																									
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	
1	E	H	Z	D	U	H	W	K	H	L	G	H	V	R	I	P	D	U	F	K							
2	Key#:																										
3	1	d	g	y	c	t	g	v	j	g	k	f	g	u	q	h	o	c	t	e	j						
4	2	c	f	x	b	s	f	u	i	f	j	e	f	t	p	g	n	b	s	d	i						
5	3	b	e	w	a	r	e	t	h	e	i	d	e	s	o	f	m	a	r	c	h						
6	4	a	d	v	z	q	d	s	g	d	h	c	d	r	n	e	l	z	q	b	g						
7	5	z	c	u	y	p	c	r	f	c	g	b	c	q	m	d	k	y	p	a	f						
8	6	y	b	t	x	o	b	q	e	b	f	a	b	p	l	c	j	x	o	z	e						
9	7	x	a	s	w	n	a	p	d	a	e	z	a	o	k	b	i	w	n	y	d						
10	8	w	z	r	v	m	z	o	c	z	d	y	z	n	j	a	h	v	m	x	c						
11	9	y	q	u	l	y	n	b	y	c	x	y	m	i	z	g	u	l	w	b							
12	10	u	x	p	t	k	x	m	a	x	b	w	x	l	h	y	f	t	k	v	a						
13	11	t	w	o	s	j	w	l	z	w	a	v	w	k	g	x	e	s	j	u	z						
14	12	s	v	n	r	i	v	k	y	v	z	u	v	j	f	w	d	r	i	t	y						
15	13	r	u	m	q	h	u	j	x	u	y	t	u	i	e	v	c	q	h	s	x						
16	14	q	t	l	p	g	t	i	w	t	x	s	t	h	d	u	b	p	g	r	w						
17	15	p	s	k	o	f	s	h	v	s	w	r	s	g	c	t	a	o	f	q	v						
18	16	o	r	j	n	e	r	g	u	r	v	q	r	f	b	s	z	n	e	p	u						
19	17	n	q	i	m	d	q	f	t	q	u	p	q	e	a	r	y	m	d	o	t						
20	18	m	p	h	l	c	p	e	s	p	t	o	p	d	z	q	x	l	c	n	s						
21	19	l	o	g	k	b	o	d	r	o	s	n	o	c	y	p	w	k	b	m	r						
22	20	k	n	f	j	a	n	c	q	n	r	m	n	b	x	o	v	j	a	l	q						
23	21	j	m	e	i	z	m	b	p	m	q	l	m	a	w	n	u	i	z	k	p						
24	22	i	l	d	h	y	l	a	o	l	p	k	l	z	v	m	t	h	y	j	o						
25	23	h	k	c	g	x	k	z	n	k	o	j	k	y	u	l	s	g	x	i	n						
26	24	g	j	b	f	w	j	y	m	j	i	j	x	t	k	r	f	w	h	m							
27	25	f	i	a	e	v	i	x	l	i	m	h	i	w	s	j	q	e	v	g	l						
28	26	e	h	z	d	u	h	w	k	h	l	g	h	v	r	i	p	d	u	f	k						

EXHIBIT 7.5 Brute Force Attack on the Caesar Cipher

cipher, the number of possible keys is relatively small, so the work involved in carrying out the attack can be completed very quickly, which is highly significant. Time is often the most important factor in practical cryptanalysis. Being able to decrypt messages within 24 hours is of little use if the information pertains to events that are measured in minutes, such as orders to buy and sell stock, or to launch air raids. If the cipher consisted entirely of random letter substitutions, like this:

Plaintext: abcdefghijklmnopqrstuvwxyz
Ciphertext: UTWFRAQOYSEDCKJVXGZIPHLNM

The number of possible keys (the *keyspace*) is now $26!$, or $\sim 4.03 \times 10^{26}$, which looks even more daunting when written out:

403,291,461,126,606,000,000,000,000

Imagine a brute force attack using a computer that can perform 1 million decryptions per microsecond (considerably more number crunching than the average personal computer can perform). Using a single processor, it could take over 10 million years to execute a brute force attack on this code. Fortunately for the code breaker, there are other ways of cracking substitution ciphers, as discussed in a moment. The point is that, while brute force attacks are possible, they are not always practical.

Although it is true that by the central limit theorem of statistics, the most likely number of trials required to hit on the correct key is one-half the total keyspace, the average reduction by a factor of 2 is negligible in the face of computational periods measured in years and the difficulty of identifying cleartext in the morass of incorrect decryptions.

BASIC CRYPTOGRAPHY 7 · 11

Functionally, brute force attacks depend on knowing which encryption algorithm is behind the ciphertext. Practically, they depend on the feasibility of successes within an appropriate time frame. They also depend on the third item of information in the list above: knowledge of the language of the plaintext. The solution to the Caesar cipher in Exhibit 7.5 tends to jump out because it is closer to plain English than any of the other solutions. However, without knowing what constitutes plaintext, a brute force attack will, at best, be inefficient, and, at worst, unsuccessful. This part of cryptanalysis, recognizing a positive result, is less amenable to automation than any other. The difficulty is compounded by encryption of purely numerical results where the correct cleartext can be impossible to determine without extensive additional knowledge.

7.2.5 Monoalphabetic Substitution Ciphers. Both the Caesar cipher and the random substitution cipher shown are examples of monoalphabetic ciphers. This means that one letter of ciphertext stands for one letter of plaintext. This renders such codes susceptible to an attack quite different from brute force. Suppose a customs officer attempts to discover when and how an illegal weapons shipment will be entering the country. The following message is intercepted:

YZYGJ KZORZ OYXZR RKZRK XUXRJ XRZXU YKQQQ

The person who encoded this text clearly substituted new letters for the original letters of the message. To the experienced code breaker or cryptanalyst, the task of deciphering this message is quite a simple one. First count how many times each letter occurs in the text. This produces a list like this:

Ciphertext:	R	Z	X	Y	K	J	U	O	G
Frequency:	6	6	5	4	4	2	2	2	1

Note that the last three letters are discounted as they are merely filling out the five-letter grouping. Next refer to a table of frequencies, which shows the relative frequency with which the letters of the alphabet occur in a specific language or dialect of that language. One such list is shown in Exhibit 7.6. This list was created for this example and proposes that the most commonly used letters in English in descending order of frequency are *e, t, r, i, n, o, s, h, a, d, l, u*, the order of keys on the English Linotype machine from the nineteenth century, although the precise order of frequencies can vary according to the region of origin or subject matter of the text.

Assuming that the original message is in English, a list that matches code letters to plaintext letters is easily derived.

Ciphertext:	R	Z	X	Y	K	J	U	O	G
Frequency:	6	6	5	4	4	2	2	2	1
Plaintext:	e	t	r	i	n	o	a	h	s

The result is:

Ciphertext: YZYGJ KZORZ OYXZR RKZRK XUXRJ XRZXU YKQQQ
Plaintext: itiso nthet hirte enten rareo retrra inqqq

This is readable as “it is on the thirteen ten rare ore train.” Although this example obviously was contrived to make a point, it clearly illustrates an important cryptographic

7 · 12 ENCRYPTION

EXHIBIT 7.6 Frequency Lists for English

English by Letter				English by Frequency			
A	7.25	N	7.75	E	12.75	U	3.00
B	1.25	O	7.50	T	9.25	M	2.75
C	3.50	P	2.75	R	8.50	P	2.75
D	4.25	Q	0.50	I	7.75	Y	2.25
E	12.75	R	8.50	N	7.75	G	2.00
F	3.00	S	6.00	O	7.50	V	1.50
G	2.00	T	9.25	A	7.25	W	1.50
H	3.50	U	3.00	S	6.00	B	1.25
I	7.75	V	1.50	D	4.25	K	0.50
J	0.25	W	1.50	L	3.75	Q	0.50
K	0.50	X	0.50	C	3.50	X	0.50
L	3.75	Y	2.25	H	3.50	J	0.25
M	2.75	Z	0.25	F	3.00	Z	0.25

tool that can quickly decipher something that at first seems to be very forbidding. The encryption in the previous example could have been based on a simple substitution cipher. For example, after using the password “TRICK” followed by the regular alphabet minus the letters in the password for the plaintext, the ciphertext is the alphabet written backward:

Plaintext: TRICKABDEFGHJLMNOPQSUVWXYZ
Ciphertext: ZYXWVUTSRQPONMLKJIHGFCBA

Frequency analysis also works if the substitution is entirely random, as in the example shown earlier, the key for which is entirely random. The specialized tools, such as frequency tables, that are required to break codes point out a basic trade-off: If a basic level of protection is needed, it is easy to get but also easy to break, at least for an expert. The qualification “for an expert” is important because users of encryption need to keep its role in perspective. The salient questions are: Who can gain from decrypting the data, and what means do they have at their disposal? There is no point investing in powerful encryption hardware or software if those likely to attempt to read your files are not particularly sophisticated, dedicated, or well equipped. For example, a person who mails a postcard knows it can be read by anyone who sees it. Envelopes can be used to prevent this, hardly the ultimate in confidentiality, but widely used and relatively successful nonetheless.

7.2.6 Polyalphabetical Substitution Ciphers. Even when the plaintext uses a wider range of letters than the contrived example, substitution ciphers can be cracked by frequency analysis. A powerful technique is to concentrate on the frequency of two-letter combinations, which are known as *digraphs*, the most common of which in English is “TH.” One way to counter frequency analysis is to use multiple substitutes for the more frequent letters. This cannot be done with a straightforward alphabetic coding. However, if using numbers for letters, it is possible to assign multiple numbers to some letters, such as 13 17 19 23 for E, which would help dilute the natural frequency of this letter. It would appear that supplying multiple substitutions, known as *homophones*, in proportion to the frequency of each letter would effectively

BASIC CRYPTOGRAPHY 7 · 13

counter frequency analysis. However, some of the underlying structure of the plaintext still survives, notably digraphs, which the cryptanalyst can use to crack the code.

In Europe during the Middle Ages, advances in cryptography were being made by the Papal States and Italian city-states to protect diplomatic messages. Then, in 1379, an Italian man named Gabriele de Lavinde created the first European manual on cryptography. “This manual, now in the Vatican archives, contains a set of keys for 24 correspondents and embraces symbols for letters, nulls, and several two-character code equivalents for words and names.”¹⁵ The nomenclature described by Lavinde’s manual “was to hold sway over all Europe and America for the next 450 years.”¹⁶

Several other notable advances emerged in Europe during the period of Lavinde’s manual. First, in 1470, Leon Battista Alberti published the first description of a cipher disk.¹⁷ Next, in 1563, Giambattista della Porta provided the first example of a digraphic cipher in which two letters are represented by one symbol.¹⁸

One method of decreasing the extent to which the structure of the plaintext is reflected in the ciphertext is to encrypt multiple letters of the plaintext. For example, “AR” might be encrypted as “CM.” This is the theory behind what is known as the Playfair cipher, which was invented in 1854 by a British scientist, Sir Charles Wheatstone, but that was named after his friend Baron Playfair who fought for its adoption by the British Foreign Office.¹⁹ Although the Playfair cipher remained in use through both world wars, it does not do enough to disguise the plaintext and cannot withstand a concerted frequency analysis.

7.2.7 The Vigenère Cipher. A particularly important technique in the evolution of polyalphabetic ciphers has its roots in the sixteenth century. In 1586, Blaise de Vigenère published a square encryption/decryption table, named after him as the Vigenère Square, and descriptions of the first plaintext and ciphertext autokey systems.²⁰ The Vigenère cipher involves a table of letters, like the one shown in Exhibit 7.7, that are used with a key to provide different monoalphabetic substitutions as the encryption proceeds through the plaintext. Thus, each letter of the ciphertext has a different relationship with the plaintext, like this:

Key:	doomsdaydoomsdaydoomsdaydoomsday
plaintext:	sellentireportfolionowandbuygold
ciphertext:	VSZXWQTGJIAZVGYWCMDBFPPBOJIKUKLQ

The message is enciphered by looking at the row in the table that begins with the first letter of the key. Then go along that row until the column headed by the first letter of the plaintext. The ciphertext substitution is the letter at that intersection in the table. Thus, row d, column s, yields V. Then proceed to the second letter, and so on. Note that the first time the letter e is encrypted the cipher is S, but the second time it is W. The two ls in sell are encoded as Z and X, respectively, and so on.

Does this cipher completely obscure the structure of the plaintext? Stallings notes: “If two identical sequences of plaintext letters occur at a distance that is an integer multiple of the keyword length, they will generate identical ciphertext sequences.”²¹ This means that the cryptanalyst can determine the length of the keyword. Once this is done, the cipher can be treated as a number of monoalphabetic substitutions, that number being equal to the key length. Frequency tables are again brought into play, and the code can be cracked. The cryptographer’s response to this weakness is to use a longer key so that it repeats less often. In fact, one technique, *autokey*, invented by

7 · 14 ENCRYPTION

a	b	c	d	e	f	G	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x	y	z	
a	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
b	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	
c	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z		
d	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	
e	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	
f	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	
g	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	
h	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	
i	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	
j	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	
k	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	
l	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	
m	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	
n	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	
o	O	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	
p	P	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	
q	Q	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	
r	R	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	
s	S	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	
t	T	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	
u	U	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	
v	V	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	
w	W	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	
x	X	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	
y	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	Q	
z	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	Y	

EXHIBIT 7.7 Vigenère Table

Vigenère, is to form the key from the plaintext itself, together with one code word, like this

Key: doomsdaysellentireportfolionowan
plaintext: sellentireportfolionowandbuygold
ciphertext: VSZXWQTGJIAZVGYWCMDBFPFBOJILUKLQ

This system is very powerful, but it still can be attacked by statistical analysis based on frequencies, because the letters of the plaintext and key share roughly the same frequency distribution. The next level of defense is to use a keyword that is as long as the plaintext but bears no statistical relationship to it. This approach, which is of great cryptographic significance, was not hit upon until the twentieth century arrived, bringing with it binary code and global warfare.

7.2.8 Early-Twentieth-Century Cryptanalysis. The advent of modern cryptography began with the invention and development of the electromagnetic telegraph system and the introduction of the Morse code. Samuel Morse brought a system of dots and dashes that allowed near real-time long-distance communication. He envisioned this system as a means of secure communications. It would be up to others to devise systems to encrypt telegraphic communications. Anson Stager, the supervisor

BASIC CRYPTOGRAPHY 7 · 15

of the U.S. Military Telegraph during the Civil War, devised 10 ciphers for the Union Army that were never broken by the Confederacy.²²

The use of telegraphic ciphers and codes continued into the two world wars. In fact, one of the most famous early successes of cryptanalysis prompted the entrance of the United States into World War I. When the war first started, the German transatlantic telegraph cable had been cut by the British, forcing all of Germany's international communications to route through the United Kingdom before being sent on to the Swedish or American transatlantic lines.²³ In 1917, "British cryptographers deciphered a telegram from German Foreign Minister Arthur Zimmermann to the German Minister to Mexico, von Eckhardt."²⁴ It promised Mexico ownership over territory that belonged to the United States (e.g., California), if Mexico joined the German cause and attacked the United States. The British informed President Wilson of their discovery, giving him a complete copy of the telegram, thus resulting in the United States declaring war on Germany.²⁵ That telegram has become famous in the history of cryptanalysis as the Zimmerman Telegram.

World War II saw several Allied victories over the Axis powers by use of advanced cryptographic systems. Few of these victories are more widely known and celebrated than the cracking of the German Enigma cipher machine, described next:

Following the decryption of the Zimmerman Telegram during World War I and the effects that weak ciphers had on that war's outcome, Germany was looking for an unbreakable cipher and was interested in leveraging automation and the use of machinery to replace traditional paper and pencil techniques. The Enigma machine consisted of a basic keyboard, a display that would reveal the cipher text letter, and a scrambling mechanism such that each plain text letter entered as input via the keyboard was transcribed to its corresponding cipher text letter. The machine was modular in design and multiple scrambling disks were employed to thwart attempts at frequency analysis.²⁶

A British cryptanalysis group, with the help of a group of Polish cryptanalysts, first broke the Enigma early in World War II, and some of the first uses of computers were for decoding Enigma ciphers intercepted from the Germans. Breaking Enigma was a major victory for the Allies, and in order to keep exploiting it, they kept the fact that they had cracked it a secret.²⁷

Thus far, the encryption schemes or devices described have encrypted messages consisting of words and nothing more. However, the emergence of the computer, even in its initial rudimentary form, revolutionized cryptology "to an extent even greater than the telegraph or radio."²⁸ Most cryptologic advances since World War II have involved, or made use of, computers. In the last few decades, cryptographic algorithms have advanced to the point where computing them by hand would be unfeasible, and only computers can do the required mathematics.²⁹ Relying on computers has broadened the kind of information that can benefit from encryption. Computers use a unique language that transforms all information stored into bits, each a 1 or a 0.³⁰ "This, in effect, means that plaintext is binary in form, and can therefore be anything; a picture, a voice, an email or even a video—it makes no difference, a string of binary bits can represent any of these."³¹

7.2.9 Adding up XOR. In 1917, an engineer at AT&T, Gilbert Vernam, was working on a project to protect telegraph transmissions from the enemy. At that time, teletypewriters were used, based on a version of Morse code called Baudot code, after its French inventor. In Baudot code, each character of the alphabet is allotted five units, each of which is either an electrical current or absence of current, known as a

7 · 16 ENCRYPTION

mark or a space. For example, the letter *a* is represented by *mark, mark, space, space, space*. In binary terms, each unit constitutes a bit that is either 0 or 1 (the five-bit code for *a* would be 11000). This system of pulses allowed teletype machines to convert text to and from telegraph signals using a keyboard and punched paper tape for input (a hole represents a mark because it allows the reading device to make electrical contact and create a pulse, whereas a space is represented by leaving the paper intact). Anyone with a suitable machine could intercept and read the transmission.

The 32 possible combinations (2^5) in this code were assigned to the 26 letters plus six “shunts” that did various things like shift to capitals or go down to the next line. Vernam’s brilliant idea was to use a tape of random characters in Baudot code as a key that could be electromechanically added to the plaintext. Kahn describes the method of addition like this:

If the key and plaintext pulses are both marks or both spaces, the ciphertext pulse will be a space. If the key pulse is a space and the plaintext pulse is a mark, or vice-versa (in other words, if the two are different), the ciphertext will be a mark.³¹

Today, this is known as Exclusive-Or, sometimes referred to as bit-wise XOR or just XOR for short (see Exhibit 7.8). XOR is widely used in computerized encryption schemes. Consider what happens when encoding the letter *a* using *B* as the key:

Plaintext:	1 1 0 0 0 (=a)
Key:	1 0 0 1 1 (=B)
Ciphertext:	0 1 0 1 1

In the first column, $1 + 1 = 0$, as indicated in Exhibit 7.8. To decipher the encrypted character, simply perform the same operation, but add the ciphertext to the key:

Ciphertext:	0 1 0 1 1
Key:	1 0 0 1 1 (=B)
Plaintext:	1 1 0 0 0 (=a)

At the time of its discovery, the significance of this method lay in its capacity for automation. The operator could feed the plaintext and key tapes into the teletype machine, and it would transmit an encrypted message with no further human input. No offline preparation was required. Furthermore, as long as the receiver had the key tape, the teletype at the receiving end automatically printed out plaintext. This made Vernam’s system the first to integrate encryption into the communication process, an essential feature of encryption systems for today’s computer-based communications.

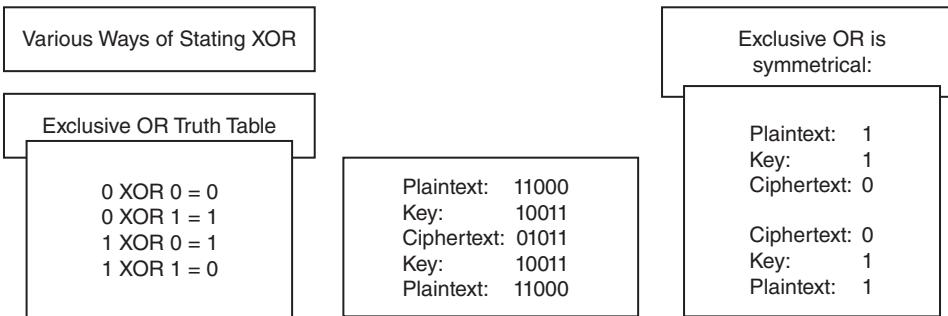


EXHIBIT 7.8 Diagram of XOR

DES AND MODERN ENCRYPTION 7 · 17

7.3 DES AND MODERN ENCRYPTION. Although the use of XOR predated computers, the fact that it worked so well with binary code ensured that it would become an essential item in the modern cryptographer’s toolkit. And so the focus of this chapter turns to modern cryptography and two of the most widely used cryptosystems today. The first is Data Encryption Standard (DES) and the second is Rivest, Shamir, Adleman (RSA).

7.3.1 Real Constraints. As the preceding overview of the evolution of encryption suggests, major advances, which are few and far between, often are linked with the individuals who made them, such as Vigenère, Playfair, and Vernam, none of whom had the benefit of computers. Today’s computerized encryption schemes typically employ a number of classic techniques that, when combined, eliminate or minimize the shortcomings of any single method. Several techniques will be discussed here, including transposition and rotors, that point the way to the most widely used encryption scheme to date: DES. First, however, consider the practical problems encountered by Vernam’s otherwise brilliant scheme.

Vernam proposed a key that was a long series of random characters. This was coded on a loop of paper tape that eventually repeated (the tape held about 125 characters per foot). The length of the key made cryptanalysis of intercepted messages extremely difficult, but not impossible, because eventually the key repeated. With sufficient volume of ciphertext, the code would yield to frequency analysis. (Bear in mind that during time of war, or even military exercises, hundreds of thousands of words may be encrypted per day, providing a solid basis for cryptanalysis.)

7.3.2 One-Time Pad. Several improvements then were suggested to avoid the impracticality of simply creating longer and longer key tapes. Another AT&T engineer, Lyman Morehouse, suggested using two key tapes of about eight feet in length, containing some 1,000 characters, to generate over 999,000 combinations of characters that could be fed into the encryption process as the key. This was an improvement in terms of practicality and security, but, as Major Joseph Mauborgne of the U.S. Army Signal Corps pointed out, heavy message traffic encrypted in this way still could be decoded. It was Mauborgne who realized that the only unbreakable cipher would use keys that are, as Kahn puts it “endless and senseless.”³² Thus he came up with what we know as the one-time system, the one unbreakable encryption scheme.

The one-time system sometimes is referred to as a *one-time pad*,³³ because this is the way it has been deployed by intelligence agents in the field. The agent is issued a pad that aligns columns and rows of entirely random characters, as shown in Exhibit 7.9. The first letter of the plaintext is encrypted using the appropriate ciphertext from row 1, the second letter is encrypted from row 2, and so on. The result is ciphertext that contains no statistical relationship to the plaintext. When the message is encrypted the pad is destroyed. The recipient, who has a copy of the pad, uses it to reverse the process and decrypt the message.

The one-time pad essentially is a polyalphabetic substitution cipher, but with the same number of alphabets as there are characters in the message, thus defeating any kind of frequency analysis. A brute force attack is defeated by the fact that every possible result is as statistically significant as every other. As Kahn points out, a four-letter group of ciphertext could just as easily yield *kiss, fast, slow*, or any other possible four-letter combination.

7 · 18 ENCRYPTION

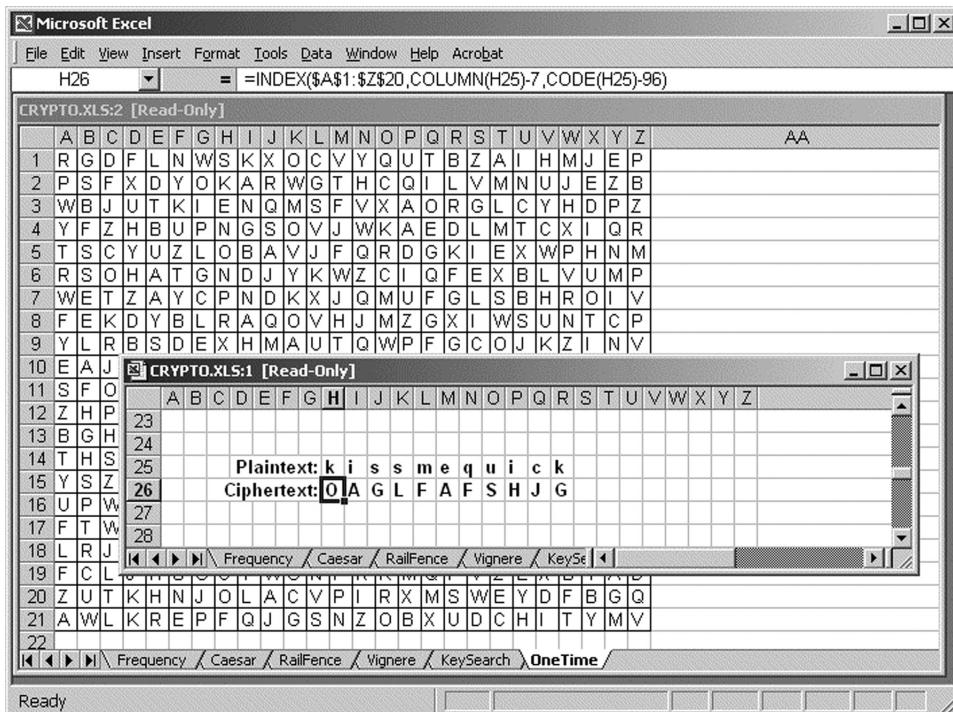


EXHIBIT 7.9 One-Time Pad

So why is the unbreakable one-time system not in universal use? Well, it remains a favorite of intelligence agents in the field who have an occasional need to send short messages. However, for large-scale commercial or military encryption, it fails to solve the key size problem that Vernam's system brought to light. The key has to be as large as the total volume of encrypted information, and there is a constant demand for new keys. Furthermore, both sender and receiver have to hold and defend identical copies of this enormous key.

7.3.3 Transposition, Rotors, Products, and Blocks. A completely different technique from substitution is *transposition*. Instead of substituting ciphertext characters for plaintext, the transposition cipher rearranges the plaintext characters. The simplest example is referred to as *rail fence*. For example, to encrypt “sell entire portfolio now and buy gold” each character is written on alternate lines, like this:

```
sletrprflloodugl
elnieotoinwnbyod
```

which results in this ciphertext:

SLETRPRFLOODUGLELNIEOTOINWNBYOD

So far, this does not present a serious challenge. More challenging is the next transposition into rows and columns that are numbered by a key (in this case, 37581426)

DES AND MODERN ENCRYPTION 7 · 19

so that the first set of ciphertext characters are under 1, the second under 2, and so on:

Key:	3	7	5	8	1	4	2	6
Plaintext:	s	e			e	n	t	i
	r	e	p	o	r	t	f	o
		i	o	n	o	w	a	n
	d	b	u	y	g	o		d

Ciphertext: EROGTFALSRLDNTWOLPOUIONDEEIBLONY

Although more complex, this transposition will still yield to cryptanalysis because it retains the letter frequency characteristics of the plaintext. The analyst also would look for digraphs and trigraphs while playing around with columns and rows of different length. (Kahn describes French code breakers during World War I literally cutting text into strips and sliding them up and down against each other to break German transposition ciphers.)

What makes transposition difficult to decipher is additional stages of encryption. For example, if the previous ciphertext is run through the system again, using the same key, all semblance of pattern seems to disappear.

Key:	3	7	5	8	1	4	2	6
Plaintext:	e	r	o	g	t	f	a	
	s	r		d	n	t	w	o
		p	o	u	i	o	n	d
	e	e	i	b		o	n	y

Ciphertext: TNILAWNNELEFTOOOLOILOYRRPEGDUB

The development of increasingly complex multiple-transposition ciphers pointed out the positive effects of multiple stages of encryption, which also apply to substitution ciphers. The prime examples of this are the rotor machines used by the Germans and Japanese in World War II. Some of the insights gained during the attack on German codes, such as Alan Turing's 1940 work on the application of information statistics to cryptanalysis, were considered so important that they remained classified for more than 50 years.

Although they eventually were defeated by Allied cryptanalysts, electromechanical systems such as Enigma were not only the most sophisticated precomputer encryption systems, but the effort to crack them was also a major catalyst in the development of computer systems themselves. When people started applying computer systems to code making rather than code breaking, they quickly hit on the idea of chopping plaintext into pieces, or blocks, for easier handling. The term "block cipher" is used to describe ciphers that encrypt one block (e.g., 8 bytes of data) at a time, one block after another. Another result of computerizing the encryption process is a class of ciphers known as *product ciphers*. A product cipher has been defined as "a block cipher that iterates several weak operations such as substitution, transposition, modular addition/multiplication [such as XOR], and linear transformation."³⁴

The mathematics of product ciphers are beyond the scope of this chapter, but it is useful to note that "[n]obody knows how to prove mathematically that a product cipher is completely secure . . . [A] product cipher should act as a 'mixing' function

7 · 20 ENCRYPTION

which combines the plaintext, key, and ciphertext in a complex nonlinear fashion.”³⁵ The parts of the product cipher that perform the rounds of substitution are referred to as *S-boxes*. The product cipher called Lucifer has two of these S-boxes, while DES encryption has eight S-boxes. The ability of a product cipher to produce truly random, nonlinear ciphertext depends on careful design of these S-boxes.

Examples of modern product ciphers include Lucifer (developed by IBM), DES (developed by IBM/NSA), LOKI (Brown, Pieprzyk, and Seberry), and FEAL (Shimizu and Miyaguchi). A class of product ciphers called *Feistel ciphers* operates on half of the ciphertext at each round, then swaps the ciphertext halves after each round. Examples of Feistel ciphers include Lucifer and DES, both of which are commercial systems, the subject of the next section of this chapter.

7.3.4 Data Encryption Standard. Traditionally, the primary markets for code makers and computer makers have been the same: governments and banks. After World War II, computers were developed for both military and commercial purposes. By the mid-1960s, the leading computer maker was IBM, which could see that the growing role of electronic communications in commerce would create a huge market for reliable encryption methods. Over a period of years, mathematicians and computer scientists, including Horst Feistel at the IBM research lab in Yorktown Heights, New York, developed a cipher called Lucifer that was sold to Lloyds of London in 1971 for use in a cash-dispensing system.³⁶

The U.S. National Security Agency (NSA) was in close touch with the Lucifer project, making regular visits to the lab (the constant flow of personnel between the NSA, IBM, and the mathematics departments of the major American universities tended to ensure that all new developments in the field were closely monitored). At roughly the same time, the National Bureau of Standards (NBS) was developing standard security specifications for computers used by the federal government. In 1973, the NBS invited companies to submit candidates for an encryption algorithm to be adopted by the government for the storage and transmission of unclassified information. (The government handles a lot of information that is sensitive but not sufficiently relevant to national security to warrant classification.)

IBM submitted a variation of its Lucifer cipher to the NBS, and after extensive testing by the NSA, this cipher was adopted as the nation’s Data Encryption Standard (DES). The acronym actually refers to a document published as Federal Information Processing Standards Publication 46, or FIPS PUB 46 for short. This was published on January 15, 1977, and DES became mandatory for all “federal departments and agencies, for any . . . nonnational-security data.”³⁷ The federal mandate also stated that commercial and private organizations were to be encouraged to use DES.³⁸ As a result, DES became widely used, especially in the banking industry.³⁹ The heart of DES is the Data Encryption Algorithm (DEA), which is described in a publication of the American National Standards Institute, titled *American National Standard for Information Systems—Data Encryption Algorithm—Modes of Operation, 1983*, referred to as ANSI X3.106-1983.

7.3.5 DES Strength. DES became, and remained, the de facto standard for commercial encryption until the late 1990s, when doubts about its strength relative to the rapid advances in computer hardware and software led to a quest for an eventual replacement. However, DES is still widely deployed, so more detailed discussion of its use is needed before discussing its replacement. The first thing to note is that the only

PUBLIC KEY ENCRYPTION 7 · 21

known method of deciphering data encrypted with DES without knowledge of the key is the use of brute force. This involves the computerized comparison of plaintext data with encrypted versions of the same data, using every possible key until both versions of the data match. With DES, the number of possible combinations is about 70 quadrillion. That is a very big number, and trying all those combinations within anything less than years requires relatively expensive hardware (or the carefully orchestrated application of large amounts of cheap hardware).

Technically speaking, the DEA is a combined substitution/transposition cipher, a product cipher that operates on blocks of data 64 bits, or 8 bytes, in length. Using 56 bits for the key produces a keyspace of 2^{56} , or 72,057,594,037,927,940, a number in the region of 70 quadrillion. A diagram of DES is shown in Exhibit 7.10.

The difficulty of attacking DES can be increased fairly easily if double or triple encryption is used, but despite this, there has always been something of a cloud over DES. At the time the DEA was approved, two Stanford University professors who are preeminent in twentieth-century cryptography, Martin Diffie and Whifffield Hellman, pointed out that the algorithm, as approved by the NBS, would be increasingly vulnerable to attack as computer equipment increased in power and came down in cost.

7.3.6 DES Weakness. As the author George Sassoon writes, “Although both the U.S. Department of Commerce and IBM deny it vigorously, everyone in the know insists that the NSA enforced a halving of the DES key length to ensure that they themselves could break the ciphers even if nobody else could.” Although the NBS dismissed such criticisms, and the NSA flatly denied that they were behind any attempts to weaken the cipher, this opinion received some support from the NSA in 1986 when the agency announced it would no longer certify the DEA for nonclassified use, less than 10 years after the DES was approved. This move was prompted by the rapid development of parallel computers, which achieve amazing processing capabilities by using hundreds or even thousands of multiple processors, working in parallel. These machines offer enormous power at considerably less cost than traditional supercomputers. Perhaps the NSA could see the inevitability of something like the EFF DES Cracker, which was built in 1998 for less than \$250,000 and broke a DES-encrypted message in fewer than three days.

The original Lucifer cipher used data blocks of 128 bits and a key of 112 bits. If this had been adhered to in the DEA, the difference in the number of possible key combinations would have been staggering. Although 2^{56} , the current keyspace, is a number greater than 7 with 16 zeroes behind it, 2^{112} is greater than 5 with 33 zeroes behind it. The practical consequence of this weakness in the DEA meant that the demand for stronger algorithms remained, and promising new ones emerged, such as Bruce Schneier’s Blowfish.

There are still some positive aspects to DES that make it viable for some commercial uses. As was mentioned earlier, the cryptographic weakness of DES can easily be strengthened by double encryption, which doubles the difficulty of decryption, taking the task well into the realm of supercomputers and purpose-built, massively parallel machines. The fact that DES has been a standard for so long means that DES now is available in many forms, such as single-chip implementations that can be inserted into ROM sockets and integrated into all manner of hardware, such as expansion cards, PCMCIA cards, and smart cards.

7.4 PUBLIC KEY ENCRYPTION. Even with a longer key, the DEA still would have a major weakness, one that it shares with all of the other private key encryption

7 · 22 ENCRYPTION

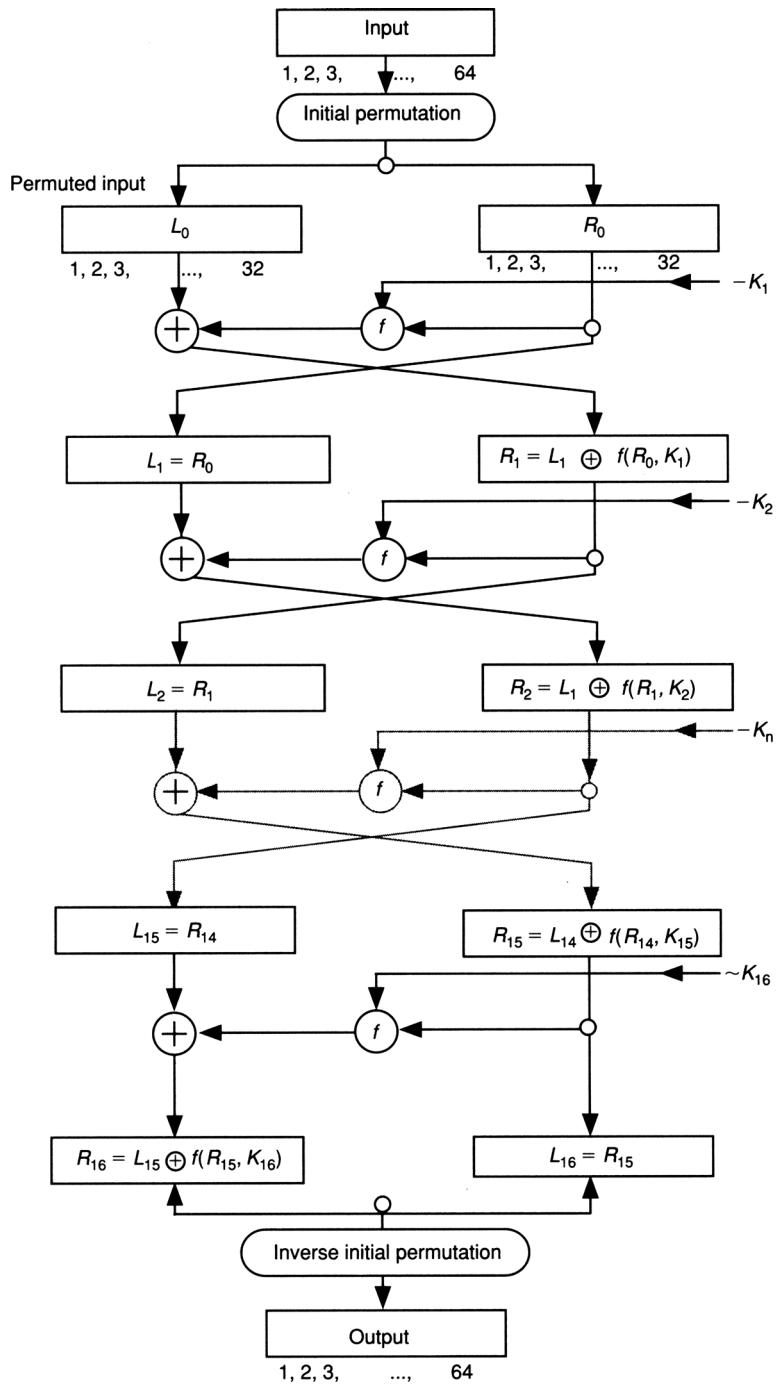


EXHIBIT 7.10 Diagram of DES

PUBLIC KEY ENCRYPTION 7 · 23

systems mentioned so far. That weakness is the need to keep the key secret. In this section we examine this problem, and the “public key” solutions that are now available.

7.4.1 Key-Exchange Problem. When password-protected data are sent from one place to another, either electronically or by hand, the need to transmit the password to the recipient presents serious obstacles. In cryptography, these are known collectively as the *key-exchange problem*. This is the way it is described by the Crypt Cabal⁴⁰:

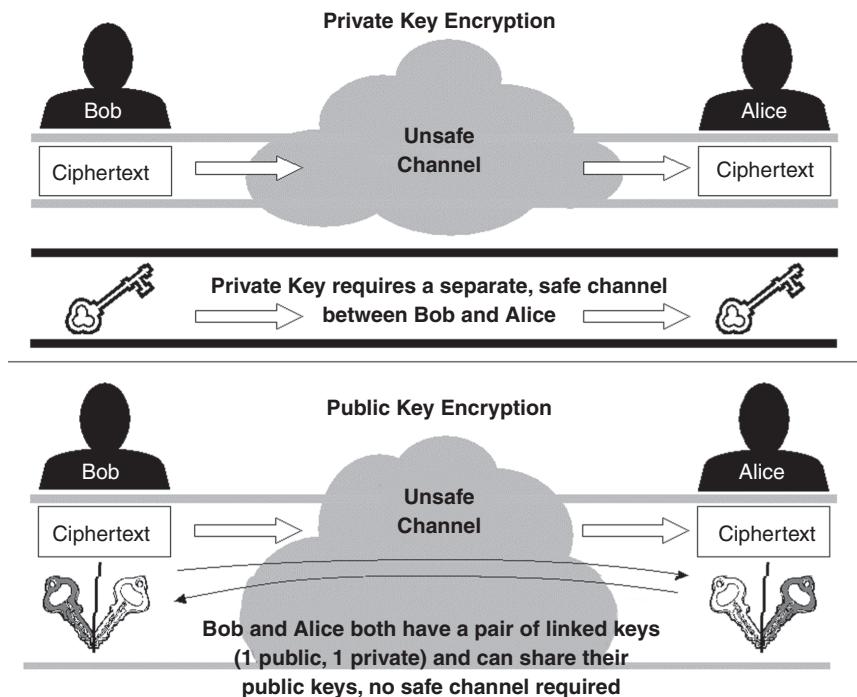
If you want your friends to be able to send secret messages to you, you have to make sure nobody other than them sees the key. . . . [This is] one of the most vexing problems of all prior cryptography: the necessity of establishing a secure channel for the exchange of the key. To establish a secure channel, one uses cryptography, but private-key cryptography requires a secure channel!

So, even when using very powerful private key systems, such as DES, password or key distribution is a major problem. After all, the reason for encrypting valuable information in the first place is because it is assumed someone is trying to steal it or tamper with it. This implies a motivated and skilled adversary. Such an adversary is likely to use every opportunity to discover the password that will unlock the information. The password is perhaps most at risk from such an adversary when it is passed from one person to another. Although it sounds like the stuff of Bond movies, it actually is a very real and practical problem that had to be faced in many areas of legitimate organized activity, from businesses to public institutions, even when a powerful DEA-based computerized encryption system became available.

Suppose an encrypted file of sensitive accounting data needs to get to the head office. How does the recipient know the password needed to access the file? The sender could make a phone call. But will it be overheard? How is the identity of the person at the other end to be verified? A courier could be dispatched with a sealed envelope. The password could be encrypted. But all of these channels present problems. How to guarantee that the courier is honest or that the envelope will arrive intact? And if the password is encrypted, it will need a password itself, which will have to be transmitted. The recipient of the file can be provided with the password before the message is encrypted, but this is no guarantee that the password will not be intercepted. There are ways of making matters more difficult for the attacker, but the ideal solution would be to use a key that was useless to the attacker. This possibility is diagrammed in Exhibit 7.11.

7.4.2 Public Key Systems. A public key encryption system offers encryption that does not depend on the decryption key remaining a secret. It also allows the receiver of keys and messages to verify the source. The first published description of a public key cryptosystem appeared in 1976, authored by Stanford University professor Martin Hellman and researcher Whitfield Diffie. Ralph Merkle independently arrived at a similar system.

Ralph Merkle first proposed the idea of public key cryptography in 1974, and Martin Hellman and Whitfield Diffie brought the same idea to the public forum in 1976.⁴¹ The idea was considered a seminal breakthrough, “for it had not occurred to anyone else in the long history of cryptology that the deciphering key could be anything other than the inverse of the enciphering key.”⁴² The Diffie-Hellman system employs a form of mathematics known as modular arithmetic. “Modular arithmetic is a way of restricting the outcome of basic mathematical operations to a set of integers with

7 · 24 ENCRYPTION**EXHIBIT 7.11** Comparison of Private and Public Key Encryption

an upper bound.⁴³ An excellent example of this mathematical principle is found by examining a military clock:

Consider a clock on military time, by which hours are measured only in the range from zero to 23, with zero corresponding to midnight and 23 to 11 o'clock at night. In this system, an advance of 25 hours on 3 o'clock brings us not to 28 o'clock, but full circle to 4 o'clock (because $25 + 3 = 28$ and $28 - 24 = 4$). In this case, the number 24, an upper bound on operations involving the measurement of hours, is referred to as a modulus. When a calculation involving hours on a clock yields a large number, we subtract the number 24 until we obtain an integer between 0 and 23, a process known as modular reduction. This idea can be extended to moduli of different sizes.⁴⁴

The Diffie-Hellman protocol allows two users to exchange a symmetric key over an unsecure medium without having any prior shared secrets. The protocol has two publicly known and widely distributed system parameters: p , a large prime integer that is 1,024 bits in length,⁴⁵ and g , an integer less than p . The two users wishing to communicate are referred to as Alice and Bob for simplicity's sake. They proceed in this way.

First, Alice generates a random private value a , and Bob generates a random private value b . Both a and b are [less than p]. Then they derive their public values using parameters p and g and their private values. Alice's public value is $g^a \bmod p$ and Bob's public value is $g^b \bmod p$. They then exchange their public values. Finally, Alice computes $g^{ab} = (g^b)^a \bmod p$, and Bob computes $g^{ba} = (g^a)^b \bmod p$. Since $g^{ab} = g^{ba} = k$, Alice and Bob now have a shared secret key k .⁴⁶

This protocol introduced a concept to cryptography known as the discrete log problem. "The discrete log problem is stated as follows: given g , p , and $g^x \bmod p$,

PUBLIC KEY ENCRYPTION 7 · 25

what is x ?⁴⁷ It is generally accepted throughout the mathematical and cryptologic communities that the discrete log problem is difficult to solve, difficult enough for algorithms to rely on it for security.⁴⁸

An algorithm to perform public key encryption was published in 1977 by Ronald Rivest of MIT, Adi Shamir of the Weizmann Institute in Israel, and Leonard Adleman of the University of Southern California. These three men formed the RSA Data Security Company, which was granted an exclusive license to the patent that MIT obtained on their algorithm. A large number of companies licensed software based on this algorithm, from AT&T to IBM and Microsoft. The RSA algorithm is currently at work in everything from online shopping to cell phones. Because it resolved the secret key dilemma, public key cryptography was hailed by many as a revolutionary technology, “representing a breakthrough that makes routine communication encryption practical and potentially ubiquitous,” according to the Sci.Crypt FAQ, which states:

In a public-key cryptosystem, E_K can be easily computed from some public key X, which in turn is computed from K. X is published, so that anyone can encrypt messages. If decryption D_K cannot be easily computed from public key X without knowledge of private key K, but readily with knowledge of K, then only the person who generated K can decrypt messages.⁴⁹

The mathematical principles that make this possible are beyond the scope of this chapter. Somewhat more detail can be found in the RSA Laboratories’ “Frequently Asked Questions About Today’s Cryptography,” which is distributed by RSA Data Security, the company that markets products based on the RSA algorithm. In brief, public key encryption is possible because some calculations are difficult to reverse, something pointed out by Diffie and Hellman, who first published the idea of public key encryption. Here is how RSA describes the calculations that make it possible (with minor clarification from the author):

Suppose Alice wants to send a private message, m , to Bob. Alice creates the ciphertext c by exponentiating:

$$c = m^e \bmod n$$

where e and n are Bob’s public key. To decrypt, Bob also exponentiates:

$$m = c^d \bmod n$$

where d is Bob’s private key. Bob recovers the original message, m ; the relationship between e and d ensures that Bob correctly recovers m . Because only Bob knows d , only Bob can decrypt.

This is diagrammed in Exhibit 7.12, which follows the scenario described. The lower part of the diagram uses numbers taken from an example given by Stallings. These numbers are much smaller than the actual numbers used by RSA. The point is that, given the ciphertext (c) and the public key (e, n) and knowledge of the algorithm, it is still impractical to decipher the message (m). This is because n is created by multiplying two prime numbers (normally represented as p and q) and e is derived from n combined with the secret key, d . To break the cipher, you need to factor a large number into a pair of prime numbers. How large? More than 150 digits in length (that is digits, not bits).

This cryptanalysis is very hard to do in a meaningful period of time, even with a very powerful computer. Large networks of computers have successfully factored

7 · 26 ENCRYPTION

1. Private key, chosen Select two prime numbers, p and q	$p = 7$ and $q = 17$
2. Public key, calculated Calculate $n = pq$	$7 \times 17 = 119$
3. Public key, chosen Calculate $\phi(n) = (p - 1)(q - 1) = 96$ Select e , such that e is relatively prime to $\phi(n)$ and $< \phi(n)$	$e = 5$
4. Private key, calculated Determine d , such that $de = 1 \bmod 96$ and $d < 96$ Because $77 \times 5 = 385 = 4 \times 96 + 1$	$d = 77$
Result; Public key, KU = 5,119 Private key, KR = 77,119	

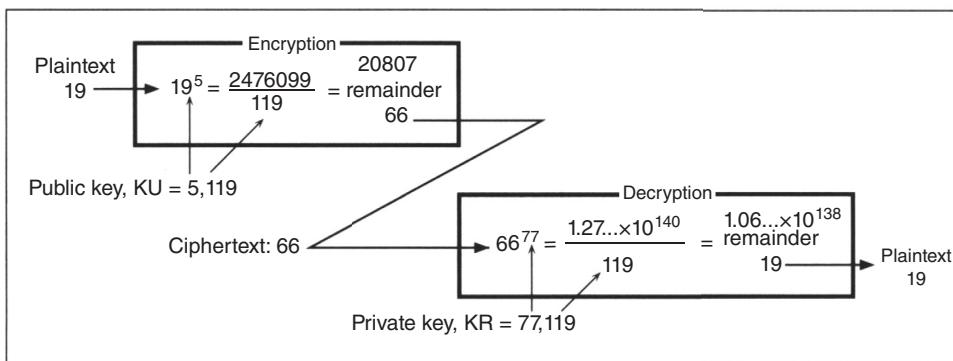


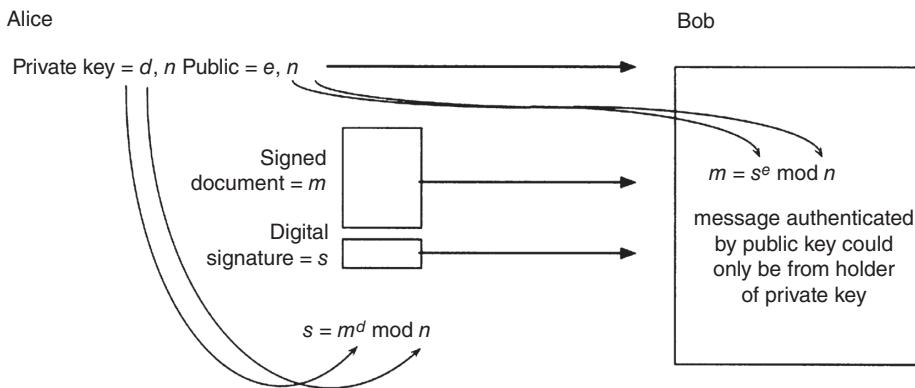
EXHIBIT 7.12 Public Key Diagram

a 100-digit number into two primes, but the RSA algorithm can use numbers even bigger if computer power and factoring algorithms start to catch up to the current implementations.

7.4.3 Authenticity and Trust. The point of the public key cryptosystems is to provide a means of encrypting information that is not compromised by the distribution of passwords, but public key encryption does not solve all problems associated with key exchange. Because the keys are considered public knowledge, some means “must be developed to testify to authenticity, because possession of keys alone (sufficient to encrypt intelligible messages) is no evidence of a particular unique identity of the sender,” according to Sci.Crypt FAQ.⁵⁰

This has led to key-distribution mechanisms that assure listed keys are actually those of the given entities. Such mechanisms rely on a *trusted authority*, which may not actually generate keys but does employ some mechanism which guarantees that “the lists of keys and associated identities kept and advertised for reference by senders and receivers are ‘correct.’ ”⁵¹ Another approach has been popularized by the program called Pretty Good Privacy, or PGP. This is the “Web of trust” approach that relies on users to distribute and track each other’s keys and trust in an informal, distributed fashion.

Here is how RSA can be used to send evidence of the sender’s identity in addition to an encrypted message. First, some information is encrypted with the sender’s private key. This is called the *signature* and is included in the message sent under the public key encryption to the receiver. The receiver can “use the RSA algorithm *in reverse* to

PRACTICAL ENCRYPTION 7 · 27**EXHIBIT 7.13** Authentication with RSA

verify that the information decrypts sensibly, such that only the given entity could have encrypted the plaintext by use of the secret key.”⁵²

What does “decrypts sensibly” mean? The answer involves something called a *message digest*, which is “a unique mathematical ‘summary’ of the secret message.”⁵³ In theory, only the sender of the message could generate his or her valid signature for that message, thereby authenticating it for the receiver. Here is how RSA describes authentication, as diagrammed in Exhibit 7.13.

Suppose Alice wants to send a signed document m to Bob. Alice creates a digital signature s by exponentiating: $s = m^d \text{ mod } n$, where d and n belong to Alice’s key pair. She sends s and m to Bob. To verify the signature, Bob exponentiates and checks that the message m is recovered: $m = s^e \text{ mod } n$, where e and n belong to Alice’s public key.

7.4.4 Limitations and Combinations. As mentioned earlier, many products use RSA today, including Microsoft Windows, Lotus Notes, Adobe Acrobat, Netscape Navigator, Internet Explorer, and many more. In most of these examples, RSA is used for its authentication capabilities rather than for large-scale data encryption. That is because public key systems have one very noticeable downside: They are slow. This is balanced by the fact that they are harder to break. According to RSA, DES generally is at least 100 times as fast as RSA when implemented in software. In hardware, DES is between 1,000 and 10,000 times as fast, depending on the implementations. RSA may narrow the gap in coming years as more specialized chips are developed. However, public key algorithms are unlikely to ever match the performance of private key ciphers such as DES. Fortunately, there is a simple solution: Use a fast private key algorithm for the data encryption, but use a public key system to handle the key exchange and authentication, as diagrammed in Exhibit 7.14.

The private key encryption system might be DES or a system such as RC2 and RC4, both of which are available from RSA Data Security, or Schneier’s Blowfish, which is freely available. Just as there are other private key systems besides DES, there are other public systems besides RSA. One method, called SEEK, is patented, trademarked, and marketed by Cylink of Sunnyvale, California. This method uses an alternative algorithm for public key distribution. Cylink manufactures a range of DES encryptors that use SEEK for key distribution.

7.5 PRACTICAL ENCRYPTION. The primary market for encryption systems and devices is communications. However, the development of Internet commerce has

7 · 28 ENCRYPTION

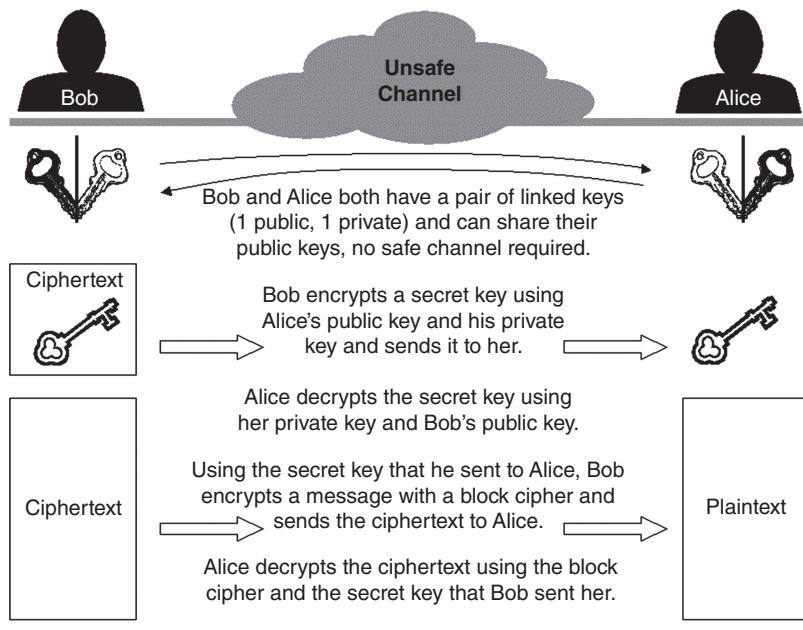


Exhibit 7.14 Combining Public and Private Key Encryption

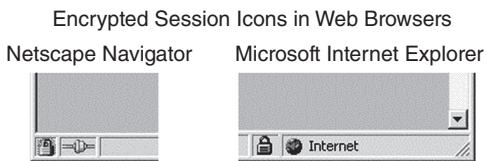
resulted in a number of new and interesting crypto components that have considerable value for computer security.

7.5.1 Communications and Storage. If you look at the commercial products on the National Institute of Standards and Technology (NIST)'s list of approved DES implementations, most are designed to protect information when it is being communicated, not when it is sitting on a machine for local use. This is understandable when you look at the development of computing, which has spread outward from "fortress mainframe." Centralized data storage facilities lend themselves to physical access control. Encrypting data that stays behind walls and locked doors may be overkill in that scenario, particularly when there is a performance penalty involved.

Encryption was reserved for data in transit, between computers, across wires. This philosophy was extended to file servers on networks. File encryption on the server was not considered a priority as people assumed the server would be protected. Data encryption on stand-alone machines and removable media is a relatively recent development, particularly as more and more confidential data are packed into physically smaller and smaller devices. There are now many products with which to implement file encryption.

7.5.2 Securing the Transport Layer. One of the most visible examples of encryption at work in computer security today is the security icon people see in their Web browser; see Exhibit 7.15 for examples of Netscape Navigator and Microsoft Internet Explorer. This is an example of something called transport layer security, which uses protocols that go by the name of SSL and TLS.

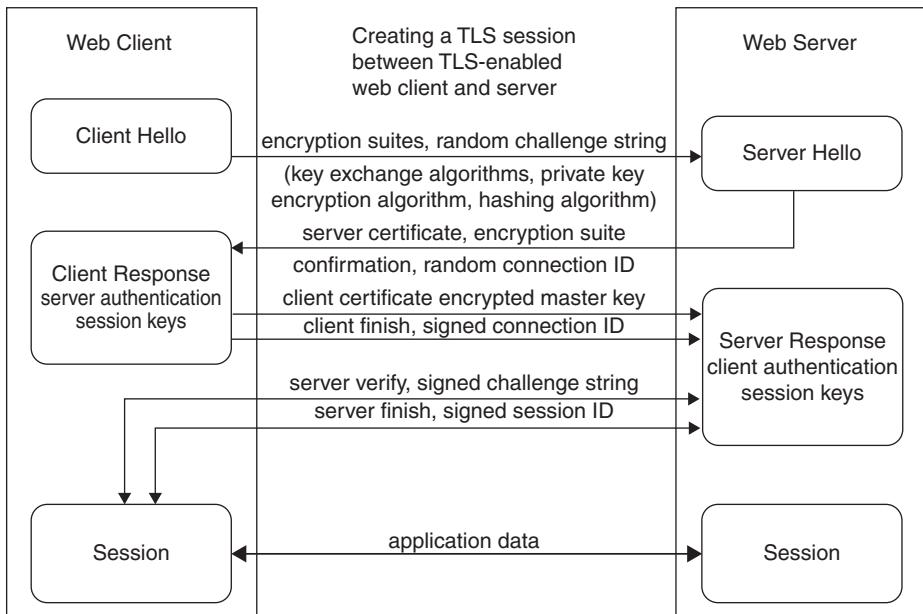
7.5.2.1 Popular Protocols. SSL stand for Secure Sockets Layer, the software encryption protocol developed by Netscape and originally implemented in Netscape

PRACTICAL ENCRYPTION 7 · 29**EXHIBIT 7.15** SSL 3.0 in Action

Secure Server and the Netscape Navigator browser. SSL is also supported by Microsoft Internet Explorer and a number of other products. TLS stands for Transport Layer Security, the name given to an Internet standard based on SSL, by the IETF (as in Internet Engineering Task Force, RFC 2246). There are minor differences between SSLv3.0 and TLSv1.0 but no significant differences as far as security strength is concerned, and both protocols interoperate with each other.

The TLS is a protocol, a standardized procedure for regulating data transmission between computers. It is actually composed of two layers of protocol. At the lowest level is the TLS Record Protocol, which is layered on top of some reliable transport protocol, typically the TCP in TCP/IP, the set of protocols that run the Internet. The TLS Record Protocol provides connection security that is both private (using symmetric cryptography for data encryption) and reliable (using a message integrity check). Above the TLS Record Protocol, encapsulated by it, is the TLS Handshake Protocol. This allows the server and client to authenticate each other, a major role for TLS in various forms of e-commerce, such as Internet banking. The TLS Handshake Protocol can also negotiate an encryption algorithm and cryptographic keys before any application protocol sitting on top of it, such as HTTP, transmits or receives its first byte of data (see Exhibit 7.16).

7.5.2.2 Properties of TLS. In providing connection security, the TLS Handshake Protocol delivers three basic properties. The identity of the parties can be

**EXHIBIT 7.16** Creating a TLS Session

7 · 30 ENCRYPTION

authenticated using public key cryptography (such as RSA). This authentication can be made optional, but typically it is required for at least one of the parties (e.g., the Yahoo! Travel server authenticates itself to the user's browser client, but the user's client does not authenticate itself to the Yahoo! Travel server, a distinction discussed in a moment).

The second and third basic properties of the TLS Handshake Protocol are that a shared secret can be securely negotiated, unavailable to eavesdroppers, even by an attacker who can place itself in the middle of the connection; and the protocol's negotiation is reliable. In the words of RFC 2246: "no attacker can modify the negotiation communication without being detected by the parties to the communication."

TLS can use a variety of encryption algorithms. For the symmetric encryption that is part of the Record protocol, DES or RC4 can be used. The keys for this symmetric encryption are generated uniquely for each connection and are based on a secret negotiated by another protocol (such as the TLS Handshake Protocol). The record protocol includes a message integrity check using a keyed MAC, with secure hash functions such as SHA and MD5, used for MAC computations. The encryption suite to be used for a specific connection is specified during the initial exchange between client and server, as shown in Exhibit 7.16.

7.5.2.3 Tested in the Real World. TLS/SSL has been widely used and extensively tested in the real world, and thoroughly probed by real cryptographers. Some of the caveats and limitations noted by these and other experts follow. The first is that neither a good standard nor a good design can guarantee a good implementation. For example, if TLS is implemented with a weak random number seed, or a random number generator that is not sufficiently random, the theoretical strength of the design will do nothing to protect the data that are thus exposed to potential compromise. (Although beyond the scope of this chapter, Pseudo-Random Number Generators, or PRNGs, play a vital part in many cryptographic operations, and they are surprisingly difficult to create; unless they closely simulate true randomness, an attacker will be able to predict the numbers they generate, thus defeating any scheme that relies on their "random" quality.)

The second major caveat is that, if clients do not have digital certificates, the client side of the TLS session is not authenticated. This presents numerous problems. Most of today's "secure" Web transactions, from airline tickets booked at Yahoo Travel to shares traded at most online brokerages, represent a calculated risk on the part of the vendor. Although the client doing the buying is assured, by means of the merchant certificate, that the merchant at www.amazon.com really is Amazon, the merchant has no digital assurance that the client computer belongs to, or is being operated by, the person making the purchase. Of course, there are other assurances, such as the match between the credit card that the purchaser supplies and the other personal details that go along with it, such as billing address. But the merchant is still risking a charge-back and possibly other penalties for a fraudulent transaction.

In the case of larger and more sensitive financial transactions, the need to be assured of the client's identity is greater. A digital certificate is a step in the right direction, but it is a step many merchants have not yet taken, for several reasons. The first is the cost of issuing certificates to customers, and the second is the difficulty of getting those certificates onto their systems. Some merchants have decided that the cost and effort are worth it. For example, the Royal Bank of Scotland took this approach with its online banking system back in 1998.

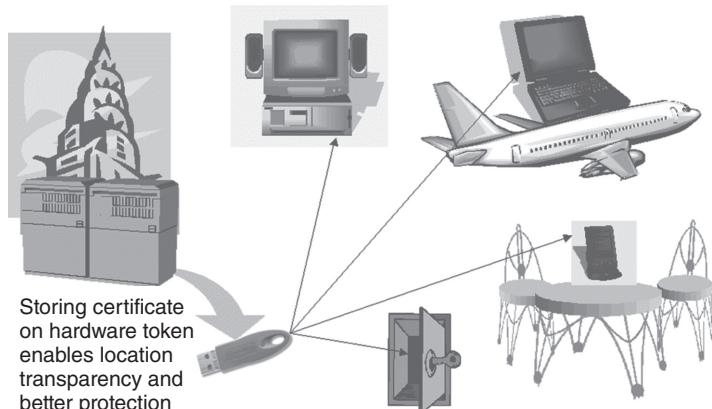
PRACTICAL ENCRYPTION 7 · 31

EXHIBIT 7.17 Using a Hardware Token for Digital Signatures

There are other issues. The user needs to protect the certificate, even from such threats as hardware failure (user reformats the drive, loses the certificate) or unauthorized use (a family member uses the computer and thus has access to the certificate). Furthermore, the user needs to be able to move the certificate, for example, onto a laptop computer so that the bank account can be accessed while traveling. The obvious answer is to place the certificate on a robust removable medium (see Exhibit 7.17). Such media are generically referred to as hardware tokens. A standard for tokens has not yet emerged. Smart cards are an obvious choice, but card readers need to be deployed. There are alternatives, such as putting the certificate on a floppy disk or on a small key fob that plugs into a USB port.

7.5.2.4 Cost of Secured Transactions. For companies looking to perform highly secure transactions today, using SSL without client-side authentication is proving acceptable in the short term, at least for some categories of transaction. Even then it can be costly, in terms of either dollars or processing power. Although TLS is an open standard, and Netscape has provided crucial parts of the technology royalty free, there is still the question of which algorithms to use. Some algorithms are more expensive than others, and not always in obvious ways. For example, you have to license RC4, whereas DES is free, but RC4 is optimized for a 32-bit processor and DES is not.

Furthermore, research shows that the amount of “hits” that a Web server can handle drops dramatically when those hits require TLS (and it drops a whole lot more when processing client authentication as well as server authentication). The answer here may be specialized hardware. Several companies, such as IBM and Rainbow Technologies, make crypto-coprocessor cards that relieve the server’s CPU of the specialized math processing involved in crypto. They are cheaper than adding another server to keep up with the very demanding task of providing secure Web transactions.

7.5.3 X.509v3 Certificate Format. Another example of encryption widely used in computer security today is X.509. This is not a rocket ship but a standard for digital certificates, described earlier in this chapter. The International Telecommunication Union’s Telecommunication Standardization Sector (ITU-T)’s X.509 standards document states: “Virtually all security services are dependent upon the identities of the communicating parties being reliably known, i.e. authentication.” Consider how this affects Web transactions. The preceding section described how SSL can encrypt

7 · 32 ENCRYPTION**EXHIBIT 7.18** Digital Certificate

Web pages sent from Web server to Web client, and vice versa, but it cannot assure the identity of the parties involved. The X.509 standard helps to address this problem, which negatively impacts the profitability of Web-based businesses.

When a Web user asks for assurance that the bn.com Web site is actually Barnes & Noble, it can be provided by way of a digital certificate (see Exhibit 7.18). This means that an entity, known as a certificate authority (CA), has taken considerable pains to reliably identify, and consequently certify, the merchant as the rightful owner of an encryption key. This key is the public half of a uniquely and mathematically related public/private key pair, such that a message encrypted with the public key can only be decrypted with the corresponding private key.

Individuals, as well as merchants, can have a public/private key pair. A bank might then access that public key, and use it, plus the bank's private key, to encrypt the account details it sends to customers over the Web. Only the customer with the right private key can decrypt this information, using the bank's public key. At the same time, customers know the statement information can only have come from the bank (otherwise the bank's public key would not work to decrypt it). Customers also know,

PRACTICAL ENCRYPTION 7 · 33**EXHIBIT 7.19 X.509 Certificate Format**

Version	Identifies the Certificate Format
Certificate Serial Number	Number that is unique within the issuing CA
Signature Algorithm Identifier	Identifies the algorithm used to sign the certificate, together with any necessary parameters
Issuer	X.500 name of the issuing CA
Validity Period	Pair of dates between which the certificate is valid
Subject	X.500 name of the holder of the private key corresponding to the public key certified by the certificate
Subject Public Key Information	Public key for the subject, plus an identifier for the algorithm with which this public key is to be used

thanks to an encrypted message digest (a digital fingerprint of the message contents), that the data they get from the bank has not been altered. Thus, it is very difficult for either party to claim that it never took place. In this way, digital certificates can enhance confidentiality, integrity, and nonrepudiation.

7.5.3.1 ISO/IEC/ITU 9594-8 a.k.a. X.509. The management of public keys is the task of Public Key Infrastructure (PKI), of which the X.509 standard is an important part. For example, an organization's employees can perform secure business communications over the Internet, such as contract negotiation, using PKI. To engage in a secure transaction with someone, it is necessary to find and access the other person's public key, and vice versa. The answer is to publish public keys in the form of a digital certificate, then use some form of directory to locate them. In order for different systems to interoperate, standards for directories have been developed, notably X.500. This standard applies such elements of directory standardization as a hierarchical naming convention:

Country, Organization, Common Name.

So Fred Jones of Megabank might have the X.500 name:

[Country = US, Organization = Megabank, Inc., Common Name = Fred Jones]

A means of locating digital certificates to verify identities was a logical extension of the standard, thus X.509 was developed, officially known as ITU-T X.509 (formerly CCITT X.509) and also ISO/IEC/ITU 9594-8. In X.509 there is a definition of a basic certificate format, which consists of seven fields shown in Exhibit 7.19.

The certificate format has evolved considerably since 1988. The original format is now referred to as X.509v1. When X.500 itself was revised in 1993, two more fields were added to support directory access control, resulting in the X.509v2 format.⁵⁴ X.509v2 added unique identifiers for the issuer and the subject, optional bit strings used to make the issuer and subject names unambiguous in the event that the same name is later reassigned to different entities. Suppose that Fred Jones, whose assigned X.500 name was given earlier, is an executive vice president of Megabank, but is then hired away by a competitor. Megabank deassigns his name, but if a different Fred Jones, a programmer, then comes to work for Megabank, he is effectively reassigned the same X.500 name:

[Country = US, Organization = Megabank, Inc., Common Name = Fred Jones]

7 · 34 ENCRYPTION

This poses authorization problems for any access control lists attached to X.500 data objects, due to the difficulty of identifying all of the access control lists that grant privileges to a particular user's name. The unique identifier field added in X.509v2 provides somewhere to put a new value whenever a name is reused. In fact, a better solution is to use a better distinguisher in the X.500 name, such as:

[Common Name = Fred Jones, Employee Number = 1000000]

In 1993, when the Internet Privacy Enhanced Mail (PEM) RFCs were published, they included specifications for a public key infrastructure based on X.509v1 certificates. Attempts to deploy PEM, however, revealed deficiencies in the Version 1 and 2 certificate formats. Consequently, ISO/IEC/ITU and ANSI X9 developed the X.509v3 format, which greatly extends the capabilities of the format by providing extension fields and broader naming options in X.509v3.

7.5.3.2 Extending the Standard. Extensions were added in Version 3 to address problems discovered while implementing Version 1 and 2 certificates. These can be seen in the diagram in Exhibit 7.20. Particular extension field types may now be specified in standards or defined and registered by any organization or community. Each extension field is assigned a type by means of an object identifier, registered in the same way that an algorithm is registered. Although theoretically anyone can define an extension type, to achieve practical interoperability, common extension types need to be understood by different implementations. Thus, the most important extension types are standardized. But when X509v3 is used within a closed group—for example,

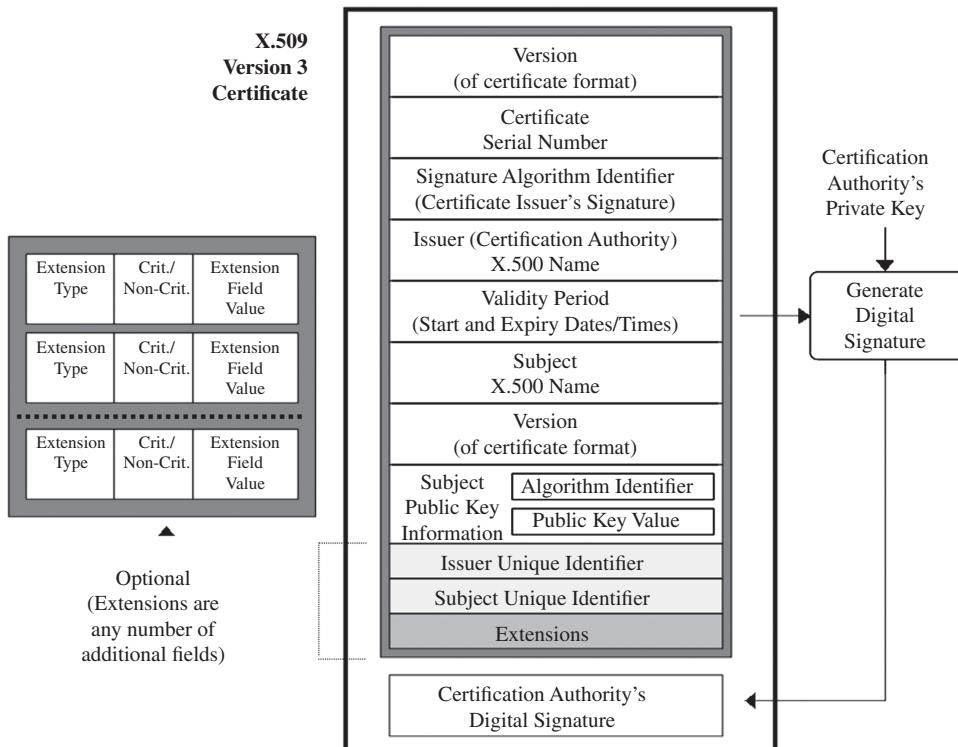


EXHIBIT 7.20 X.509v3 Certificate

PRACTICAL ENCRYPTION 7 · 35

a group of business partners—it is possible to define unique extension types to satisfy specific needs.

7.5.3.3 X.509 Sources, Issues, and CAs. Someone managing an e-commerce project does not necessarily need to know X.509 in detail but should at least read the Arsenault and Turner document (see Section 7.8, “Further Reading,” at the end of the chapter); it clearly describes not only X.509 but the role it plays in PKI (which they define as “the set of hardware, software, people, policies and procedures needed to create, manage, store, distribute, and revoke certificates based on public-key cryptography”⁵⁵). Also very helpful are the presentations by VeriSign’s Warwick Ford, which NIST has online at its Web site. For the e-commerce developer who wants more detail, the next step is Ford’s book, coauthored with fellow VeriSign executive Michael Baum, *Secure Electronic Commerce*.⁵⁶ This documents other important aspects of X.509, such as the Certificate Revocation List, used to revoke certificates before they expire (e.g., if the private key has been compromised). A copy of the standard, available online, at the International Telecommunication Union (ITU) Web site (www.itu.int), is also valuable.

The extensions and improvements in the X.509v3 certificate format greatly increase its usefulness, but providing a uniform method of going beyond the standard does raise the specter of a lack of standardization. This is something that the IETF’s PKIX working group is addressing. And there are other issues to consider when evaluating X.509 as a security technology, many of which are raised by Ed Gerck of the Meta-Certificate Group. Articles at the group’s Web site point out that X.509 does not address “the level of effort which is needed to validate the information in a certificate.”⁵⁷ In other words, some security issues are beyond the scope of X.509, but they do need to be considered when deploying systems that rely on these certificates. For example, it does not make sense to rely on a digital certificate if the measures taken to assure the identity of the owner and user of the certificate are not commensurate with the risk involved in relying on the certificate. Furthermore, transactions that do not use certificates on both sides will remain inherently problematic.

These issues point to the importance of the role played by the CA. As mentioned earlier, CAs are the entities that issue and sign certificates. Each has a public key that is listed in the certificate. The CA is responsible for scheduling an expiration date and for revoking certificates when necessary. The CA maintains and publishes a Certificate Revocation List (CRL).

In other words, ensuring the validity of certificates entails a lot of maintenance. The CRL, for example, is crucial if certificates are compromised or found to be issued fraudulently. This happened in 2001 when a number of VeriSign certificates were found to be issued in error to someone posing as Microsoft. Because some computer users now rely on certificates to guarantee the authenticity of software upgrades and components, failure to check the revocation list before downloading certified code could result in malicious code attacks.

Problems with certificates have the potential for widespread impact because the authority in certificates is hierarchical, as shown in Exhibit 7.21. When a CA issues a certificate, it signs it with its own key. Anyone relying on certificates issued by that CA needs to know by what authority the CA is issuing that certificate. To simplify, there are two possible answers. The CA is self-certifying, that is, providing its own “root” key, or it is relying on another CA for the root key. Clearly, any compromise of the root key undermines all certificates that gain their authority from it.

See Chapter 37 of this *Handbook* for a more extensive discussion of PKI and CAs.

7 · 36 ENCRYPTION

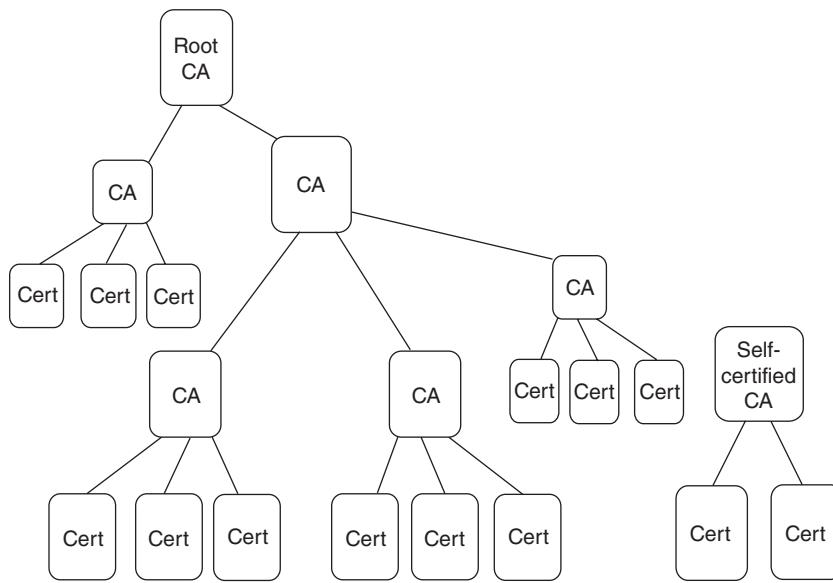


EXHIBIT 7.21 Certificate Authorities and the Root Key

7.6 BEYOND RSA AND DES. Cryptography research and development did not stop with the development of the RSA algorithms. Events in the last two decades of the twentieth century and the first decade of the twenty-first, and their implications, are discussed in this final section of the chapter, which concludes with some warnings on implementing encryption.

7.6.1 Elliptic Curve Cryptography. In 1985, Neal Koblitz from the University of Washington and Victor Miller of IBM independently discovered the application of elliptic curve systems to cryptography. When applied to public key cryptography, elliptic curve arithmetic has been found to offer certain advantages over first-generation public key techniques such as Diffie-Hellman and RSA.

The security of elliptic curve algorithms is based on the same principle as the Diffie-Hellman algorithm, the discrete log problem, as described in Section 7.4.2. The advantages to elliptic curve algorithms lie in the key size needed to achieve certain levels of security. As one scales security upward over time to meet the evolving threat posed by eavesdroppers and hackers with access to greater computing resources, elliptic curves begin to offer dramatic savings over the old, first-generation techniques.⁵⁸

Until 2010, public key systems used 1,024 bits or 2,048 bits for creating keys. NIST recommended that after 2010, these systems be upgraded to a system that can provide adequate security. One way of doing this would be to increase the key size that is used. However, systems that are in place today become increasingly cumbersome the larger the key size. The NSA is endorsing elliptic curve cryptography, stating on its Web site that it has implemented elliptic curve public key cryptography systems to protect both classified and unclassified information.⁵⁹ Elliptic curve systems offer a way to increase key size moderately when more security is required. Exhibit 7.22 shows the NIST recommended key size that RSA or Diffie-Hellman should use to protect the transportation of symmetric keys of various sizes as well as the corresponding elliptic curve key size.

BEYOND RSA AND DES 7 · 37**EXHIBIT 7.22** NIST Recommended Key Sizes

Symmetric Key Size (bits)	RSA and Diffie-Hellman Key Size (bits)	Elliptic Curve Key Size (bits)
80	1024	160
112	2048	224
128	3072	256
192	8192	384
256	15360	521

Source: National Security Agency, "The Case for Elliptic Curve Cryptography," www.nsa.gov/business/programs/elliptic_curve.shtml.

Thus, in order to use RSA to protect a 256-bit AES key, one should use a key of 15,360 bits, which is an order of magnitude larger than the key sizes currently in use throughout the Internet. However, an elliptic curve key would need to be only 521 bits. Elliptic curve algorithms can use smaller keys, because the math involved makes the inverse, or decryption, operations harder as the key length increases.⁶⁰

Another feature that makes elliptic curves appealing is the fact that they are more efficient than the current implementations of public key cryptography, which tend to be relatively slow, causing them to be used more as key distribution methods than data encryption methods. Exhibit 7.23 shows the ratio of Diffie-Hellman computations versus elliptic curve computations for each of the key sizes listed in Exhibit 7.22.⁶¹

7.6.2 RSA Patent Expires. On September 6, 2000, RSA Security released the RSA public key encryption algorithm into the public domain. This means that anyone can now create products that incorporate this algorithm (provided it is their own implementation and not one licensed from RSA). In effect, RSA Security waived its rights to enforce the patent for any development activities that include the RSA algorithm occurring after September 6, 2000. The U.S. patent for the RSA algorithm actually expired on September 20, 2000. The result has been an even broader use of public key encryption, at lower cost.

The RSA patent was always somewhat controversial, because it applied to a piece of mathematics, which is not what most people think of when they think of an invention. The owners of the patent were never able to expand protection beyond the United States. As a result, versions of public key encryption based on alternatives to the RSA algorithm were developed and marketed outside the country, by companies

EXHIBIT 7.23 Relative Computation Costs of Diffie-Hellman and Elliptic Curves

Security Level (bits)	Ratio of DH Cost : EC Cost
80	3:1
112	6:1
128	10:1
192	32:1
256	64:1

Source: National Security Agency, "The Case for Elliptic Curve Cryptography," www.nsa.gov/business/programs/elliptic_curve.shtml.

7 · 38 ENCRYPTION

like Ireland's Baltimore Technologies, Finland's F-Secure, and Israel's Algorithmic Research. Now encryption companies can dispense with the costly maintenance of multiple versions of their public key products (U.S. and non-U.S.). In addition, U.S. companies can develop and market RSA-based products. Large companies actually can "roll their own" public key encryption schemes for internal use, based on a proven, royalty-free algorithm.

7.6.3 DES Superseded. RSA Security, the company that tried to make the RSA algorithm synonymous with public key encryption, played a leading role in the other watershed crypto event of 2000, the naming of a successor to DES, the Data Encryption Standard. As noted earlier, projects like the EFF DES Cracker showed that a computer built for less than \$250,000 could decipher a DES-encrypted message in fewer than three days. In fact, this was part of the "DES Challenges" sponsored by RSA Security. DES Challenge I was won by Rocke Verser of Loveland, Colorado, who led a group of Internet users in a distributed brute force attack. The project, code-named DESCHALL, began on March 13, 1997, and was successfully completed some 90 days later. DES Challenge II consisted of two contests posted on January 13 and July 13, 1998. The first contest was cracked by a distributed computing effort coordinated by distributed.net, which met the challenge in 39 days. The second contest was the one solved by EFF's purpose-built DES Cracker.

The effect of these projects was to focus attention on the need for stronger encryption. Companies and government agencies wanting to archive sensitive data need it to remain secure for decades, not days. However, as predicted in the 1970s, advances in computer power rendered "obsolete" the DEA, the widely used private key algorithm that forms the basis of the DES. Of course, the term "obsolete" is relative in this context. DES is not obsolete when applications need to encrypt bulk data to keep it confidential for a limited period of time, and a lot of data falls into this category. As Exhibit 7.24 shows, there is a direct relationship among time, technology, and the degree of protection that any ciphersystem provides.

In 1997, the U.S. Government began the process of establishing a more powerful standard than DES, known as Advanced Encryption Standard (AES). This is a Federal Information Processing Standard (FIPS) Publication, FIPS 197, specifying "a

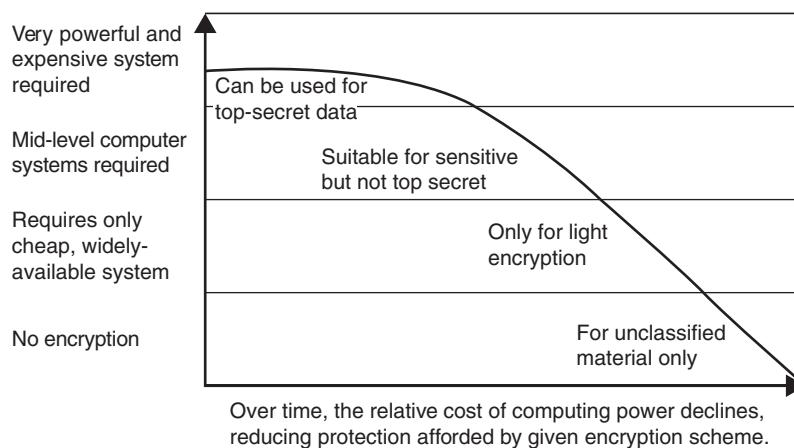


EXHIBIT 7.24 Relationship among Time, Technology, and Protection

BEYOND RSA AND DES 7 · 39

cryptographic algorithm for use by U.S. Government organizations to protect sensitive (unclassified) information.” The government anticipated correctly that AES would be “widely used on a voluntary basis by organizations, institutions, and individuals outside of the U.S. Government—and outside of the United States—in some cases.”

In essence, a competition was held to find the best possible algorithm for the job, and the winner, chosen in October 2000, was Rijndael (pronounced “Rhine Doll”). This algorithm was developed specifically for the AES by two cryptographers from Belgium, Dr. Joan Daemen and Dr. Vincent Rijmen. Rijndael is a block cipher with a variable block length and key length. So far, keys with a length of 128, 192, or 256 bits have been specified to encrypt blocks with a length of 128, 192, or 256 bits. (All nine combinations of key length and block length are possible.) However, both block length and key length can be extended very easily in multiples of 32 bits. Rijndael can be implemented very efficiently in hardware, even on smart cards.

7.6.4 Quantum Cryptography. A new basis for computation will profoundly affect cryptographic strength in the coming decades. This section provides a brief and nontechnical summary of the science of quantum computation and quantum cryptography.

7.6.4.1 Historical Perspective. The entirety of this chapter has focused on the status of cryptography as it currently exists. The classic computer has been sufficient to perform the computations and processes required of AES, RSA, and all of the cryptographic systems and algorithms that have been explored since the advent of cryptography. Although modern computers are fundamentally the same as they were in the 1950s, the machines we use today are significantly faster.⁶² Even though the speed has increased, the primary task of computers has remained the same: “to manipulate and interpret an encoding of binary bits into a useful computational result.”⁶³ To push the bounds of computer performance ever forward, computer scientists’ goal has “been the reduction of size in the transistors used in modern processors.”⁶⁴

Early computers were constructed of gates and storage “bits” made of many thousands of molecules. The components of today’s processors are moving in the direction of a few hundred molecules. The computing industry has always known that miniaturization would reach a barrier below which circuits could not be built, because their fundamental physical behavior would change.⁶⁵

The components of modern computers are reaching this barrier; should transistors become much smaller, they will “eventually reach a point where individual elements would be no larger than a few atoms.”⁶⁶ Computer scientists are concerned about this continual shrinking, because at the atomic level, the laws of quantum mechanics will govern the properties and behavior of circuits, not the laws of classical mechanics.⁶⁷

The science of quantum mechanics is not fully understood by scientists; it was initially thought to be a major limitation to the evolution of computer technology.⁶⁸ It was not until 1982 that the scientific community saw any benefit from the unusual effects associated with quantum mechanics. That year, Richard Feynman theorized about a new type of computer that would harness the effects of quantum mechanics and use these effects to its advantage.⁶⁹ In 1985, David Deutsch of the University of Oxford published a “ground breaking theoretical paper describing how any physical process could be modeled perfectly (in theory) using a quantum computing system.”⁷⁰ He further argued that a quantum system would be able to execute tasks that no modern computer could perform, such as *true* random number generation.⁷¹ “After

7 · 40 ENCRYPTION

Deutsch published this paper, the search began to find interesting applications for such a machine.”⁷²

7.6.4.2 Fundamentals. A “quantum” is “the smallest amount of a physical quantity that can exist independently, especially a discrete quantity of electromagnetic radiation.”⁷³ Quantum mechanics explains the physics and behaviors of particles, atoms, and energy.⁷⁴ The idea of a quantum computer is based on the phenomena that occur at the atomic and subatomic level, which are explained by quantum mechanics and defy all classical laws of physics.⁷⁵ These phenomena will be covered in more detail shortly; it is necessary at this point, however, to explain several fundamental differences between classical modern computers and the idea of a quantum computer.

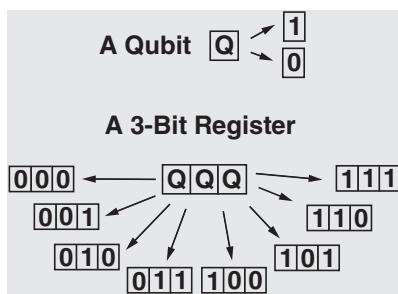
Classical computers store and process information in units called bits, represented as a zero (0) or a one (1) in a computer’s transistors. Bits are then organized into bytes, a series of eight bits. Thus, the information stored on a computer is stored as individual bits grouped into bytes. Therefore, a document “comprised of n-characters stored on the hard drive of a typical computer is accordingly described by a string of $8n$ zeros and ones.”⁷⁶ It is important to emphasize that bits “can *only* exist in one of two distinct states, a ‘0’ or a ‘1.’”⁷⁷ This leads to the first difference between classical computers and quantum computers.

Quantum computers store and process information in units called quantum bits, referred to as “qubits.” “Qubits represent atoms, ions, photons or electrons and their respective control devices that are working together to act as computer memory and a processor.”⁷⁸ Similar to a classic bit, a qubit is represented as a 0 or a 1. Unlike a classic bit, a qubit can also exist in a *superposition* of both a 0 *and* a 1. In other words, it is possible for a single qubit to exist as a 0, a 1, or simultaneously as both a 0 *and* a 1. A qubit that is in two positions at once is said to be in its coherent state.⁷⁹ This can be explained more coherently with an example:

If a coin is flipped in a darkened room, the result of the coin being flipped is mathematically just as likely to be heads or tails. While the light is off, the coin is in a superposition—whereby it is both heads and tails at once, because an [observer] cannot see which it is. If [the observer] turns on the light, [he or she] “collapses” the superposition, and forces the coin to be either heads or tails by measuring it. Measuring something destroys the superposition, forcing it into being in just one classical state.⁸⁰

This coherent state leads to the phenomenon that would make a quantum computer exponentially more powerful than any computer to date; this is the phenomenon called “quantum parallelism.”⁸¹ Essentially, because a qubit in a coherent state holds two values at once, a single operation done on such a qubit would act on both values at the same time.⁸² “Likewise, a two-qubit system would perform the operation on four values and a three-qubit system on eight [values].”⁸³ To summarize, an operation done on a system of n qubits would act on 2^n values simultaneously.⁸⁴ Exhibit 7.25 shows this concept using a system containing three qubits, which represent eight states simultaneously.

“The very property that makes quantum computing so powerful also makes it very fragile and hard to control.”⁸⁵ In order to harness the power of quantum parallelism, scientists need to be able to read and measure the output from the operations performed on groups of qubits. Herein lies the problem of decoherence. When a qubit in the coherent state measurably interacts with the environment, it will immediately decohere and resume one of the two classical states, either a 0 or a 1, and it will no longer exhibit

BEYOND RSA AND DES 7 · 41**EXHIBIT 7.25** Three-Qubit System

Source: Simon Bone and Matias Castro, "A Brief History of Quantum Computing," Imperial College, London, www.doc.ic.ac.uk/~nd/surprise97/journal/vol4/spb3/.

its dual-state ability. In other words, simply looking at a qubit can cause it to decohere, and this makes measuring qubits directly impossible.⁸⁶

If scientists are unable to measure something directly, then they must find a way to measure indirectly, or a practical quantum computer will never be made. One possible answer lies in another property of quantum mechanics called entanglement. Entanglement is an obscure attribute that involves two or more atoms or particles. When certain conditions are met or certain forces are applied to two or more particles, then they can become entangled, whereby the particles exhibit opposite properties. The entangled particles will remain entangled, no matter the physical distance between them, and one entangled particle will always be able to communicate with its partner. Particles spin either up or down, and this spin is how scientists measure information about the particles. The property of coherence tells us that a particle will spin both up and down simultaneously until a scientist looks at it and measures it. "The spin state of the particle being measured is . . . communicated to the correlated particle, which simultaneously assumes the opposite spin direction to that of the measured particle."⁸⁷ Thus, entanglement could allow scientists to know the value of a qubit without actually looking at one. Scientists admit that entanglement is a difficult notion; they are still exploring the concept.⁸⁸ They also acknowledge that it could be years before a workable solution to the problem of measuring information in a quantum system is discovered.⁸⁹

7.6.4.3 Impacts. Although quantum computers, in theory, can perform any task that a classical computer can, this does not necessarily mean that a quantum computer will always outperform a classical computer. Multiplication is an often-cited example of something that would be done just as quickly on a classical computer as on a quantum computer.⁹⁰ From the early stages of quantum computing, scientists knew that to demonstrate the superior computing power, new algorithms would have to be designed to exploit the phenomenon of quantum parallelism. Such algorithms are complex and difficult to devise, but two are driving the development of this highly theorized field: Shor's algorithm and Grover's algorithm.⁹¹

Peter Shor of Bell Labs designed the first quantum algorithm in 1994. Shor's algorithm allows for rapid factoring of very large numbers into their prime factors. For example, scientific estimates state that it would take a modern computer 10^{24} years to

7 · 42 ENCRYPTION

factor a 1,000-digit number; it would take a quantum computer about 20 minutes.⁹² The implications of this quantum algorithm on classic algorithms that depend on the difficulty of factoring for security, such as the widely used RSA algorithm, are immense. “The ability to break the RSA coding system will render almost all current channels of communications unsecure.”⁹³

Lov Grover, also of Bell Labs, invented the second quantum algorithm in 1996. Grover’s algorithm allows a quantum computer to search databases of all kinds much more quickly than any capability existing today. Grover notes that the greatest benefit is gained when his algorithm is used on an unsorted database.⁹⁴ On average, it takes a classical computer $n/2$ number of searches to find a specific entry in a database of n entries. Grover’s algorithm allows the same search to be done in the square root of n number of searches. For example, in a database of 1 million entries, it would take a computer today on average of 500,000 searches to find the right answer; it would take a quantum computer using Grover’s algorithm only 1,000 searches. This could have implications for symmetric key algorithms such as DES, because this algorithm would allow an exhaustive search of all possible keys to occur quite rapidly.⁹⁵

7.6.4.4 Current Status. Encouraged by the repercussions of quantum computing and the related algorithms on the security of information and cryptography, governments around the world are funding efforts to build a practical quantum computing system. The United States has many initiatives on-going. In 2001, the Defense Advanced Research Projects Agency (DARPA) of the Department of Defense launched a \$100 million effort that would last five years. In addition, the National Science Foundation has \$8 million in grant money for researching quantum capabilities. DARPA’s Quantum Information Science and Technology initiative will now exist indefinitely; it became a fully funded and permanent program in 2006.⁹⁶ A number of other governments, primarily within Europe and Asia, are involved in quantum computation research and development. In 2000, the European Commission launched a comprehensive research effort with \$20 million budgeted over three years. In Japan, the Ministry of Post and Telecommunications began an initiative in 2001 that will last 10 years with a total requested budget of \$400 million. There are several commercial enterprises also involved in quantum projects. This includes IBM, Bell Labs, the Japanese firms of Fujitsu, Ltd., NEC Corporation, and Nippon Telephone and Telegraph Corporation.⁹⁷ This list is by no means exhaustive, as there are universities and other organizations throughout the world with research efforts in full swing.

Because of the worldwide effort to understand quantum computing more thoroughly, several key advancements have been made. In 1998, researchers at Los Alamos National Laboratory and MIT were able to spread a qubit over three nuclear spins of certain types of molecules. According to the experiments, spreading the information (qubit) out made it more difficult to corrupt, or decohere. The researchers were able to accomplish this using a technique called nuclear magnetic resonance (NMR), which allows the manipulation and control of a nucleus’s spin. This technique allowed the researchers to use the property of entanglement to analyze indirectly the quantum information.⁹⁸

In 2000, researchers at IBM developed a five-qubit computer, also using the nuclei of a liquid. The nuclei were programmed by radio frequency pulses and then detected by NMR techniques. Using this technique, the team was able to find the period of a particular function, or the length of the shortest interval over which it repeats its values. This problem would take a classical computer several repeated cycles to compute; the team at IBM was able to do it in one step. In 2001, a combined group of scientists from

BEYOND RSA AND DES 7 · 43

IBM and Stanford University demonstrated Shor's algorithm and were able to find the prime factors of 15. The seven-qubit computer correctly deduced that the prime factors were 3 and 5.⁹⁹

In February 2007, a Canadian company called D-Wave claimed to demonstrate the first commercial quantum computer. It is a "supercooled, superconducted niobium chip housing an array of 16 qubits."¹⁰⁰ D-Wave chose not to focus on cryptographic efforts when building the Orion, as the computer is called. Instead, Orion focuses its energy on solving pattern-matching problems and nondeterministic polynomial problems (NP-complete problems). NP-complete problems are decision problems that contain searching and optimization problems, and are used when someone needs to know if a certain solution for a certain problem exists. Examples of such problems include database searches, pattern matching, identifying diseases from symptoms, and finding matches for genetic material.¹⁰¹ The company's demonstrations were done via a television feed from a remote location, due to the sensitive nature of the machine and the difficulty in transporting equipment that is cooled to just above absolute zero. Despite the demonstrations and the claims of D-Wave, scientists are skeptical that Orion is actually performing quantum computations. Even the chief executive of D-Wave said that, although all evidence indicates that Orion is performing quantum computations, there is some uncertainty. Nevertheless, D-Wave announced plans to boost the Orion to 1,000 qubits by 2008.¹⁰²

In July 2007, scientists from NIST (United States) and the Rutherford Appleton Laboratory (United Kingdom) teamed up to explore magnetic quantum effects. This team reports having chained together "100 atoms of yttrium barium nickel oxide into a quantum spin-chain that, in effect, turn[ed] the 30-nanometer long magnetic molecule into a single element."¹⁰³ This discovery is an important step toward putting qubits onto solid-state circuits. Thirty nanometers is well beyond the atomic length scale, and it is unusual to see quantum coherence beyond the atomic level. However, the team did report stable coherent states at this size, which is large enough for the lithographic techniques used to create circuit boards and conductors of classical computers.¹⁰⁴

In April 2013, researchers "successfully transmitted a secure quantum code through the atmosphere from an aircraft to a ground station." The author continues,

"This demonstrates that quantum cryptography can be implemented as an extension to existing systems," says LMU's Sebastian Nauerth. In the experiment, single photons were sent from the aircraft to the receiver on the ground. The challenge was to ensure that the photons could be precisely directed at the telescope on the ground in spite of the impact of mechanical vibrations and air turbulence. "With the aid of rapidly movable mirrors, a targeting precision of less than 3 m over a distance of 20 km was achieved," reports Florian Moll, project leader at the DLR's Institute for Communication and Navigation. With this level of accuracy, William Tell could have hit the apple on his son's head even from a distance of 500 m.

With respect to the rate of signal loss and the effects of air turbulence, the conditions encountered during the experiment were comparable to those expected for transmission via satellite. The same holds for the angular velocity of the aircraft. The success of the experiment therefore represents an important step towards secure satellite-based global communication.¹⁰⁵

Even with the advances just mentioned, skeptics believe that practical quantum computers that outperform classical computers are still years, or even decades, away. After conducting many hours of research on the topic of quantum computing, this author's opinion is that it is not a matter of *if* quantum computing will become a reality but a matter of *when*. That scientists have been able to demonstrate a few theoretical

7 · 44 ENCRYPTION

quantum computations on systems comprised of only a few qubits is highly promising. Yet scientists need to overcome many obstacles. Systems containing hundreds or thousands of qubits will be needed to perform useful computations. In addition, precise controls will be required to accomplish operations while avoiding decoherence; in fact, decoherence is perhaps the biggest obstacle to the creation of a quantum system. Until scientists can reliably measure information produced by qubits at work, it is unlikely that a practical quantum system will be built in the near future.¹⁰⁶

7.6.5 Snake Oil Factor. As encryption vendors and cryptographers come to grips with the implementation and extended testing of new algorithms, it is important to note these words from the AES competition requirements:

A complete written specification of the algorithm shall be included, consisting of all necessary mathematical equations, tables, diagrams, and parameters that are needed to implement the algorithm.

In other words, there is no secret about how the AES will make things secret, just as there is no secret about how DES works. This often strikes the crypto-novice as illogical. Why not keep the algorithm secret? Surely that will make any messages encrypted with it that much harder to decrypt. Not really. Any reliance on the secrecy of the algorithm inserts a weak link in the chain of security. Encrypting data does not guarantee that it will remain confidential. The keys must be kept secret, and the identity of persons requesting authorized access must be verified to ensure they are authentic, and so on. This is true of public key encryption as well as private key encryption.

This principle is known as Kerckhoffs' Principle, based on an 1883 publication by military cryptographer Auguste Kerckhoffs:

1. The system must be practically, if not mathematically, indecipherable;
2. It must not be required to be secret, and it must be able to fall into the hands of the enemy without inconvenience;
3. Its key must be communicable and retainable without the help of written notes, and changeable or modifiable at the will of the correspondents;
4. It must be applicable to telegraphic correspondence;
5. It must be portable, and its usage and function must not require the concourse of several people;
6. Finally, it is necessary, given the circumstances that command its application, that the system be easy to use, requiring neither mental strain nor the knowledge of a long series of rules to observe.¹⁰⁷

There is no benefit to be gained by relying on an algorithm that has not been subject to open review, particularly when strong, reviewed algorithms exist. Beware of encryption vendors, or producers of any security products, that claim strength based on secret algorithms. Such claims are often a case of snake oil. (For more on bogus claims for crypto products, see Curtin's "Snake Oil FAQ," included in Section 7.8, "Further Reading," at the end of this chapter.)

7.7 STEGANOGRAPHY. Instead of scrambling data through cryptography, one can also insert data covertly into other data streams. *Steganography* (literally *covered writing* in Greek) generally uses the low-order bits of a data stream—typically an

FURTHER READING 7 · 45

image—to convey the cleartext. In today’s high-resolution representations of color images, modifying the least significant bits of a pixel makes a negligible change in color, at least to the human eye. The steganographic software can make the changes and then extract them from the modified image.¹⁰⁸

Such modified images are difficult to identify, but steganography detection tools, which rely on detecting abnormal patterns in the pixels of a carrier image, do exist.¹⁰⁹ For example, StegoHunt™ and StegoAnalyst software from Wetstone Technologies can identify and analyze steganographically modified data; StegoBreak can extract the cleartext from the carrier file.¹¹⁰

7.8 FURTHER READING. As stated at the outset, this chapter was not designed to be an extensive treatise on cryptography or a complete guide to the implementation of encryption technology. There are many resources available to help readers deepen their understanding of this fundamental area of information security.

Books and Articles

- Bishop, M. *Computer Security: Art and Science*. Upper Saddle River, NJ: Addison-Wesley/Pearson Education, 2003.
- Hinsley, F. H., and A. Stripp, eds. *Codebreakers: The Inside Story of Bletchley Park*. Oxford, UK: Oxford University Press, 2001.
- Cobb, C. *Cryptography for Dummies*. Hoboken, NJ: John Wiley & Sons, 2003.
- Gilbert, G., Y. S. Weinstein, and M. Hamrick. *Quantum Cryptography*. World Scientific Publications, 2013.
- Goldreich, O. *Foundations of Cryptography: Volume I, Basic Tools*. New York, NY: Cambridge University Press, 2007.
- Juels, Ari. “Encryption Basics.” In H. Bidgoli, ed., *Handbook of Information Security*, Vol. 2. Hoboken, NJ: John Wiley & Sons, 2006.
- Kahn, D. *The Codebreakers: The Comprehensive History of Secret Communication from Ancient Times to the Internet*, Revised Edition. New York: Scribner, 1996.
- Katz, J., and Y. Lindell. *Introduction to Modern Cryptography*, Second Edition. London: Chapman & Hall/CRC, 2014.
- Mao, W. *Modern Cryptography: Theory and Practice*. Upper Saddle River, NJ: Prentice-Hall, 2003.
- Mel, H. X., and D. Baker. *Cryptography Decrypted*. Addison-Wesley, Upper Saddle River, NJ 2000.
- Schneier, B. *Applied Cryptography: Protocols, Algorithms, and Source Code in C*, 2nd ed. New York, NY: John Wiley & Sons, 1996.
- Seberry, J., and J. Pieprzyk. *Cryptography: An Introduction to Computer Security*. Englewood Cliffs, NJ: Prentice-Hall, 1989.
- Spillman, R. J. *Classical and Contemporary Cryptology*. Upper Saddle River, NJ: Prentice-Hall, 2004.
- van Tilborg, H. C. A. & S. Jojodia, eds. *Encyclopedia of Cryptography and Security*, 2nd ed. Springer, 2013.
- Yan, S. Y. *Quantum Attacks on Public-Key Cryptosystems*. Springer, 2013.

Web Resources

- Arsenault, A., and S. Turner. “Internet X.509 Public Key Infrastructure: PKIX Roadmap,” 2000; www.ietf.org/proceedings/00jul/I-D/pkix-roadmap-05.txt

7 · 46 ENCRYPTION

- Bacard, A. “Non-Technical PGP (Pretty Good Privacy) FAQ,” 2002; www.andrebacard.com/pgp.html
- Beezer, R. “Cryptography Independent Study,” 2002; <http://buzzard.ups.edu/courses/2002spring/iscryptos2002.html>
- Cate, V. Vince Cate’s Cryptorebel/Cypherpunk Page, www.offshore.com.ai/security/
Cryptography Research Inc. Research Links: www.cryptography.com/resources/researchlinks.html (URL inactive)
- Curtin, M. “Snake-Oil FAQ/Snake Oil Warning Signs: Encryption Software to Avoid,” 1998; www.interhack.net/people/cmcurtin/snake-oil-faq.html
- Electronic Frontier Foundation. “Frequently Asked Questions (FAQ) About the Electronic Frontier Foundation’s ‘DES Cracker’ Machine,” 1999; http://w2.eff.org/Privacy/Crypto/Crypto_misc/DESCracker/HTML/19980716_eff_des_faq.html
- Electronic Frontier Foundation RSA. “Code-Breaking Contest Again Won by Distributed.Net and Electronic Frontier Foundation (EFF). DES Challenge III Broken in Record 22 Hours,” 1999; http://w2.eff.org/Privacy/Crypto/Crypto_misc/DESCracker/HTML/19990119_deschallenge3.html
- Gerck, E. “Why Is Certification Harder than It Looks?” 1999; <http://mcwg.org/mcg-mirror/whycert.htm>
- ICSA Labs’ Cryptography Community. www.icsalabs.com/icsa/main.php?pid=vjgj7567 (URL inactive)
- International PGP Home Page. 2002; www.pgpi.org
- Kessler, G. “An Overview of Cryptography,” 2004; www.garykessler.net/library/crypto.html
- PGP Home. www.pgp.com/index.php (URL inactive)
- RSA Security Content Library. www.rsasecurity.com/doc_library/index.asp
- Schneier, B. *Crypto-Gram* newsletter archive, 1998–2008; www.schneier.com/crypto-gram-back.html

7.9 NOTES

1. David Kahn, *The Codebreakers: The Story of Secret Writing* (New York: Scribner, 1996), pp. 980–984.
2. *The American Heritage® New Dictionary of Cultural Literacy*, 3rd ed. s.v. “cryptography,” <http://dictionary.reference.com/browse/cryptography>
3. RSA Laboratories, “What Is Cryptanalysis?” www.rsa.com/rsalabs/node.asp?id=2200 (URL inactive)
4. J. Seberry and J. Pieprzyk, *Cryptography: An Introduction to Computer Security* (Englewood Cliffs, NJ: Prentice-Hall, 1989).
5. Seberry and Pieprzyk, *Cryptography*
6. Kahn, *Codebreakers*, pp. 71–72.
7. *Encyclopaedia Britannica Online Academic Edition*, s.v. “cryptology,” <http://search.eb.com.library.norwich.edu/eb/article-25638>. This article only available within the Norwich University system, or with a paid subscription to the Encyclopedia.
8. Brigitte Collard, “La cryptographie dans l’Antiquité gréco-romaine. III. Le chiffrement par transposition,” *Folia Electronica Classica* (Louvain-la-Neuve) 7 (January-June 2004): section II(2), “Définition de la scytale”; <http://bcs.fltr.ucl.ac.be/FE/07/CRYPT/Crypto44-63.html#42047>

NOTES 7 · 47

9. Brad Stark, "A Closer Look at Cryptography," Bucknell University, www.facstaff.bucknell.edu/udaapp/090/w3/brads.htm (URL inactive)
10. Kahn, *Codebreakers*, pp. 83–84.
11. "Time Table/Time-Travel through Cryptography and Cryptanalysis," www.cryptool.com/menu_zeittafel.en.html
12. Oliver Pell, "Cryptology," www.ridex.co.uk/cryptology/#_Toc439908853
13. "Time Table/Time-Travel."
14. *Encyclopaedia Britannica Online Academic Edition*, s.v. "cryptology."
15. *Encyclopaedia Britannica Online Academic Edition*, s.v. "cryptology."
16. Kahn, *Codebreakers*, p. 107.
17. *Encyclopedia Britannica Online Academic Edition*, s.v. "cryptology."
18. National Security Agency, "The Rare Books Collection: Giovanni Battista Porta," www.nsa.gov/about/cryptologic_heritage/center_crypt_history/publications/rare_books.shtml#giovanni
19. D. Salomon, *Coding for Data and Computer Communications* (New York: Springer, 2005), p. 218; <http://tinyurl.com/2dsmc8>
20. *Encyclopedia Britannica Online Academic Edition*, s.v. "cryptology."
21. Stallings, W. *Network and Internetwork Security Principles and Practices*. Prentice-Hall, January, 1995
22. Kevin Romano, "The Stager Ciphers and the US Military's First Cryptographic System," www.gordon.army.mil/AC/Wntr02/stager.htm
23. Cypher Research Laboratories, "A Brief History of Cryptography," www.cypher.com.au/crypto_history.htm
24. National Archives, "Teaching with Documents: The Zimmermann Telegram," www.archives.gov/education/lessons/zimmermann
25. National Archives, "Teaching with Documents."
26. Jacob Mathai, "History of Cryptography and Secrecy Systems," Fordham University, www.dsm.fordham.edu/~mathai/crypto.html#ENIGMA
27. Judson Knight, "Cryptology, History," www.espionageinfo.com/Cou-De-Cryptology-History.html
28. Knight, "Cryptology, History."
29. Oli Cooper, "Cryptography," University of Bristol, www.cs.bris.ac.uk/cooper/Cryptography/crypto.html
30. Cooper, "Cryptography."
31. Kahn, *Codebreakers*,
32. Kahn, *Codebreakers*,
33. Jacob Mathai, "History of Cryptography and Secrecy Systems," Fordham University, www.dsm.fordham.edu/~mathai/crypto.html#OneTimePad
34. Stallings, "Network and Internetwork Security"
35. SCI.CRYPT FAQ §5.2, www.faqs.org/faqs/cryptography-faq/part5
36. Kahn, *Codebreakers*, p. 979.
37. Ari Juels, "Encryption Basics," in *Handbook of Information Security*, Vol. 2, ed. H. Bidgoli, (Hoboken, NJ: John Wiley & Sons, 2006), p. 980.
38. Juels, "Encryption Basics," p. 981.
39. Juels, "Encryption Basics," p. 471.

7 · 48 ENCRYPTION

40. SCI.CRYPT FAQ, www.faqs.org/faqs/cryptography-faq/part6/
41. Bruce Schneier, *Applied Cryptography*, 2nd ed. (New York: John Wiley & Sons, 1996), p. 461.
42. Kahn, *Codebreakers*, p. 982.
43. Juels, “Encryption Basics,” p. 474.
44. Juels, “Encryption Basics,” p. 474.
45. Juels, “Encryption Basics,” p. 474.
46. RSA Laboratories, “What Is Diffie-Hellman?” www.rsa.com/rsalabs/node.asp?id=2248
47. Charlie Kaufman, “IPsec: IKE (Internet Key Exchange),” vol. 1, *Handbook for Information Security* (Hoboken, NJ: John Wiley & Sons, 2006), p. 974.
48. Juels, “Encryption Basics,” pp. 474–475.
49. SCI.CRYPT FAQ, www.faqs.org/faqs/cryptography-faq/part6/
50. SCI.CRYPT FAQ, www.faqs.org/faqs/cryptography-faq/part6/
51. SCI.CRYPT FAQ, www.faqs.org/faqs/cryptography-faq/part6/
52. SCI.CRYPT FAQ, www.faqs.org/faqs/cryptography-faq/part6/
53. SCI.CRYPT FAQ, www.faqs.org/faqs/cryptography-faq/part6/
54. For more on this evolution, see the excellent IETF document, “Internet X.509 Public Key Infrastructure: PKIX Roadmap,” by A. Arsenault and S. Turner. www.ietf.org/proceedings/00jul/I-D/pkix-roadmap-05.txt
55. A. Arsenault and S. Turner, “Internet X.509 Public Key Infrastructure: PKIX Roadmap.”
56. Michael Baum, *Secure Electronic Commerce* (Prentice-Hall, 1997)
57. E. Gerck, “Why Is Certification Harder than It Looks?” 1999, <http://mcwg.org/mcg-mirror/whycert.htm>
58. National Security Agency, “The Case for Elliptic Curve Cryptography,” www.nsa.gov/business/programs/elliptic_curve.shtml
59. National Security Agency, “The Case for Elliptic Curve Cryptography.”
60. Certicom, “An Elliptic Curve Cryptography (ECC) Primer,” www.certicom.com/pdfs/WP-ECCprimer_login.pdf
61. National Security Agency, “The Case for Elliptic Curve Cryptography.”
62. Jacob West, “The Quantum Computer,” www.cs.rice.edu/~taha/teaching/05F/210/news/2005_09_16.htm
63. West, “Quantum Computer.”
64. Simon Bone and Matias Castro, “A Brief History of Quantum Computing,” Imperial College, London, www.doc.ic.ac.uk/~nd/surprise_97/journal/vol4/spb3
65. Quantum Information Partners LLP, “Short History of Quantum Information Processing,” www.qipartners.com/publications/Short_History_of_QC.pdf. (URL inactive)
66. West, “Quantum Computer.”
67. West, “Quantum Computer.”
68. Bone and Castro, “Brief History.”
69. Bone and Castro, “Brief History.”
70. Simon Bone, “The Hitchhiker’s Guide to Quantum Computing,” www.doc.ic.ac.uk/~nd/surprise_97/journal/vol1/spb3/

NOTES 7 · 49

71. Bone, "Hitchhiker's Guide."
72. West, "Quantum Computer."
73. *American Heritage® Dictionary of the English Language, Fourth Edition*, s.v. "quanta," <http://dictionary.reference.com/browse/quantum>
74. Genomics and Proteomics, "Glossary: Quantum Mechanics," www.genpromag.com/Glossary.aspx?LETTER=Q (URL inactive)
75. Stephen Jenkins, "Some Basic Ideas About Quantum Mechanics," University of Exeter, newton.ex.ac.uk/research/qsystems/people/jenkins/mbody/mbody2.html
76. West, "Quantum Computer."
77. Bone and Castro, "Brief History."
78. Bonsor and Strickland, "How Quantum Computers Work," computer. www.howstuffworks.com/quantum-computer3.htm
79. Bone and Castro, "Brief History."
80. Duncan McKimm, "Quantum Entanglement," www.abc.net.au/science/features/quantum
81. West, "Quantum Computer."
82. Bone and Castro, "Brief History."
83. Bone and Castro, "Brief History."
84. Bone and Castro, "Brief History."
85. Bone and Castro, "Brief History."
86. Bonsor and Strickland, "How Quantum Computers Work."
87. SearchSMB.com, "Entanglement," <http://searchcio-midmarket.techtarget.com/definition/entanglement>
88. McKimm, "Quantum Entanglement."
89. SearchSMB.com, "Entanglement."
90. Bone and Castro, "Brief History."
91. Bone and Castro, "Brief History."
92. Bone and Castro, "Brief History."
93. Bone, "Hitchhiker's Guide."
94. Lov Grover, "What's a Quantum Phone Book?" www.bell-labs.com/user/feature/archives/lkgrover
95. Bone and Castro, "Brief History."
96. Confidential Source #3, "Quantum Computing," Internal Research branch Web site, August 10, 2007.
97. Confidential Source #3, "Quantum Computing."
98. West, "Quantum Computer."
99. Bonsor and Strickland, "How Quantum Computers Work."
100. R. Colin Johnson "Quantum Computer 'Orion' Debuts," EETimes, www.eetimes.com/electronics-news/4069654/Quantum-computer-Orion-debuts
101. Johnson, "Quantum Computer 'Orion' Debuts."
102. Jordon Robertson, "Scientists Dubious of Quantum Computer Claims," abcnews.go.com/Technology/wireStory?id=2875656.
103. R. Colin Johnson, "Circuit-Sized Quantum Effect Observed," EETimes, www.eetimes.com/electronics-news/4073585/Circuit-sized-quantum-effect-observed

7 · 50 ENCRYPTION

104. Johnson, "Circuit-Sized Quantum Effect Observed."
105. ScienceDaily. "Quantum Cryptography: On Wings of Light." *ScienceDaily Science News*. 04 13, 2013. www.sciencedaily.com/releases/2013/04/130403071950.htm (accessed May 16, 2013).
106. Johnson, "Circuit-Sized Quantum Effect Observed."
107. Petitcolas, Fabien. *la cryptographie militaire*. 2012. www.petitcolas.net/fabien/kerckhoffs/#english (accessed May 21, 2013).
108. Kessler, Gary C. *Steganography: Hiding Data Within Data*. 2002. www.garykessler.net/library/steganography.html (accessed May 21, 2013).
109. National Institute of Justice. *Digital Evidence Analysis: Steganography Detection*. 11 05, 2010. www.nij.gov/topics/forensics/evidence/digital/analysis/steganography.htm (accessed 05 21, 2013).
110. Wetstone Technologies. "StegoHunt." *Wetstone*. 2013. www.wetstonetech.com/product/stegohunt/ (accessed May 21, 2013).

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 8

USING A COMMON LANGUAGE FOR COMPUTER SECURITY INCIDENT INFORMATION

John D. Howard

8.1 INTRODUCTION	8·1	8.4.3 Full Incident Information Taxonomy	8·15
8.2 WHY A COMMON LANGUAGE IS NEEDED	8·2	8.5 ADDITIONAL INCIDENT INFORMATION TERMS	8·16
8.3 DEVELOPMENT OF THE COMMON LANGUAGE	8·3	8.5.1 Success and Failure 8.5.2 Site and Site Name 8.5.3 Other Incident Terms	8·17 8·17 8·17
8.4 COMPUTER SECURITY INCIDENT INFORMATION TAXONOMY	8·4	8.6 HOW TO USE THE COMMON LANGUAGE	8·18
8.4.1 Events	8·4		
8.4.2 Attacks	8·11	8.7 NOTES	8·20

8.1 INTRODUCTION. A computer security *incident* is some set of events that involves an attack or series of attacks at one or more sites. (See Section 8.4.3 for a more formal definition of the term “incident.”) Dealing with these incidents is inevitable for individuals and organizations at all levels of computer security. A major part of dealing with these incidents is recording and receiving incident information, which almost always is in the form of relatively unstructured text files. Over time, these files can end up containing a large quantity of very valuable information. Unfortunately, the unstructured form of the information often makes incident information difficult to manage and use.

This chapter presents the results of several efforts over the last few years to develop and propose a method to handle these unstructured computer security incident records. Specifically, this chapter presents a *tool* designed to help individuals and organizations record, understand, and share computer security incident information. We call the tool the *common language for computer security incident information*. This common language contains two parts:

1. A set of “high-level” incident-related terms
2. A method of classifying incident information (a taxonomy)

8 · 2 USING A COMMON LANGUAGE FOR COMPUTER SECURITY INCIDENT

The two parts of the common language, the terms and the taxonomy, are closely related. The taxonomy provides a structure that shows how most common-language terms are related. The common language is intended to help investigators improve their ability to:

- Talk more understandably with others about incidents
- Gather, organize, and record incident information
- Extract data from incident information
- Summarize, share, and compare incident information
- Use incident information to evaluate and decide on proper courses of action
- Use incident information to determine effects of actions over time

This chapter begins with a brief overview of why a common language is needed, followed by a summary of how the incident common language was developed. We then present the common language in two parts: (1) incident terms and taxonomy and (2) additional incident information terms. The final section contains information about some practical ways to use the common language.

8.2 WHY A COMMON LANGUAGE IS NEEDED. When the first edition of this *Handbook* was published more than 30 years ago, computer security was a small, obscure, academic specialty. Because there were only a few people working in the field, the handling of computer security information could largely take place in an ad hoc way. In this environment, individuals and groups developed their own terms to describe computer security information. They also developed, gathered, organized, evaluated, and exchanged their computer security information in largely unique and unstructured ways. This lack of generalization has meant that computer security information has typically not been easy to compare or combine, or sometimes even to talk about in an understandable way.

Progress over the years in agreeing on a relatively standard set of terms for computer security (a common language) has had mixed results. One problem is that many terms are not yet in widespread use. Another problem is that the terms that are in widespread use often do not have standard meanings. An example of the latter is the term “computer virus.” We hear the term frequently, not only in academic forums but also in the news media and popular publications. It turns out, however, that even in academic publications, “computer virus” has no accepted definition.¹ Many authors define a computer virus to be “a code fragment that copies itself into a larger program.”² They use the term “worm” to describe an independent program that performs a similar invasive function (e.g., the Internet Worm in 1988). But other authors use the term “computer virus” to describe *both* invasive code fragments and independent programs.

Progress in developing methods to gather, organize, evaluate, and exchange computer security information also has had limited success. For example, the original records (1988–1992) of the Computer Emergency Response Team (now the CERT Coordination Center or CERT/CC) are simply a file of email and other files sent to the CERT/CC. These messages and files were archived together in chronological order, without any other organization. After 1992, the CERT/CC and other organizations developed methods to organize and disseminate their information, but the information remains difficult to combine or compare because most of it remains almost completely textual information that is uniquely structured for the CERT/CC.

DEVELOPMENT OF THE COMMON LANGUAGE 8 · 3

Such ad hoc terms and ad hoc ways to gather, organize, evaluate, and exchange computer security information are no longer adequate. Far too many people and organizations are involved, and there is far too much information to understand and share. Today computer security is an increasingly important, relevant, and sophisticated field of study. Numerous individuals and organizations now regularly gather and disseminate computer security information. Such information ranges all the way from the security characteristics and vulnerabilities of computers and networks, to the behavior of people and systems during security incidents—far too much information for each individual and organization to have its own unique language.

One of the key elements to making systematic progress in any field of inquiry is the development of a consistent set of terms and taxonomies (principles of classification) that are used in that field.³ This is a necessary and natural process that leads to a growing *common language*, which enables gathering, exchanging, and comparing information. In other words, the more a field of inquiry such as computer security grows, the more a common language is needed to understand and communicate with one another.

8.3 DEVELOPMENT OF THE COMMON LANGUAGE. Two of the more significant efforts in the process of developing this common language for computer security incident information were (1) a project to classify more than 4,300 Internet security incidents completed in 1997,⁴ and (2) a series of workshops in 1997 and 1998 called the *Common Language Project*. Workshop participants included people primarily from the Security and Networking Research Group at the Sandia National Laboratories, Livermore, California, and from the CERT/CC at the Software Engineering Institute, Carnegie Mellon University, Pittsburgh, Pennsylvania. Additional participation and review came from people in the Department of Defense (DoD) and the National Institute of Standards and Technology (NIST).

These efforts to develop the common language were *not* efforts to develop a comprehensive dictionary of terms. Instead, the participants were trying to develop both a *minimum* set of “high-level” terms to describe computer security attacks and incidents, and a structure and classification scheme for these terms (a *taxonomy*), which could be used to classify, understand, exchange, and compare computer security attack and incident information.

Participants in the workshops hoped this common language would gain wide acceptance because of its usefulness. There is already evidence that this acceptance is taking place, particularly at incident response teams and in the DoD.

In order to be complete, logical, and useful, the common language for computer security incident information was based initially and primarily on theory (i.e., it was *a priori* or nonempirically based).⁵ Classification of actual Internet security incident information was then used to refine and expand the language. More specifically, the common language development proceeded in six stages:

1. Records at the CERT/CC for incidents reported to them from 1988 through 1995 were examined to establish a preliminary list of terms used to describe computer security incidents.
2. The terms in this list, and their definitions, were put together into a structure (a preliminary taxonomy).
3. This preliminary taxonomy was used to classify the information in the 1988 through 1995 incident records.
4. The preliminary taxonomy and classification results were published in 1997.⁶

8 · 4 USING A COMMON LANGUAGE FOR COMPUTER SECURITY INCIDENT

5. A series of workshops was conducted from 1997 through 1998 (the *Common Language Project*) to make improvements to the taxonomy and to add additional terms.
6. The results of the workshops (the “common language for security incidents”) were first published in 1998.

A *taxonomy* is a classification scheme (a structure) that partitions a body of knowledge and defines the relationship of the pieces.⁷ Most of the terms in this common language for security incident information are arranged in such a taxonomy, as presented in the next section. *Classification* is the process of using a taxonomy for separating and ordering. As discussed earlier, classification of information using a taxonomy is necessary for computer security incident information because of the rapidly expanding amount of information and the nature of that information (primarily text). Classification using the common-language taxonomy is discussed in the final section of this chapter.

Our experience has shown that satisfactory taxonomies have classification categories with these six characteristics⁸:

1. **Mutually exclusive.** Classifying in one category excludes all others because categories do not overlap.
2. **Exhaustive.** Taken together, the categories include all possibilities.
3. **Unambiguous.** The taxonomy is clear and precise, so that classification is not uncertain, regardless of who is doing the classifying.
4. **Repeatable.** Repeated applications result in the same classification, regardless of who is doing the classifying.
5. **Accepted.** It is logical and intuitive, so that categories can become generally approved.
6. **Useful.** The taxonomy can be used to gain insight into the field of inquiry.

These characteristics were used to develop and evaluate the common-language taxonomy. A taxonomy, however, is merely an approximation of reality, and as such, even the best taxonomy will fall short in some characteristics. This may be especially true when the characteristics of the data being classified are imprecise and uncertain, as is typical for computer security incident information. Nevertheless, classification is an important, useful, and necessary prerequisite for systematic study of incidents.

8.4 COMPUTER SECURITY INCIDENT INFORMATION TAXONOMY. We have been able to structure most of the terms in the common language for security incident information into a taxonomy. These terms and the taxonomy are presented in this section. Additional terms that describe the more general aspects of incidents are presented in Section 8.5.

8.4.1 Events. The operation of computers and networks involves innumerable *events*. In a general sense, an event is a discrete change of state or status of a system or device.⁹ From a computer security viewpoint, these changes of state result from *actions* that are directed against specific *targets*. An example is a user taking action to log in to the user’s account on a computer system. In this case, the action taken by the user is to *authenticate* to the login program by claiming to have a specific identity and

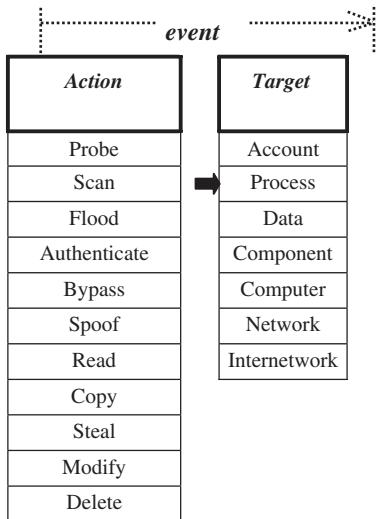
COMPUTER SECURITY INCIDENT INFORMATION TAXONOMY 8 · 5

EXHIBIT 8.1 Computer and Network Events

then presenting the required verification. The target of this action would be the user's *account*. Other examples include numerous actions that can be targeted toward:

- *Data* (e.g., actions to *read*, *copy*, *modify*, *steal*, or *delete*)
- A *process* (e.g., actions to *probe*, *scan*, *authenticate*, *bypass*, or *flood* a running computer process or execution thread)
- A *component*, *computer*, *network*, or *internetwork* (e.g., actions to *scan* or *steal*)

Exhibit 8.1 presents a matrix of actions and targets that represent possible computer and network events (although not all of the possible combinations shown are feasible). A computer or network event is defined as:

Event—action directed at a target that is intended to result in a change of state, or status, of the target.¹⁰

Several aspects of this definition are important to emphasize. First, in order for there to be an event, there must be an action that is taken, and it must be directed against a target, but the action does *not* have to succeed in actually changing the state of the target. For example, if a user enters an incorrect user name and password combination when logging in to an account, an authentication event has taken place, but the event was not successful in verifying that the user has the proper credentials to access that account.

A second important aspect is that an event represents a *practical* linkage between an action and a specific target against which the action is directed. As such, it represents the way people generally conceptualize events on computers and networks, and not all of the individual steps that actually take place during an event. For example, when a user logs in to an account, we classify the action as *authenticate* and the target as *account*. The actual action that takes place is for the user to access a process (e.g., a "login" program) in order to authenticate. Trying to depict all of the individual steps

8 · 6 USING A COMMON LANGUAGE FOR COMPUTER SECURITY INCIDENT

is an unnecessary complication; the higher-level concepts presented here can describe correctly and accurately the event in a form well understood by people. In other words, it makes sense to abstract the language and its structure to the level at which people generally conceptualize the events.

By all means, supporting evidence should be presented so the evidence provides a complete idea of what happened. Stated another way, abstraction, conceptualization, and communication should be applied as close to the evidence as possible. For example, if a network switch is the target of an attack, then the target should normally be viewed as a computer or as a component (depending on the nature of the switch), and not the network, because assuming the network is the target may be an inaccurate interpretation of the evidence.

Another aspect of the definition of event is that it does not make a distinction between *authorized* and *unauthorized* actions. Most events that take place on computers or networks are both routine and authorized and, therefore, are not of concern to security professionals. Sometimes, however, an event is part of an attack or is a security concern for some other reason. This definition of event is meant to capture both authorized and unauthorized actions. For example, if a user authenticates properly, by giving the correct user identification and password combination while logging in to an account, that user is given access to that account. It may be the case, however, that this user is masquerading as the actual user, after having obtained the user identification and password from snooping on the network. Either way, this is still considered authentication.

Finally, an important aspect of events is that not all of the possible events (the action–target combinations depicted in Exhibit 8.1) are considered likely or even possible. For example, an action to *authenticate* is generally associated with an *account* or a *process* and not a different target, such as *data* or a *component*. Other examples include *read* and *copy*, which are generally targeted toward *data*; *flooding*, which is generally targeted at an *account*, *process*, or *system*; or *stealing*, which is generally targeted against *data*, a *component*, or a *computer*.

We define *action* and *target* as follows:

Action—step taken by a user or process in order to achieve a result,¹¹ such as to probe, scan, flood, authenticate, bypass, spoof, read, copy, steal, modify, or delete.

Target—computer or network logical entity (account, process, or data) or a physical entity (component, computer, network or internetwork).

8.4.1.1 Actions. The actions depicted in Exhibit 8.1 represent a spectrum of activities that can take place on computers and networks. An action is a step taken by a *user* or a *process* in order to achieve a result. Actions are initiated by accessing a target, where *access* is defined as:

Access—establish logical or physical communication or contact.¹²

Two actions are used to gather information about targets: *probe* and *scan*. A probe is an action to determine one or more characteristics of a specific target. This is unlike a scan, which is an action where a user or process accesses a set of targets systematically, in order to determine which targets have one or more characteristics.

“Probe” and “scan” are terms commonly used by incident response teams. As a result, they have common, accepted definitions. Despite this, there is a logical ambiguity: A scan could be viewed as multiple probes. In other words, if an attacker is testing

COMPUTER SECURITY INCIDENT INFORMATION TAXONOMY 8 · 7

Test for One or More Characteristics	
Test a Single Host	Probe
Nonsystematically Test Multiple Hosts	Multiple probes
Systematically Test a Set of Hosts	Scan

EXHIBIT 8.2 Probe Compared to Scan

for one or more characteristics on multiple hosts, this can be (a) multiple attacks (all *probes*), or (b) one attack (a *scan*). This point was discussed extensively in the Common Language Project workshops, and the conclusion was that the terms in the common language should match, as much as possible, their common usage. This common usage is illustrated in Exhibit 8.2.

With probes and scans, it is usually obvious what is taking place. The attacker is either “hammering away” at one host (a *probe*), randomly testing many hosts (multiple *probes*), or using some “automatic” software to look for the same characteristic(s) systematically across a group of hosts (a *scan*). As a practical matter, incident response teams do not usually have a problem deciding what type of action they are dealing with.

One additional point about *scan* is that the term “systematic” is not meant to specify some specific pattern. The most sophisticated attackers try to disguise the systematic nature of a scan. A scan may, at first, appear to be multiple probes. For example, an attacker may randomize a scan with respect to hosts and with respect to the characteristic(s) being tested. If the attack can be determined to involve testing of one or more characteristics on a group of hosts with some common property (e.g., an Internet Protocol [IP] address range) or if tests on multiple hosts appear to be otherwise related (e.g., having a common origin in location and time), then the multiple probes should be classified as a scan.

Unlike probe or scan, an action taken to *flood* a target is not used to gather information about the target. Instead, the desired result of a flood is to overwhelm or overload the target’s capacity by accessing the target repeatedly. An example is repeated requests to open connections to a port over a network or repeated requests to initiate processes on a computer. Another example is a high volume of email messages, which may exceed the resources available for the targeted account.

Authenticate is an action taken by a user to assume an identity. Authentication starts with a user accessing an authentication process, such as a login program. The user must claim to have a certain identity, such as by entering a user name. Usually verification is also required as a second authentication step. For verification, the user must prove knowledge of some secret (e.g., a password), prove the possession of some token (e.g., a secure identification card), and/or prove to have a certain characteristic (e.g., a retinal scan pattern). Authentication can be used not only to log in to an account but also to access other objects, such as to operate a process or to access a file. In other words, the target of an authentication action is the entity (e.g., account, process, or data) that the user is trying to access, not the authentication process itself.

Two general methods might be used to defeat an authentication process. First, a user could obtain a valid identification and verification pair that could be used to authenticate, even though it does not belong to that user. For example, during an

8 · 8 USING A COMMON LANGUAGE FOR COMPUTER SECURITY INCIDENT

incident, an attacker might use a process operating on an Internet host computer that captures user name, password, and IP address combinations that are sent in clear text across the Internet. The attacker could then use this captured information to authenticate (log in) to accounts that belong to other users. It is important to note, as mentioned earlier, that this action is still considered *authenticate*, because the attacker presents valid identification and verification pairs, even though they have been stolen.

The second method that might be used to defeat an authentication process is to exploit a vulnerability in order to bypass the authentication process and access the target. *Bypass* is an action taken to avoid a process by using an alternative method to access a target. For example, some operating systems have vulnerabilities that an attacker could exploit to gain privileges without actually logging in to a privileged account.

As was discussed with respect to *authenticate*, an action to *bypass* does not necessarily indicate that the action is unauthorized. For example, some programmers find it useful to have a shortcut (“back-door”) method to enter an account or run a process, particularly during development. In such a situation, an action to bypass may be considered authorized.

Authenticate and *bypass* are actions associated with users identifying themselves. In network communications, processes also identify themselves to each other. For example, each packet of information traveling on a network contains addresses identifying both the source and the destination, as well as other information. “Correct” information in these communications is assumed, since it is automatically generated. Thus, no action is included on the list to describe this normal situation. Incorrect information could, however, be entered into these communications. Supplying such false information is commonly called an action to *spoof*. Examples include IP spoofing, mail spoofing, and Domain Name System (DNS) spoofing.

Spoofing is an active security attack in which one machine on the network masquerades as a different machine. . . . [It] disrupts the normal flow of data and may involve injecting data into the communications link between other machines. This masquerade aims to fool other machines on the network into accepting the imposter as an original, either to lure the other machines into sending it data or to allow it to alter data.¹³

Some actions are closely associated with data found on computers or networks, particularly with files: *read*, *copy*, *modify*, *steal*, and *delete*. There has been some confusion over these terms because their common usage in describing the physical world sometimes differs from their common usage describing the electronic world. For example, if I say that an attacker *stole* a computer, then you can assume I mean the attacker took possession of the target (computer) and did not leave an identical computer in that location. If I say, however, that the attacker stole a computer *file*, what does that actually mean? It is often taken to mean that the attacker *duplicated* the file and now has a copy, but also it means that the original file is still in its original location. In other words, “*steal*” sometimes means something different in the physical world than it does in the electronic world.

It is confusing for there to be differences in the meaning of actions in the physical world and the electronic world. Workshop participants attempted to reconcile these differences by carefully defining each term (*read*, *copy*, *modify*, *steal*, or *delete*) so it would have a very specific and mutually exclusive meaning that matches the physical-world meaning as much as possible.

Read is defined as an action to obtain the content of the data contained within a file or other data medium. This action is distinguished conceptually from the actual physical steps that may be required to read. For example, in the process of reading a computer

COMPUTER SECURITY INCIDENT INFORMATION TAXONOMY 8 · 9

file, the file may be copied from a storage location into the computer's main memory and then displayed on a monitor to be read by a user. These physical steps (copy the file into memory and then onto the monitor) are not part of the abstract concept of *read*. In other words, to read a target (obtain the content in it), copying of the file is not necessarily required, and it is conceptually not included in our definition of *read*.

The same separation of concepts is included in the definition of the term "copy." In this case, we are referring to acquiring a copy of a target without deleting the original. The term "copy" does not imply that the *content* in the target is obtained, just that a *copy* has been made and obtained. To get the content, the file must be *read*. An example is copying a file from a hard disk to a floppy disk. This copying is done by duplicating the original file while leaving the original file intact. A user would have to open the file and look at the content in order to *read* it.

Copy and *read* are both different concepts from *steal*, which is an action that results in the attacker taking possession of the target and the target also becoming unavailable to the original owner or user. This definition agrees with our concepts about physical property, specifically that there is only one object that cannot be copied. For example, if someone steals a car, then that person has deprived the owner of his or her possession. When dealing with property that is in electronic form, such as a computer file, often the term "steal" is used, when *copy* is what actually is meant. The term "steal" specifically means that the original owner or user has been denied access or use of the target. On the other hand, stealing also could mean *physically* taking a floppy disk that has the file located on it or stealing an entire computer.

Two other actions involve changing the target in some way. The first are actions to *modify* a target. Examples include changing the content of a file, changing the password of an account, sending commands to change the characteristics of an operating process, or adding components to an existing system. If the target is eliminated entirely, the term "delete" is used to describe the action.

As stated earlier, differences in usage of terms between the physical world and the electronic world are undesirable. As such, we tried to be specific and consistent in our usage. The resulting set of terms is exhaustive and mutually exclusive, but goes against the grain in some common usage for the electronic world, particularly with respect to the term "steal." The situation seems unavoidable. Here are some examples that might clarify the terms:

- A user clicks on a link with the browser and sees the content of a Web page on the computer screen. We would classify this as a *read*. While what actually happens is that the content of the page is stored in volatile memory, copied to the cache on the hard drive, and displayed on the screen, from a *logical* (i.e., user) point of view, the Web page has *not* been copied (nor stolen). Now, if a user copies the content of the Web page to a file or prints it out, then the user *has* copied the Web page. Again, this would be a logical classification of the action, from the user's point of view.
- A user duplicates a file that is encrypted. We would classify this as *copy*, not *read*. In this case, the file was reproduced, but the content not obtained, so it was not *read*.
- A user deletes several entries in a password or group file. Should this action be described as several *delete* actions or as one action to *modify*? We would describe this action as *modify*, and the target is *data*. There is no ambiguity here because of the definition of "data." Data are defined to be either a stationary file or a file in transit (see the next section). If a user deletes a line out of the password file, then the file has been modified. The action would be described as *delete* only if

8 · 10 USING A COMMON LANGUAGE FOR COMPUTER SECURITY INCIDENT

the whole file was deleted. If we had defined data to include part of a file, then we would indeed have an ambiguity.

- A user copies a file and deletes the original. We would classify this as *steal*. Although the steps actually include a *copy* followed by a *delete*, that is the electronic way of stealing a file, and therefore it is more descriptive to describe the action as *steal*.

In reality, the term “steal” is rarely used (correctly) because attackers who copy files usually do not delete the originals. The term “steal” often is used *incorrectly*, as in “stealing the source code,” when in fact the correct term is *copy*.

The list of actions was hashed over in numerous group discussions, off and on, for several years before being put into the common language. Most people who participated in these discussions were not entirely happy with the list, but it is the best we have seen so far. Specifically, the list seems to capture all of the common terms with their common usage (*probe*, *scan*, *flood*, *spoof*, *copy*, *modify*, and *delete*) and the other terms are logical (to the people who participated in the discussion groups) and are necessary to make the action category exhaustive (*authenticate*, *bypass*, *read*, and *steal*).

Here is a summary of our definitions of the actions shown in Exhibit 8.1.

Probe—access a target in order to determine one or more of its characteristics.

Scan—access a set of targets systematically in order to identify which targets have one or more specific characteristics.¹⁴

Flood—access a target repeatedly in order to overload the target’s capacity.

Authenticate—present an identity to a process and, if required, verify that identity, in order to access a target.¹⁵

Bypass—avoid a process by using an alternative method to access a target.¹⁶

Spoof—masquerade by assuming the appearance of a different entity in network communications.¹⁷

Read—obtain the content of data in a storage device or other data medium.¹⁸

Copy—reproduce a target leaving the original target unchanged.¹⁹

Steal—take possession of a target without leaving a copy in the original location.

Modify—change the content or characteristics of a target.²⁰

Delete—remove a target or render it irretrievable.²¹

8.4.1.2 Targets. Actions are considered to be directed toward seven categories of targets. The first three of these are “logical” entities (*account*, *process*, and *data*), and the other four are “physical” entities (*component*, *computer*, *network*, and *internetwork*).

In a multiuser environment, an *account* is the domain of an individual user. This domain includes the files and processes the user is authorized to access and use. A special program that records the user’s account name, password, and use restrictions controls access to the user’s account. Some accounts have increased or special permissions that allow access to system accounts, other user accounts, or system files and processes, and often are called *privileged*, *superuser*, *administrator*, or *root* accounts.

Sometimes an action may be directed toward a *process*, which is a program executing on a computer or network. In addition to the program itself, the process includes the program’s data and stack; its program counter, stack pointer, and other registers; and all other information needed to execute the program.²² The action may then be to supply information to the process or command the process in some manner.

COMPUTER SECURITY INCIDENT INFORMATION TAXONOMY 8 · 11

The target of an action may be *data* that are found on a computer or network. Data are representations of facts, concepts, or instructions in forms that are suitable for use by either users or processes. Data may be found in two forms: files or data in transit. *Files* are data that are designated by name and considered as a unit by the user or by a process. Commonly we think of files as being located on a storage medium, such as a storage disk, but files also may be located in the volatile or nonvolatile memory of a computer. *Data in transit* are data being transmitted across a network or otherwise emanating from some source. Examples of the latter include data transmitted between devices in a computer and data found in the electromagnetic fields that surround computer monitors, storage devices, processors, network transmission media, and the like.

Sometimes we conceptualize the target of an action as not being a *logical* entity (account, process, or data) but rather as a *physical* entity. The smallest of the physical entities is a *component*, which is one of the parts that make up a computer or network. A *network* is an interconnected or interrelated group of computers, along with the appropriate switching elements and interconnecting branches.²³ When a computer is attached to a network, it is sometimes referred to as a *host computer*. If networks are connected to each other, then they are sometimes referred to as an *internetwork*.

Here is a summary of our definitions of the targets shown in Exhibit 8.1.

Account—domain of user access on a computer or network that is controlled according to a record of information which contains the user's account name, password, and use restrictions.

Process—program in execution, consisting of the executable program, the program's data and stack, its program counter, stack pointer and other registers, and all other information needed to execute the program.²⁴

Data—representations of facts, concepts, or instructions in a manner suitable for communication, interpretation, or processing by humans or by automatic means.²⁵ Data can be in the form of *files* in a computer's volatile memory or nonvolatile memory, or in a data storage device, or in the form of *data in transit* across a transmission medium.

Component—one of the parts that make up a computer or network.²⁶

Computer—device that consists of one or more associated components, including processing units and peripheral units, that is controlled by internally stored programs and that can perform substantial computations, including numerous arithmetic operations or logic operations, without human intervention during execution. Note: may be stand-alone or may consist of several interconnected units.²⁷

Network—interconnected or interrelated group of host computers, switching elements, and interconnecting branches.²⁸

Internetwork—network of networks.

8.4.2 Attacks. Sometimes an event that occurs on a computer or network is part of a series of steps intended to result in something that is not authorized to happen. This event is then considered part of an *attack*. An attack has three elements.

1. It is made up a series of steps taken by an *attacker*. Among these steps is an action directed at a target (an *event*, as described in the previous section) as well as the use of some *tool* to exploit a *vulnerability*.

8 · 12 USING A COMMON LANGUAGE FOR COMPUTER SECURITY INCIDENT

2. An attack is intended to achieve an *unauthorized result* as viewed from the perspective of the owner or administrator of the system involved.
3. An attack is a series of *intentional* steps initiated by the attacker. This differentiates an attack from something that is inadvertent.

We define an attack in this way:

Attack—a series of steps taken by an attacker to achieve an unauthorized result.

Exhibit 8.3 presents a matrix of possible attacks, based on our experience. Attacks have five parts that depict the logical steps an attacker must take. An attacker uses a (1) *tool* to exploit a (2) *vulnerability* to perform an (3) *action* on a (4) *target* in order to achieve an (5) *unauthorized result*. To be successful, an attacker must find one or more paths that can be connected (attacks), perhaps simultaneously or repeatedly. The first two steps in an attack, *tool* and *vulnerability*, are used to cause an *event* (*action* directed at a *target*) on a computer or network. The logical end of a successful attack is an *unauthorized result*. If the logical end of the previous steps is an *authorized result*, then an attack has not taken place.

The concept of *authorized* versus *unauthorized* is key to understanding what differentiates an attack from the normal events that occur. It is also a system-dependent concept in that what may be authorized on one system may be unauthorized on another.

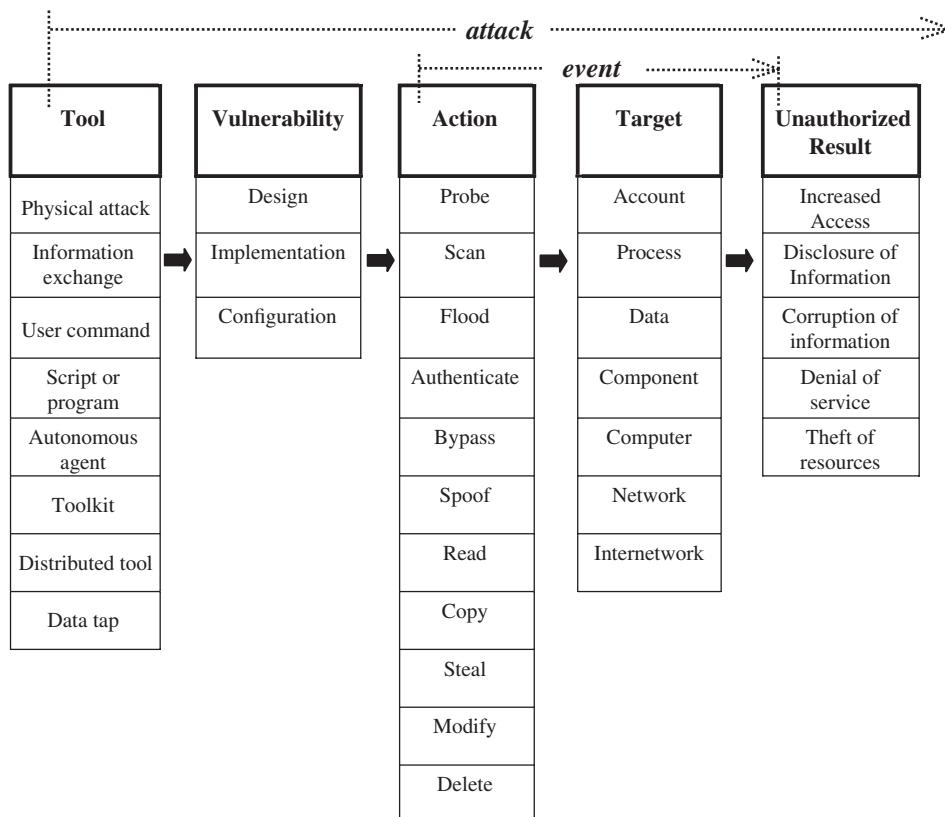


EXHIBIT 8.3 Computer and Network Attacks

COMPUTER SECURITY INCIDENT INFORMATION TAXONOMY 8 · 13

For example, some services, such as anonymous File Transfer Protocol (FTP), may be enabled on some systems and not on others. Even actions that are normally viewed as hostile, such as attempts to bypass access controls to gain entry into a privileged account, may be authorized in special circumstances, such as during an approved test of system security or in the use of a “back door” during development. System owners or their administrators make the determination of what actions they consider authorized for their systems by establishing a security policy.²⁹ Here are the definitions for authorized and unauthorized.

Authorized—approved by the owner or administrator.

Unauthorized—not approved by the owner or administrator.

The steps *action* and *target* in Exhibit 8.1 are the two parts of an event as discussed in Section 8.4.1. The following sections discuss the other steps: *tool*, *vulnerability*, and *unauthorized result*.

8.4.2.1 Tool. The first step in the sequence that leads attackers to their unauthorized results is the *tool* used in the attack. A tool is some means that can be used to exploit a vulnerability in a computer or network. Sometimes a tool is simple, such as a user command or a physical attack. Other tools can be very sophisticated and elaborate, such as a Trojan horse program, computer virus, or distributed tool. We define *tool* in this way.

Tool—means of exploiting a computer or network vulnerability.

The term “tool” is difficult to define more specifically because of the wide variety of methods available to exploit vulnerabilities in computers and networks. When authors make lists of methods of attack, often they are actually making lists of tools. Based on our experience, these categories of tools are currently an exhaustive list. (See Exhibit 8.3)

Physical attack—means of physically stealing or damaging a computer, network, its components, or its supporting systems (e.g., air conditioning, electric power, etc.).

Information exchange—means of obtaining information either from other attackers (e.g., through an electronic bulletin board) or from the people being attacked (commonly called social engineering).

User command—means of exploiting a vulnerability by entering commands to a process through direct user input at the process interface. An example is entering UNIX commands through a telnet connection or commands at a protocol’s port.

Script or program—means of exploiting a vulnerability by entering commands to a process through the execution of a file of commands (script) or a program at the process interface. Examples are a shell script to exploit a software bug, a Trojan horse log-in program, or a password-cracking program.

Autonomous agent—means of exploiting a vulnerability by using a program or program fragment that operates independently from the user. Examples are computer viruses or worms.

8 · 14 USING A COMMON LANGUAGE FOR COMPUTER SECURITY INCIDENT

Toolkit—software package that contains scripts, programs, or autonomous agents that exploit vulnerabilities. An example is the widely available toolkit called *rootkit*.

Distributed tool—tool that can be distributed to multiple hosts, which then can be coordinated to anonymously perform an attack on the target host simultaneously after some time delay.

Data tap—means of monitoring the electromagnetic radiation emanating from a computer or network using an external device.

With the exception of the physical attack, information exchange, and data tap categories, each of the tool categories may contain the other tool categories *within* itself. For example, toolkits contain scripts, programs, and sometimes autonomous agents. So when a *toolkit* is used, the *script or program* category is also included. *User commands* also must be used for the initiation of scripts, programs, autonomous agents, toolkits, and distributed tools. In other words, there is an order to some of the categories in the tools block, from the simple user command category to the more sophisticated distributed tools category. In describing or classifying an attack, generally a choice must be made among several alternatives within the tools block. We chose to classify according to the *highest* category of tool used, which makes the categories mutually exclusive in practice.

8.4.2.2 Vulnerability. To reach the desired result, an attacker must take advantage of a computer or network *vulnerability*.

Vulnerability—weakness in a system allowing unauthorized action.³⁰

A vulnerability in software is an error that arises in different stages of development or use.³¹ This definition can be used to give us three categories of vulnerabilities:

Design vulnerability—vulnerability inherent in the design or specification of hardware or software whereby even a perfect implementation will result in a vulnerability.

Implementation vulnerability—vulnerability resulting from an error made in the software or hardware implementation of a satisfactory design.

Configuration vulnerability—vulnerability resulting from an error in the configuration of a system, such as having system accounts with default passwords, having “world write” permission for new files, or having vulnerable services enabled.³²

8.4.2.3 Unauthorized Result. As shown in Exhibit 8.3, the logical end of a successful attack is an *unauthorized result*. At this point, an attacker has used a tool to exploit a vulnerability in order to cause an event to take place.

Unauthorized result—unauthorized consequence of an event.

If successful, an attack will result in one of the following³³:

Increased access—unauthorized increase in the domain of access on a computer or network.

Disclosure of information—dissemination of information to anyone who is not authorized to access that information.

Corruption of information—unauthorized alteration of data on a computer or network.

COMPUTER SECURITY INCIDENT INFORMATION TAXONOMY 8 · 15

Denial of service—intentional degradation or blocking of computer or network resources.

Theft of resources—unauthorized use of computer or network resources.

8.4.3 Full Incident Information Taxonomy. Often attacks on computers and networks occur in a distinctive group that we would classify as being part of one *incident*. What makes these attacks a distinctive group is a combination of three factors, each of which we may only have partial information about.

1. There may be one attacker, or there may be several attackers who are related in some way.
2. The attacker(s) may use similar attacks, or they may be trying to achieve a distinct or similar objective.
3. The sites involved in the attacks and the timing of the attacks may be the same or may be related.

Here is the definition of *incident*:

Incident—group of attacks that can be distinguished from other attacks because of the distinctiveness of the attackers, attacks, objectives, sites, and timing.

The three parts of an incident are shown in simplified form in Exhibit 8.4, which shows that an attacker, or group of attackers, achieves objectives by performing attacks. An incident may comprise one single attack or multiple attacks, as illustrated by the return loop in the figure.

Exhibit 8.5 shows the full incident information taxonomy. It shows the relationship of events to attacks and attacks to incidents, and suggests that preventing attackers from achieving objectives could be accomplished by ensuring that an attacker cannot make any complete connections through the seven steps depicted. For example, investigations could be conducted of suspected terrorist *attackers*, systems could be searched periodically for attacker *tools*, system *vulnerabilities* could be patched, access controls could be strengthened to prevent *actions* by an attacker to access a *targeted* account, files could be encrypted so as not to *result* in disclosure, and a public education program could be initiated to prevent terrorists from achieving an *objective* of political gain.

8.4.3.1 Attackers and Their Objectives. People attack computers. They do so through a variety of methods and for a variety of objectives. What distinguishes the categories of attackers is a combination of who they are and their *objectives* (what they want to accomplish).

Attacker—individual who attempts one or more attacks in order to achieve an objective.

Objective—purpose or end goal of an incident.



EXHIBIT 8.4 Simplified Computer and Network Incident

8 · 16 USING A COMMON LANGUAGE FOR COMPUTER SECURITY INCIDENT

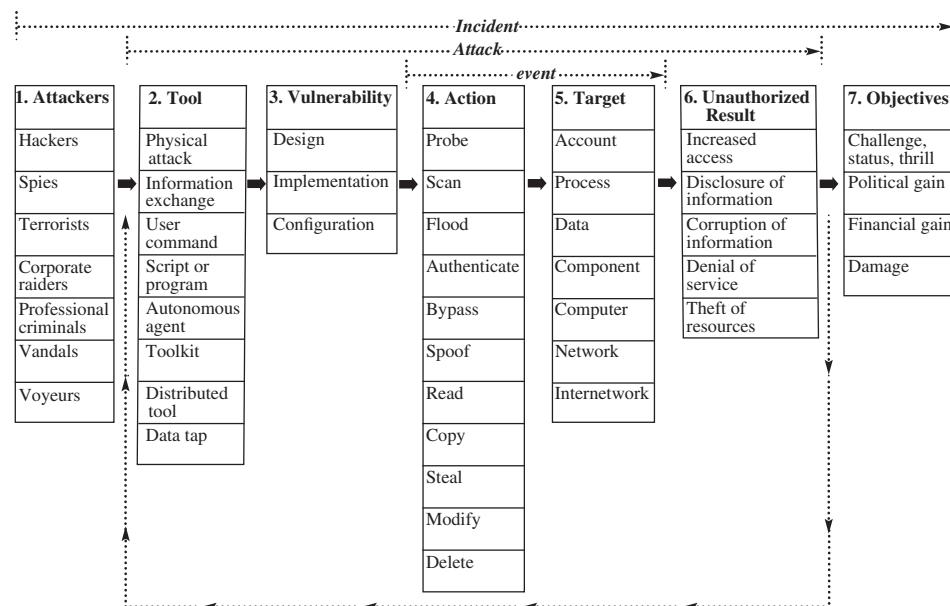


EXHIBIT 8.5 Computer and Network Incident Information Taxonomy

Based on their objectives, we have divided attackers into a number of categories:

Hackers—attackers who attack computers for challenge, status, or the thrill of obtaining access. (Note: We have elected to use the term “hacker” because it is common and widely understood. We realize that the term’s more positive connotation was once more widely accepted.)

Spies—attackers who attack computers for information to be used for political gain.

Terrorists—attackers who attack computers to cause fear, for political gain.

Corporate raiders—employees (attackers) who attack competitors’ computers for financial gain.

Professional criminals—attackers who attack computers for personal financial gain.

Vandals—attackers who attack computers to cause damage.

Voyeurs—attackers who attack computers for the thrill of obtaining sensitive information.

These seven categories of attackers and their four categories of objectives as shown in the leftmost and rightmost blocks of Exhibit 8.5 are fundamental to the difference between *incidents* and *attacks*. This difference is summed up in the phrase “attackers use attacks to achieve objectives.”

8.5 ADDITIONAL INCIDENT INFORMATION TERMS. The taxonomy of the last section presented all of the terms in the common language for computer security that describe how attackers achieve objectives during an incident. However, some other, more general terms are required to fully describe an incident. The next sections discuss these terms.

ADDITIONAL INCIDENT INFORMATION TERMS 8 · 17

8.5.1 Success and Failure. Information on success or failure can be recorded at several levels in the overall taxonomy. In the broadest sense, overall success or failure is an indication of whether one or more attackers have achieved one or more objectives. A narrower focus would be to determine the success or failure of an individual attack by evaluating whether the attack leads to an unauthorized result. Information on success or failure, however, may simply not be known. For example, an attempt to log in to the root or superuser account on a system may be classified as a *success* a *failure*, or as being *unknown*.

8.5.2 Site and Site Name. “Site” is the common term used to identify Internet organizations as well as physical locations. A “site” is also the organizational level of the site administrator or other authority with responsibility for the computers and networks at that location.

The term “site name” refers to a portion of the fully qualified domain name in the Internet’s Domain Name System (DNS). For sites in the United States, site names generally are at the second level of the DNS tree. Examples would be *cmu.edu* or *wIDGETS.com*. In other countries, the site name is the third or lower level of the DNS tree, such as *wIDGETS.co.uk*. Some site names occur even farther down the DNS tree. For example, a school in Colorado might have a site name of *myschool.k12.co.us*.

Here are the definitions of site and site name.

Site—organizational level with responsibility for security events; the organizational level of the site administrator or other authority with responsibility for the computers and networks at that location.

Site name—portion of the fully qualified domain name that corresponds to a site.

Some organizations, such as larger universities and companies, are large enough to be physically divided into more than one location, with separate administration. This separation cannot easily be determined. Therefore, often these different locations must be treated as one site.

8.5.3 Other Incident Terms. Several additional terms are necessary to fully describe actual Internet incidents. The first of these terms concern dates.

Reporting date—first date that the incident was reported to a response team or other agency or individuals collecting data.

Starting date—date of the first known incident activity.

Ending date—date of the last known incident activity.

Several terms concern the sites involved.

Number of sites—overall number of sites known to have reported or otherwise to have been involved in an incident.

Reporting sites—site names of sites known to have reported an incident.

Other sites—site names of sites known to have been involved in an incident but that did not report the incident.

For most incident response teams, actual site names are considered sensitive information. In our research, in order to protect the identities of the sites associated with an incident, we sanitize the site information by coding the site names prior to public

8 · 18 USING A COMMON LANGUAGE FOR COMPUTER SECURITY INCIDENT

release. An example would be to replace a site name, such as the fictitious *widgets.com*, with numbers and the upper-level domain name, such as *123.com*.

Response teams often use incident numbers to track incidents and to identify incident information.

Incident number—reference number used to track an incident or identify incident information.

The last term we found to be of use is *corrective action*, which indicates those actions taken in the aftermath of an incident. These actions could include changing passwords, reloading systems files, talking to the intruders, or even criminal prosecution. Information on corrective actions taken during or after an incident is difficult to obtain for incident response teams, since response team involvement generally is limited to the early stages of an incident. CERT/CC records indicate that the variety of corrective actions is extensive, and a taxonomy of corrective actions may be a desirable future expansion of the common language.

Corrective action—action taken during or after an incident to prevent further attacks, repair damage, or punish offenders.

8.6 HOW TO USE THE COMMON LANGUAGE. Two things are important to emphasize about using the common language for computer security incident information. First, the common language really is a high-level set of terms. As such, it will not settle all the disputes about everything discussed concerning computer security incidents. For example, the common language includes “autonomous agent” as a term (a category of tool). Autonomous agents include *computer viruses*, *worms*, and the like, regardless of how those specific terms might be defined. In other words, the common language does not try to settle disputes on what should or should not be considered a *computer virus* but rather deals at a higher level of abstraction (“autonomous agent”) where, it is hoped, there can be more agreement and standardization. Stated another way, participants in the Common Language Project workshops anticipated that individuals and organizations would continue to use their own terms, which may be more specific in both meaning and use. The common language has been designed to enable these lower-level terms to be classified *within* the common language structure.

The second point to emphasize is that the common language, even though it presents a taxonomy, does not classify an incident (or individual attacks) as any *one* thing. Classifying computer security *attacks* or *incidents* is difficult because attacks and incidents are a *series of steps* that an attacker must take. In other words, attacks and incidents are not just *one* thing but rather a *series* of things. That is why I say the common language provides a taxonomy for computer security incident *information*.

An example of the problem is found in the popular and simple taxonomies often used to attempt to classify incidents. They appear as a list of single, defined terms. The following terms from Icove, Seger, and VonStorch provide an example.³⁴

Covert channels	Data diddling	Degradation of service
Denial of service	Dumpster diving	Eavesdropping on emanations
Excess privileges	Harassment	IP spoofing
Logic bombs	Masquerading	Password sniffing
Salamis	Scanning	Session hijacking
Software piracy	Timing attacks	Traffic analysis
Trap doors	Trojan horses	Tunneling
Unauthorized data copying	Viruses and worms	Wiretapping

HOW TO USE THE COMMON LANGUAGE 8 · 19

Lists of terms are *not* satisfactory taxonomies for classifying actual attacks or incidents. They fail to have most of the six characteristics of a satisfactory taxonomy. First, the terms tend not to be mutually exclusive. For example, the terms “virus” and “logic bomb” are generally found on these lists, but a virus may *contain* a logic bomb, so the categories overlap. Actual attackers generally also use multiple methods so their attacks would have to be classified into multiple categories. This makes classification ambiguous and difficult to repeat.

A more fundamental problem is that, assuming that an exhaustive and mutually exclusive list could be developed, the taxonomy would be unmanageably long and difficult to apply. It also would not indicate any relationship between different types of attacks. Finally, none of these lists has become widely accepted, partly because it is difficult to agree on the definition of terms. In fact, many different definitions of terms are in common use.

The fundamental problems with these lists (and their variations) are that most incidents involve multiple attacks, and attacks involve multiple steps. As a result, information about the typical incident must be classified in multiple categories. For example, one of the attacks in an incident might be a flood of a host resulting in a denial of service. But this same incident might involve the exploitation of a vulnerability to compromise the host computer that was the specific origin of the flood. Should this be classified as a flood? As a root compromise? As a denial-of-service attack? In reality, the incident should be classified in all of these categories. In other words, this incident has multiple classifications.

In summary, in developing the common language, we have found that, with respect to *attacks* and *incidents*, we can really *only* hope to (1) present a common set of high-level terms that are in general use and have common definitions and (2) present a logical structure to the terms that can be used to classify information *about* an incident or attack *with respect to specific categories*.

Some examples may make this clear. As discussed earlier, most of the information about actual attacks and incidents is in the form of textual records. In a typical incident record at the CERT/CC, three observations might be reported:

1. We found *rootkit* on host xxx.xxx.
2. A flood of email was sent to account xxx@xxx.xxx, which crashed the mail server.
3. We traced the attack back to a teenager in Xyz city, who said he was not trying to cause any damage, just trying to see if he could break in.

For observation 1, we would classify *rootkit* in the “toolkit” category under “Tool” and the hostname in the “computer” category under “Target.” For observation 2, the “email flood” is a specific instantiation in the “flood” category under “Action” as well as in the “denial-of-service” category under “Unauthorized Result.” There is ambiguity as to the target for observation 2: Is it the account or the computer? As a practical matter, the observations would be classified as both, since information is available on both. For observation 3, it could be inferred that this is a “hacker” seeking “challenge, status, or thrill.”

What does this taxonomic process provide that is of practical value? First, the taxonomy helps us communicate to others what we have found. When we say that *rootkit* is a type of toolkit, then our common set of terms (“common language”) provides us the general understanding of what we mean. When it is said that 22 percent of incidents reported to CERT/CC from 1988 through 1995 involved various problems with

8 · 20 USING A COMMON LANGUAGE FOR COMPUTER SECURITY INCIDENT

passwords (a correct statistic³⁵), then the taxonomy has proven useful in communicating valuable information.

The application of the taxonomy, in fact, is a four-step process that can be used to determine the biggest security problems. Specifically, the process is to:

1. Take observations from fragmentary information in incident reports.
2. Classify those observations.
3. Perform statistical studies of these data.
4. Use this information to determine the best course(s) of action.

Over time, the same process can be used to determine the effects of these actions.

Two more points are important to emphasize about this taxonomy. First, an *attack* is a process that, with enough information, is *always* classified in multiple categories. For example: in a “Tool” category, in a “Vulnerability” category, in an “Action” category, in a “Target” category, and in an “Unauthorized Result” category. Second, an *incident* can involve multiple, perhaps thousands, of attacks. As such, the information gathered in an incident theoretically could be classified correctly into *all* of the taxonomy categories.

Within these guidelines, the common language for computer security incidents has proven to be a useful and increasingly accepted tool to gather, exchange, and compare computer security information. The taxonomy itself has proven to be simple and straightforward to use.

8.7 NOTES

1. E. G. Amoroso, *Fundamentals of Computer Security Technology* (Upper Saddle River, NJ: Prentice-Hall PTR, 1994), p. 2.
2. Deborah Russell and G. T. Gangemi Sr., *Computer Security Basics* (Sebastopol, CA: O’Reilly & Associates, 1991), p. 79.
3. Bill McKelvey, *Organization Systematics: Taxonomy, Evolution, Classification* (Berkeley: University of California Press, 1982), p. 3.
4. John D. Howard, “An Analysis of Security Incidents on the Internet, 1989–1995” (PhD diss., Department of Engineering and Public Policy, Carnegie Mellon University, Pittsburgh, PA, April 1997). Also available online at www.cert.org/archive/pdf/JHThesis.pdf.
5. Ivan Victor Krsul, “Software Vulnerability Analysis” (PhD diss., Computer Sciences Department, Purdue University, Lafayette, IN, May 1998), p. 12.
6. Howard, “Analysis of Security Incidents on the Internet.”
7. John Radatz, ed., *The IEEE Standard Dictionary of Electrical and Electronics Terms*, 6th ed. (New York: Institute of Electrical and Electronics Engineers, 1996), p. 1087.
8. Amoroso, *Fundamentals of Computer Security Technology*, p. 34.
9. Radatz, *IEEE Standard Dictionary*, p. 373.
10. Radatz, *IEEE Standard Dictionary*, p. 373.
11. Radatz, *IEEE Standard Dictionary*, p. 11.
12. Radatz, *IEEE Standard Dictionary*, p. 5.
13. Derek Atkins et al., *Internet Security Professional Reference* (Indianapolis: New Riders Publishing, 1996), p. 258.

NOTES 8 · 21

14. Radatz, *IEEE Standard Dictionary*, p. 947, and K. M. Jackson and J. Hruska, eds., *Computer Security Reference Book* (Boca Raton, FL: CRC Press, 1992), p. 916.
15. Merriam-Webster, *Merriam-Webster's Collegiate Dictionary*, 10th ed. (Springfield, MA: Author, 1996), pp. 77, 575, 714, and Radatz, *IEEE Standard Dictionary*, p. 57.
16. Merriam-Webster's Collegiate Dictionary, p. 157.
17. Radatz, *IEEE Standard Dictionary*, p. 630, and Atkins et al., *Internet Security*, p. 258.
18. Radatz, *IEEE Standard Dictionary*, p. 877.
19. Radatz, *IEEE Standard Dictionary*, p. 224.
20. Radatz, *IEEE Standard Dictionary*, p. 661.
21. Radatz, *IEEE Standard Dictionary*, p. 268.
22. Andrew S. Tanenbaum, *Modern Operating Systems* (Englewood Cliffs, NJ: Prentice-Hall, 1992), p. 12.
23. Radatz, *IEEE Standard Dictionary*, p. 683.
24. Tanenbaum, *Modern Operating Systems*, p. 12, and Radatz, *IEEE Standard Dictionary*, p. 822.
25. Radatz, *IEEE Standard Dictionary*, p. 250.
26. Radatz, *IEEE Standard Dictionary*, p. 189.
27. Radatz, *IEEE Standard Dictionary*, p. 192.
28. Radatz, *IEEE Standard Dictionary*, p. 683.
29. Krsul, "Software Vulnerability Analysis," pp. 5–6.
30. National Research Council, *Computers at Risk: Safe Computing in the Information Age* (Washington, DC: National Academy Press, 1991), p. 301; and Amoroso, *Fundamentals of Computer Security Technology*, p. 2.
31. Krsul, *Software Vulnerability Analysis*, pp. 10–11.
32. Atkins et al., *Internet Security*, p. 196.
33. Amoroso, *Fundamentals of Computer Security Technology*, pp. 3–4, 31; Russell and Gangemi, *Computer Security Basics*, pp. 9–10; and Frederick B. Cohen, *Protection and Security on the Information Superhighway* (New York: John Wiley & Sons, 1995), pp. 55–56.
34. David Icove, Karl Seger, and William VonStorch, *Computer Crime: A Crime-fighter's Handbook* (Sebastopol, CA: O'Reilly & Associates, 1995), pp. 31–52; Cohen, *Protection and Security on the Information Superhighway*, pp. 40–54 (39 terms); and Frederick B. Cohen, "Information System Attacks: A Preliminary Classification Scheme," *Computers and Security* 16, No. 1 (1997): 29–46 (96 terms).
35. Howard, "Analysis of Security Incidents on the Internet," p. 100.

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 9

MATHEMATICAL MODELS OF COMPUTER SECURITY

Matt Bishop

9.1 WHY MODELS ARE IMPORTANT	9·1	9.3.3 Role-Based Access-Control Models and Groups	9·7
9.2 MODELS AND SECURITY	9·3	9.3.4 Summary	9·9
9.2.1 Access-Control Matrix Model	9·3	9.4 CLASSIC MODELS	9·9
9.2.2 Harrison, Ruzzo, and Ullman and Other Results	9·5	9.4.1 Bell-LaPadula Model	9·9
9.2.3 Typed Access-Control Model	9·6	9.4.2 Biba's Strict Integrity Policy Model	9·12
		9.4.3 Clark-Wilson Model	9·14
		9.4.4 Chinese Wall Model	9·17
		9.4.5 Summary	9·18
9.3 MODELS AND CONTROLS	9·6	9.5 OTHER MODELS	9·18
9.3.1 Mandatory and Discretionary Access-Control Models	9·6	9.6 CONCLUSION	9·19
9.3.2 Originator-Controlled Access-Control Model and DRM	9·6	9.7 FURTHER READING	9·19
		9.8 NOTES	9·21

9.1 WHY MODELS ARE IMPORTANT. When you drive a new car, you look for specific items that will help you control the car: the accelerator, the brake, the shift, and the steering wheel. These exist on all cars and perform the function of speeding the car up, slowing it down, and turning it left and right. This forms a model of the car. With these items properly working, you can make a convincing argument that the model correctly describes what a car must have in order to move and be steered properly.

A model in computer security serves the same purpose. It presents a general description of a computer system (or collection of systems). The model provides a definition of “protect” (e.g., “keep confidential” or “prevent unauthorized change to”) and conditions under which the protection is provided. With mathematical models, the conditions can be shown to provide the stated protection. This provides a high degree of assurance that the data and programs are protected, assuming the model is implemented correctly.

This last point is critical. To return to our car analogy, notice the phrase “with these items properly working.” This also means that the average driver must be able to work

9 · 2 MATHEMATICAL MODELS OF COMPUTER SECURITY

them correctly. In most, if not all, cars the model is implemented in the obvious way: The accelerator pedal is to the right of the brake pedal, and speeds the car up; the brake pedal slows it down; and turning the steering wheel moves the car to the left or right, depending on the direction that the wheel is turned. The average driver is familiar with this implementation and so can use it properly. Thus, the model and the implementation together show that this particular car can be driven.

Now, suppose that the items are implemented differently. All the items are there, but the steering wheel is locked so it cannot be turned. Even though the car has all the parts that the model requires, they do not work the way the model requires them to work. The implementation is incorrect, and the argument that the model provides does not apply to this car, because the model makes assumptions—like the steering wheel turning—that are incorrect for this car. Similarly, in all the models we present in this chapter, the reader should keep in mind the assumptions that the models make. When one applies these models to existing systems, or uses them to design new systems, one must ensure that the assumptions are met in order to gain the assurance that the model provides.

This chapter presents several mathematical models, each of which serves a different purpose. We can divide these models into several types.

The first set of models is used to determine under what conditions one can prove types of systems secure. The access-control matrix model presents a general description of a computer system that this type of model uses, and it will give some results about the decidability of security in general and for particular classes of systems.

The second type of model describes how the computer system applies controls. The mandatory access-control model and the discretionary access-control model form the basis for components of the models that follow. The originator-controlled access-control model ties control of data to the originator rather than the owner, and has obvious applications for digital rights management systems. The role-based access-control model uses job function, rather than identity, to provide controls and so can implement the principle of least privilege more effectively than many models.

The next few models describe confidentiality and integrity. The Bell-LaPadula model describes a class of systems designed to protect confidentiality and was one of the earliest, and most influential, models in computer security. The Biba model's strict integrity policy is closely related to the Bell-LaPadula model and is in widespread use today; it is applied to programs to determine when their output can be trusted. The Clark-Wilson model is also an integrity model, but it differs fundamentally from Biba's model because the Clark-Wilson model describes integrity in terms of processes and process management rather than in terms of attributes of the data.

The fourth type of model is the hybrid model. The Chinese Wall model examines conflicts of interest, and is an interesting mix of both confidentiality and integrity requirements. This type of model arises when many real-world problems are abstracted into mathematical representations, for example, when analyzing protections required for medical records and for the process of recordation of real estate.¹

The main goal of this chapter is to provide the reader with an understanding of several of the main models in computer security, of what these models mean, and of when they are appropriate to use. An ancillary goal is to make the reader sensitive to how important assumptions in computer security are. Dorothy Denning said it clearly and succinctly in her speech when accepting the National Computer Systems Security Award in 1999:

The lesson I learned was that security models and formal methods do not establish security. They only establish security with respect to a model, which by its very nature is extremely

MODELS AND SECURITY 9 · 3

simplistic compared to the system that is to be deployed, and is based on assumptions that represent current thinking. Even if a system is secure according to a model, the most common (and successful) attacks are the ones that violate the model's assumptions. Over the years, I have seen system after system defeated by people who thought of something new.²

Given this, the obvious question is: Why are models important? Models provide a framework for analyzing systems and for understanding where to focus our security efforts: on either validating the assumptions or ensuring that the assumptions are met in the environment in which the system exists. The mechanisms that do this may be technical; they may be procedural. Their quality determines the security of the system. So the model provides a basis for asserting that, if the mechanisms work correctly, then the system is secure—and that is far better than simply implementing security mechanisms without understanding how they work together to meet security requirements.

9.2 MODELS AND SECURITY. Some terms recur throughout our discussion of models.

- A *subject* is an active entity, such as a process or a user.
- An *object* is a passive entity, such as a file.
- A *right* describes what a subject is allowed to do to an object; for example, the *read* right gives permission for a subject to read a file.
- The *protection state* of a system simply refers to the rights held by all subjects on the system.

The precise meaning of each right varies from actual system to system. For example, on Linux systems, if a process has *write* permission for a file, that process can alter the contents of the file. But if a process has *write* permission for a directory, that process can create, delete, or rename files in that directory. Similarly, having *read* rights over a process may mean the possessor can participate as a recipient of interprocess communications messages originating from that process. The point is that the meaning of the rights depends on the interpretation of the system involved. The assignment of meaning to the rights used in a mathematical model is called *instantiating the model*.

The first model we explore is the foundation for much work on the fundamental difficulty of analyzing systems to determine whether they are secure.

9.2.1 Access-Control Matrix Model. The *access-control matrix model*³ is perhaps the simplest model in computer security. It consists of a matrix, the rows of which correspond to subjects and the columns of which correspond to entities (subjects and objects). Each entry in the matrix contains the set of rights that the subject (row) has over the entity (column). For example, the access-control matrix in Exhibit 9.1 shows a system with two processes and two files. The first process has own rights over itself; read rights over the second process; read and execute rights over the first file; and read, write, and own rights over the second file. The second process can write to the first process; owns itself; can read, write, execute, and owns the first file; and can read the second file.

The access-control matrix captures a protection state of a system. But systems evolve; their protection state does not remain constant. So the contents of the

9.4 MATHEMATICAL MODELS OF COMPUTER SECURITY

	Process 1	Process 2	File 1	File 2
Process 1	own	read	read, execute	read, write, own
Process 2	write	own	read, write, execute, own	read

EXHIBIT 9.1 Example Access-Control Matrix with Two Processes and Two Files

access-control matrix must change to reflect this evolution. Perhaps the simplest set of rules for changing the access-control matrix are these *primitive operations*⁴:

- **Create subject** s creates a new row and column, both labeled s
- **Create object** o creates a new column labeled o
- **Enter r into $A[s, o]$** adds the right r into the entry in row s and column o ; it corresponds to giving the subject s the right r over the entity o
- **Delete r from $A[s, o]$** removes the right r from the entry in row s and column o ; it corresponds to deleting the subject s 's right r over the entity o
- **Destroy subject** s removes the row and column labeled s
- **Destroy object** o removes the column labeled o

These operations can be combined into *commands*. The next command creates a file f and gives the process p read and own rights over that file:

```
command createread(p, f)
    create object f
    enter read into A[p, f]
    enter own into A[p, f]
end.
```

A *mono-operational* command consists of a single primitive operation. For example, the command

```
command grantwrite(p, f)
    enter write into A[p, f]
end.
```

which gives p write rights over f , is mono-operational.

Commands may include conditions. For example, the next command gives the subject p execute rights over a file f if p has read rights over f :

```
command grantexec(p, f)
    if read in A[p, f] then
        enter execute into A[p, f]
end.
```

If p does not have read rights over f when this command is executed, it does nothing. This command has one condition and so is called *monoconditional*. *Biconditional* commands have two conditions joined by *and*:

```
command copyread(p, q, f)
    if read in A[p, f] and own in A[p, f] then
        enter read into A[q, f]
end.
```

MODELS AND SECURITY 9 · 5

This command gives a subject q read rights over the object f if the subject p owns f and has read rights over f .

Commands may have conditions only at the beginning, and if the condition is false, the command terminates. Commands may contain other commands as well as primitive operations.

If all commands in a system are mono-operational, the system is said to be *mono-operational*; if all the commands are monoconditional or biconditional, then the system is said to be *monoconditional* or *biconditional*, respectively. Finally, if the system has no commands that use the *delete* or *destroy* primitive operations, the system is said to be *monotonic*.

The access-control matrix provides a theoretical basis for two widely used security mechanisms: access-control lists and capability lists. In the realm of modeling, it provides a tool to analyze the difficulty of determining how secure a system is.

9.2.2 Harrison, Ruzzo, and Ullman and Other Results. The question of how to test whether systems are secure is critical to understanding computer security. Define *secure* in the simplest possible way: A system is secure with respect to a generic right r if that right cannot be added to an entity in the access-control matrix unless that square already contains it. In other words, a system is secure with respect to r if r cannot leak into a new entry in the access-control matrix. The question then becomes:

Safety Question. Is there an algorithm to determine whether a given system with initial state σ is secure with respect to a given right?

In the general case:

Theorem (Harrison, Ruzzo, and Ullman [HRU] Result).⁵ The safety question is undecidable.

The proof is to reduce the halting problem to the safety question.⁶ This means that, if the safety question were decidable, so would the halting problem be. But the undecidability of the halting problem is well known,⁷ so the safety problem must also be undecidable.⁸

These results mean that one cannot develop a general algorithm for determining whether systems are secure. One can do so in limited cases, however, and the models that follow are examples of such cases. The characteristics that classes of systems must meet in order for the safety question to be decidable are not yet known fully, but for specific classes of systems, the safety question can be shown to be decidable. For example:

Theorem.⁹ There is an algorithm that will determine whether mono-operational systems are secure with respect to a generic right r .

But these classes are sensitive to the commands allowed:

Theorem.¹⁰ The safety question for monotonic systems is undecidable.

Limiting the set of commands to biconditional commands does not help:

Theorem.¹¹ The safety question for biconditional monotonic systems is undecidable.

But limiting them to monoconditional operations:

Theorem.¹² There is an algorithm that will determine whether monoconditional monotonic systems are secure with respect to a generic right r .

In fact, adding the *delete* primitive operation does not affect this result (although the proof is different):

Theorem.¹³ There is an algorithm that will determine whether monotonic systems that do not use the *destroy* primitive operations are secure with respect to a generic right r .

9 · 6 MATHEMATICAL MODELS OF COMPUTER SECURITY

9.2.3 Typed Access-Control Model. A variant of the access-control matrix model adds type to the entities. The *typed access-control matrix model*, called TAM,¹⁴ associates a type with each entity and modifies the rules for matrix manipulation accordingly. This notion allows entities to be grouped into finer categories than merely *subject* and *object*, and enables a slightly different analysis than the HRU result suggests.

In TAM, a rule set is *acyclic* if neither an entity E nor any of its descendants can create a new entity with the same type as E . Given that definition:

Theorem.¹⁵ There is an algorithm that will determine whether acyclic, monotonic typed matrix models are secure with respect to a generic right r .

Thus, a system being acyclic and monotonic is sufficient to make the safety question decidable. But we still do not know exactly what properties are *necessary* to make the safety question decidable.

We now turn to models that have direct application to systems and environments and that focus on more complex definitions of “secure” and the mechanisms needed to achieve them.

9.3 MODELS AND CONTROLS. Models of computer security focus on control: who can access files and resources, and what types of access are allowed. The next characterizations of these controls organize them by flexibility of use and by the roles of the entities controlling the access. These are essential to understanding how more sophisticated models work.

9.3.1 Mandatory and Discretionary Access-Control Models. Some access-control methods are rule based; that is, users have no control over them. Only the system or a special user called (for example) the *system security officer* (SSO) can change them. The government classification system works this way. Someone without a clearance is forbidden to read TOP SECRET material, even if the person who has the document wishes to allow it. This rule is called *mandatory* because it must be followed, without exception. Examples of other mandatory rules are the laws in general, which are to be followed as written, and one cannot absolve another of liability for breaking the laws; or the Multics ring-based access-control mechanism, in which accessing a data segment from below the lower bound of the segment’s access bracket is forbidden regardless of the access permissions. This type of access control is called a *mandatory access control*, or MAC. These rules base the access decision on attributes of the subject and object (and possibly other information).

Other access-control methods allow the owner of the entity to control access. For example, a person who keeps a diary decides who can read it. She need not show it to anyone, and if a friend asks to read it, she can say no. Here the owner allows access to the diary at her discretion. This type of control is called *discretionary*. *Discretionary access control*, or DAC, is the most common type of access-control mechanism on computers.

Controls can be (and often are) combined. When mandatory and discretionary controls are combined to enforce a single access-control policy, the mandatory controls are applied first. If they deny access, the system denies access and the discretionary controls need never be invoked. If the mandatory rules permit access, then the discretionary controls are consulted. If both allow the accesses, access is granted.

9.3.2 Originator-Controlled Access-Control Model and DRM. Other types of access controls contain elements of both mandatory and discretionary access

MODELS AND CONTROLS 9 · 7

controls. *Originator-controlled access control*,¹⁶ or ORCON,¹⁷ mechanisms allow the originator to determine who can access a resource or data.

Consider a large government research agency that produces a study of projected hoe-handle sales for the next year. The market for hoe handles is extremely volatile, and if the results of the study leak out prematurely, certain vendors will obtain a huge market advantage. But the study must be circulated to regulatory agencies so they can prepare appropriate regulations that will be in place when the study is released. Thus, the research agency must retain control of the study even as it circulates it among other groups.

More precisely, an originator-controlled access control satisfies two conditions. Suppose an object o is marked as ORCON for organization X . X decides to release o to subjects acting on behalf of another organization Y . Then

1. The subjects to whom the copies of o are given cannot release o to subjects acting on behalf of other organizations without X 's consent; and
2. Any copies of o must bear these restrictions.

Consider a control that implements these requirements. In theory, mandatory access controls could solve this problem. In practice, the required rules must anticipate *all* the organizations to which the data will be made available. This requirement, combined with the need to have a separate rule for each possible set of objects and organizations that are to have access to the object, makes a mandatory access control that satisfies the requirements infeasible. But if the control were discretionary, each entity that received a copy of the study could grant access to its copy without permission of the originator; so originator-controlled access control is neither discretionary nor mandatory.

However, a combination of discretionary and mandatory access controls can implement this control. The mandatory access-control mechanisms forbid the owner from changing access permissions on an object o and require that every copy of that object have the same access-control permissions as are on o . The discretionary access control says that the originator can change the access-control permissions on any copy of o .

As an example of the use of this model in a more popular context, record companies want to control the use of their music. Conceptually, they wish to retain control over the music *after* it is sold in order to prevent owners from distributing unauthorized copies to their friends. Here the originator is the record company and the protected resource is the music.

In practice, originator-controlled access controls are difficult to implement technologically. The problem is that access-control mechanisms typically control access to *entities*, such as files, devices, and other objects. But originator-controlled access control requires that access controls be applied to *information* that is contained in the entities—a far more difficult problem for which there is not yet a generally accepted mechanism.

9.3.3 Role-Based Access-Control Models and Groups. In real life, job function often dictates access permissions. The bookkeeper of an office has free access to the company's bank accounts, whereas the sales people do not. If Anne is hired as a salesperson, she cannot access the company's funds. If she later becomes the bookkeeper, she can access those funds. So the access is conditioned not on the identity of the person but on the role that person plays.

9 · 8 MATHEMATICAL MODELS OF COMPUTER SECURITY

This example illustrates *role-based access control* (RBAC).¹⁸ It assigns a set of roles, called the *authorized roles of the subject s*, to each subject *s*. At any time, *s* may assume at most one role, called the *active role of s*. Then

Axiom. The *rule of role authorization* says that the active role of *s* must be in the set of authorized roles of *s*.

This axiom restricts *s* to assuming those roles that it is authorized to assume. Without it, *s* could assume any role, and hence do anything.

Extending this idea, let the predicate *canexec(s, c)* be true when the subject *s* can execute the command *c*.

Axiom. The *rule of role assignment* says that if *canexec(s, c)* is true for any *s* and any *c*, then *s* must have an active role.

This simply says that in order to execute a command *c*, *s* must have an active role. Without such a role, it cannot execute any commands. We also want to restrict the commands that *s* can execute; the next axiom does this.

Axiom. The *rule of transaction authorization* says that if *canexec(s, c)* is true, then only those subjects with the same role as the active role of *s* may also execute transaction.

This means that every role has a set of commands that it can execute, and if *c* is not in the set of commands that the active role of *s* can execute, then *s* cannot execute it.

As an example of the power of this model, consider two common problems: containment of roles and separation of duty. Containment of roles means that a subordinate *u* is restricted to performing a limited set of commands that a superior *s* can also perform; the superior may also perform other commands. Assign role *a* to the superior and role *b* to the subordinate; as everything a subject with active role *b* can do, a subject with active role *a* can do, we say that role *a* *contains* role *b*. Then we can say that if *a* is an authorized role of *s*, and *a* contains *b*, then *b* is also an authorized role of *s*. Taking this further, if a subject is authorized to assume a role that contains other (subordinate) roles, it can also assume any of the subordinate roles.

Separation of duty is a requirement that multiple entities must combine their efforts to perform a task. For example, a company may require two officers to sign a check for more than \$50,000. The idea is that a single person may breach security, but two people are less likely to combine forces to breach security.¹⁹ One way to handle separation of duty is to require that two distinct roles complete the task and make the roles mutually exclusive. More precisely, let *r* be a role and *meauth(r)*, the mutually exclusive authorization set of *r*, be the set of roles that a subject with authorized role *r* can never assume. Then separation of duty is:

Axiom. The *rule of separation of duty* says that if a role *a* is in the set *meauth(b)*, then no subject for which *a* is an authorized role may have *b* as another authorized role.

This rule is applied to a task that requires two distinct people to complete. The task is broken down into steps that two people are to complete. Each person is assigned a separate role, and each role is in the mutually exclusive authorization set of the other. This prevents either person from completing the task; they must work together, each in their respective role, to complete it.

Roles bear a resemblance to groups, but the goals of groups and roles are different. Membership in a group is defined by essentially arbitrary rules, set by the managers of the system. Membership in a role is defined by job function and is tied to a specific set of commands that are necessary to perform that job function. Thus, a role is a type of group, but a group is broader than a role and need not be tied to any particular set of commands or functions.

CLASSIC MODELS 9 · 9

9.3.4 Summary. The four types of access controls discussed in this section have different focuses. Mandatory, discretionary, and originator-controlled access controls are data-centric, determining access based on the nature or attributes of the data. Role-based access control focuses on the subject's needs. The difference is fundamental.

The principle of least privilege²⁰ says that subjects should have no more privileges than necessary to perform their tasks. Role-based access control, if implemented properly, does this by constraining the set of commands that a subject can execute. The other three controls do this by setting attributes on the data to control access to the data rather than by restricting commands. Mandatory access controls have the attributes set by a system security officer or other trusted process; discretionary access controls, by the owner of the object; and originator-controlled access controls, by the creator or originator of the data.

As noted, these mechanisms can be combined to make the controls easier to use and more precise in application. We now discuss several models that do so.

9.4 CLASSIC MODELS. Three models have played an important role in the development of computer security. The Bell-LaPadula model, one of the earliest formal models in computer security, influenced the development of much computer security technology, and it is still in widespread use. Biba, its analog for integrity, now plays an important role in program analysis. The Clark-Wilson model describes many commercial practices to preserve integrity of data. We examine each of these models in this section.

9.4.1 Bell-LaPadula Model. The Bell-LaPadula model²¹ is a formalization of the famous government classification system using UNCLASSIFIED, CONFIDENTIAL, SECRET, and TOP SECRET levels. We begin by using those four levels to explain the ideas underlying the model and then augment those levels to present the full model. Because the model involves multiple levels, it is an example of a *multilevel security model*.

The four-level version of the model assumes that the levels are ordered from lowest to highest as UNCLASSIFIED, CONFIDENTIAL, SECRET, and TOP SECRET. Objects are assigned levels based on their sensitivity. An object at a higher level is more sensitive than an object at a lower level. Subjects are assigned levels based on what objects they can access. A subject is *cleared* into a level, and that level is called the subject's *security clearance*. An object is *classified* at a level, and that level is called the object's *security classification*. The goal of the classification system is to prevent information from leaking, or flowing downward (e.g., a subject at CONFIDENTIAL should not be able to read information classified TOP SECRET).

For convenience, we write *level(s)* for a subject's security clearance and *level(o)* for an object's security classification. The name of the classification is called a *label*. So an object classified at TOP SECRET has the label TOP SECRET.

Suppose Tom is cleared into the SECRET level. Three documents, called Paper, Article, and Book, are classified as CONFIDENTIAL, SECRET, and TOP SECRET, respectively. As Tom's clearance is lower than Book's classification, he cannot read Book. As his clearance is equal to or greater than Article's and Paper's classification, he can read them.

Definition. The *simple security property* says that a subject *s* can read an object *o* if and only if $\text{level}(o) \leq \text{level}(s)$.

9 · 10 MATHEMATICAL MODELS OF COMPUTER SECURITY

This is sometimes called the *no-reads-up* rule, and it is a mandatory access control.

But that is insufficient to prevent information from flowing downward. Suppose Donna is cleared into the CONFIDENTIAL level. By the simple security property, she cannot read Article because

$$\text{level(Article)} = \text{SECRET} > \text{CONFIDENTIAL} = \text{level(Donna)}.$$

But Tom can read the information in Article and write it on Paper. And Donna can read Paper. Thus, SECRET information has leaked to a subject with CONFIDENTIAL clearance.

To prevent this, Tom must be prevented from writing to Paper:

Definition. The $*$ -property says that a subject s can write an object o if and only if $\text{level}(s) \leq \text{level}(o)$.

This is sometimes called the *no-writes-down* rule, and it too is a mandatory access control. It is also known as the *star property* and the *confinement property*.

Under this rule, as $\text{level(Tom)} = \text{SECRET} > \text{level(Paper)}$, Tom cannot write to Paper. This solves the problem.

Finally, the Bell-LaPadula model allows owners of objects to use discretionary access controls:

Definition. The *discretionary security property* says that a subject s can read an object o only if the access-control matrix entry for s and o contains the read right.

So, in order to determine whether Tom can read Paper, the system checks the simple security property and the discretionary security property. As both hold for Tom and Paper, Tom can read Paper. Similarly, the system checks the $*$ -property to determine whether Tom can write to Paper. As the $*$ -property does not hold for Tom and Paper, Tom cannot write to Paper. Note that the discretionary security property need not be checked, because the relevant mandatory access-control property (the $*$ -property) denies access.

The basic security theorem states that, if a system starts in a secure state, and every operation obeys the three properties, then the system remains secure:

Basic Security Theorem. Let a system Σ have a secure initial state σ_0 . Further, let every command in this system obey the simple security property, the $*$ -property, and the discretionary security property. Then every state σ_i , $i \geq 0$, is also secure.

We can generalize this to an arbitrary number of levels. Let L_0, \dots, L_n be a set of security levels that are linearly ordered (i.e., $L_0 < \dots < L_n$). Then the simple security property, the $*$ -property, and the discretionary security property all apply, as does the Basic Security Theorem. This allows us to have many more than the four levels described.

Now suppose Erin works for the European Department of a government agency, and Don works for the Asia Department for the same agency. Erin and Don are both cleared for SECRET. But some information Erin will see is information that Don has no need to know, and vice versa. Introducing additional security levels will not help here, because then either Don would be able to read all of the documents that Erin could, or vice versa. We need an alternate mechanism.

The alternate mechanism is an expansion of the idea of “security level.” We define a *category* to be a kind of information. A *security compartment* is a pair $(\text{level}, \text{category set})$ and plays the role that the security level did previously.

As an example, suppose the category for the European Department is EUR, and the category for the Asia Department is ASIA. Erin will be cleared into the compartment (SECRET, {EUR}) and Don into the compartment (SECRET, {ASIA}). Documents have security compartments as well. The paper EurDoc may be classified as

CLASSIC MODELS 9 · 11

(CONFIDENTIAL, {EUR}), and the paper AsiaDoc may be (SECRET, {ASIA}). The paper EurAsiaDoc contains information about both Europe and Asia, and so would be in compartment (SECRET, {EUR, ASIA}). As before, we write $level(Erin) = (\text{SECRET}, \{\text{EUR}\})$, $level(\text{EurDoc}) = (\text{CONFIDENTIAL}, \{\text{EUR}\})$, and $level(\text{EurAsiaDoc}) = (\text{SECRET}, \{\text{EUR, ASIA}\})$.

Next, we must define the analog to “greater than.” As noted earlier, security compartments are no longer linearly ordered, because not every pair of compartments can be compared. For example, Don’s compartment is not “greater” than Erin’s, and Erin’s is not “greater” than Don’s. But the classification of EurAsiaDoc is clearly “greater” than that of both Don and Erin.

We compare compartments using the relation dom , for “dominates.”

Definition. Let L and L' be security levels and let C and C' be category sets. Then

$$(L, C)dom(L', C') \text{ if and only if } L' = L \text{ and } C' \subseteq C$$

The dom relation plays the role that “greater than or equal to” did for security levels. Continuing our example, $level(Erin) = (\text{SECRET}, \{\text{EUR}\})$, $dom(\text{CONFIDENTIAL}, \{\text{EUR}\}) = level(\text{EurDoc})$, and $level(\text{EurAsiaDoc}) = (\text{SECRET}, \{\text{EUR, ASIA}\})$ $dom(\text{SECRET}, \{\text{EUR}\}) = level(Erin)$.

We now reformulate the simple security property and $*$ -property in terms of dom :

Definition. The *simple security property* says that a subject s can read an object o if and only if $level(s) dom level(o)$.

Definition. The $*$ -*property* says that a subject s can write to an object o if and only if $level(o) dom level(s)$.

In our example, assume the discretionary access controls are set to allow any subject all types of access. In that case, as $level(Erin) dom level(\text{EurDoc})$, Erin can read EurDoc (by the simple security property) but not write EurDoc (by the $*$ -property). Conversely, as $level(\text{EurAsiaDoc}) dom level(Erin)$, Erin cannot read EurAsiaDoc (by the simple security property) but can write to EurAsiaDoc (by the $*$ -property).

A logical question is how to determine the highest security compartment that both Erin and Don can read and the lowest that both can write. In order to do this, we must review some properties of dom .

First, note that $level(s) dom level(s)$; that is, dom is reflexive. The relation is also antisymmetric, because if both $level(s) dom level(o)$ and $level(o) dom level(s)$ are true, then $level(s) = level(o)$. It is transitive, because if $level(s_1) dom level(o)$ and $level(o) dom level(s_2)$, then $level(s_1) dom level(s_2)$.

We also define the *greatest lower bound (glb)* of two compartments as:

Definition. Let $A = (L, C)$ and $B = (L', C')$. Then $glb(A, B) = (\min(L, L'), C \cap C')$.

This answers the question of the highest security compartment that two subjects s and s' can read an object in. It is $glb(level(s), level(s'))$. For example, Don and Erin can both read objects in:

$$glb(level(\text{Don}), level(Erin)) = (\text{SECRET}, \emptyset).$$

This makes sense because Don cannot read an object in any compartment except those with the category set {ASIA} or the empty set, and Erin can only read objects in a compartment with the category set {EUR} or the empty set. Both are at the SECRET level, so the compartment must also be at the SECRET level.

We can define the *least upper bound (lub)* of two compartments analogously:

Definition. Let $A = (L, C)$ and $B = (L', C')$. Then $lub(A, B) = (\max(L, L'), C \cup C')$.

9 · 12 MATHEMATICAL MODELS OF COMPUTER SECURITY

We can now determine the lowest security compartment into which two subjects s and s' can write. It is $\text{lub}(\text{level}(s), \text{level}(s'))$. For example, Don and Erin can both write to objects in:

$$\text{glb}((\text{level}(\text{Don}), \text{level}(\text{Erin})) = (\text{SECRET}, \{\text{EUR}, \text{ASIA}\}).$$

This makes sense because Don cannot write to an object in any compartment except those with ASIA in the category set, and Erin can only write to objects in a compartment with EUR in the category set. The smallest category set meeting both these requirements is $\{\text{EUR}, \text{ASIA}\}$. Both are at the SECRET level, so the compartment must also be at the SECRET level.

The five properties of dom (reflexive, antisymmetric, transitive, existence of a least upper bound for every pair of elements, and existence of a greatest lower bound for every pair of elements) mean that the security compartments form a mathematical structure called a *lattice*. This has useful theoretical properties, and is important enough so models exhibiting this type of structure are called *lattice models*.

When the model is implemented on a system, the developers often make some modifications. By far the most common one is to restrict writing to the current compartment or to within a limited set of compartments. This prevents confidential information from being altered by those who cannot read it. The structure of the model can also be used to implement protections against malicious programs that alter files, such as system binaries. To prevent this, place the system binaries in a compartment that is dominated by those compartments assigned to users. By the simple security property, then users can read the system binaries, but by the $*$ -property, users cannot write them. Hence, if a computer virus infects a user's programs or documents,²² it can spread within that user's compartment but not to system binaries.

The Bell-LaPadula model is the basis for several other models. We explore one of its variants that models integrity rather than confidentiality.

9.4.2 Biba's Strict Integrity Policy Model. Biba's strict integrity policy model,²³ usually called *Biba's model*, is the mathematical dual of the Bell-LaPadula model.

Consider the issue of trustworthiness. When a highly trustworthy process reads data from an untrusted file and acts based on that data, the process is no longer trustworthy—as the saying goes, “garbage in, garbage out.” But if a process reads data more trustworthy than the process, the trustworthiness of that process does not change. In essence, the trustworthiness of the result is as trustworthy as the least trustworthy of the process and the data.

Define a set of *integrity classes* in the same way that we defined security compartments for the Bell-LaPadula model, and let $i\text{-level}(s)$ be the integrity compartment of s . Then the preceding text says that “reads down” (a trustworthy process reading untrustworthy data) should be banned, because it reduces the trustworthiness of the process. But “reads up” is allowed, because it does not affect the trustworthiness of the process. This is exactly the opposite of the simple security property.

Definition. The *simple integrity property* says that a subject s can read an object o if and only if $i\text{-level}(o) \text{ dom } i\text{-level}(s)$.

This definition captures the notion of allowing “reads up” and disallowing “reads down.”

Similarly, if a trustworthy process writes data to an untrustworthy file, the trustworthiness of the file may (or may not) increase. But if an untrustworthy process writes

CLASSIC MODELS 9 · 13

data to a trustworthy file, the trustworthiness of that file drops. s “writes down” should be allowed and “writes up” forbidden.

Definition. The $*$ -integrity property says that a subject s can write to an object o if and only if $i\text{-level}(s) \text{ dom } i\text{-level}(o)$.

This property blocks attempts to “write up” while allowing “writes down.”

A third property relates to execution of subprocesses. Suppose process *date* wants to execute the command *time* as a subprocess. If the integrity compartment of *date* dominates that of *time*, then any information *date* passes to *time* is passed to a less trustworthy process, and hence is allowed under the $*$ -integrity property. But if the integrity compartment of *time* dominates that of *date*, then the $*$ -integrity property is violated. Hence:

Definition. The *execution integrity property* says that a subject s can execute a subject s' if and only if $i\text{-level}(s') \text{ dom } i\text{-level}(s)$.

Given these three properties, one can show:

Theorem. If information can be transferred from object o_1 to object o_n , then by the simple integrity property, the $*$ -integrity property, and the execution integrity property, $i\text{-level}(o_1) \text{ dom } i\text{-level}(o_n)$.

In other words, if all the rules of Biba’s model are followed, the integrity of information cannot be corrupted because information can never flow from a less trustworthy object to a more trustworthy object.

This model suggests a method for analyzing programs to prevent security breaches. When the program runs, it reads data from a variety of sources: itself, the system, the network, and the user. Some of these sources are trustworthy, such as the process itself and the system. The user and the network are under the control of ordinary users (or remote users) and so are less trustworthy. So, apply Biba’s model with two integrity compartments, (UNTAINTED, \emptyset) (this means the set of categories in the compartment is empty) and (TAINTED, \emptyset), where (UNTAINTED, \emptyset) dom (TAINTED, \emptyset). For notational convenience, we shall write (UNTAINTED, \emptyset) as UNTAINTED and (TAINTED, \emptyset) as TAINTED; and dom as \geq . Thus, UNTAINTED \geq TAINTED.

The technique works with either static or dynamic analysis but is usually used for dynamic analysis. In this mode, all constants are assigned the integrity label UNTAINTED. Variables are assigned labels based on the data flows within the program. For example, in an assignment, the integrity label of the variable being assigned to is set to the integrity label of the expression assigned to it. When UNTAINTED and TAINTED variables are mixed in the expression, the integrity label of the expression is TAINTED. If a variable is assigned a value from an untrusted source, the integrity label of the variable is set to TAINTED.

When data are used as (for example) parameters of system calls or library functions, the system checks that the integrity label of the variable dominates that of the parameter. If it does not, the program takes some action, such as aborting, or logging a warning, or throwing an exception. This action either prevents an exploit or alerts the administrator of the attack.

For example, suppose a programmer wishes to prevent a format string attack. This is an attack that exploits a vulnerability in the C printing function *printf*. The first argument to *printf* is a format string, and the contents of that string determine how many other arguments *printf* expects. By manipulating the contents of a format string, an attacker can overwrite values of variables and corrupt the stack, causing the program to malfunction—usually to the attacker’s benefit. The key step of the attack is to input an unexpected value for the format string. Here is a code fragment with the flaw:

```
if (fgets(buf, sizeof(buf), stdin) != NULL) printf(buf);
```

9 · 14 MATHEMATICAL MODELS OF COMPUTER SECURITY

This reads a line of characters from the input into an array *buf* and immediately prints the contents of the array. If the input is “xyzzy%n”, then some element of the stack will be overwritten by the value 5.²⁴ Hence, the first parameter to *printf* must always have integrity class UNTAINTED.

Under this analysis technique, when the input function *fgets* is executed, the variable *buf* would be assigned an integrity label of TAINTED, because the input (which is untrusted) is stored in it. Then, at the call to *printf*, the integrity class of *buf* is compared to that required for the first parameter of *printf*. The former is TAINTED; the latter is UNTAINTED. But we require that the variable’s integrity class (TAINTED) dominate that of the parameter (UNTAINTED), and here TAINTED \leq UNTAINTED. Hence, the analysis has found a disallowed flow and acts accordingly.

9.4.3 Clark-Wilson Model. Lipner²⁵ identified five requirements for commercial integrity models:

1. Users may not write their own programs to manipulate trusted data. Instead, they must use programs authorized to access that data.
2. Programmers develop and test programs on nonproduction systems, using non-production copies of production data if necessary.
3. Moving a program from nonproduction systems to production systems requires a special process.
4. That special process must be controlled and audited.
5. Managers and system auditors must have access to system logs and the system’s current state.

Biba’s model can be instantiated to meet the first and last conditions by appropriate assignment of integrity levels, but the other three focus on integrity of processes. Hence, while Biba’s model works well for some problems of integrity, it does not satisfy these requirements for a commercial integrity model.

The Clark-Wilson model²⁶ was developed to describe processes within many commercial firms. There are several specialized terms and concepts needed to understand the Clark-Wilson mode; these are best introduced using an example:

- Consider a bank. If *D* are the day’s deposits, *W* the day’s withdrawals, *I* the amount of money in bank accounts at the beginning of the day, and *F* the amount of money in bank accounts at the end of the day, those values must satisfy the constraint $I + D - W = F$.
- This is called an *integrity constraint* because, if the system (the set of bank accounts) does not satisfy it, the bank’s integrity has been violated.
- If the system does satisfy its integrity constraints, it is said to be in a *consistent* state.
- When in operation, the system moves from one consistent state to another. The operations that do this are called *well-formed transactions*. For example, if a customer transfers money from one account to another, the transfer is the well-formed transaction. Its component actions (withdrawal from the first account and deposit in the second) individually are not well-formed transactions, because if only one completes, the system will be in an inconsistent state.

CLASSIC MODELS 9 · 15

- Procedures that verify that all integrity constraints are satisfied are called *integrity verification procedures* (IVPs).
- Data that must satisfy integrity constraints are called *constrained data items* (CDIs), and when they satisfy the constraints are said to be in a *valid* state.
- All other data are called *unconstrained data items* (UDIs).
- In addition to integrity constraints on the data, the functions implementing the well-formed transactions themselves are constrained. They must be *certified* to be well formed and to be implemented correctly. Such a function is called a *transformation procedure* (TP).

The model provides nine rules, five of which relate to the certification of data and TPs and four of which describe how the implementation of the model must enforce the certifications.

The first rule captures the requirement that the system be in a consistent state:

Certification Rule 1. An IVP must ensure that the system is in a consistent state.

The relation *certified* associates some set of CDIs with a TP that transforms those CDIs from one valid state to a (possibly different) valid state. The second rule captures this.

Certification Rule 2. For some set of associated CDIs, a TP transforms those CDIs from a valid state to a (possibly different) valid state.

The first enforcement rule ensures that the system keeps track of the *certified* relation and prevents any TP from executing with a CDI not in its associated *certified* set:

Enforcement Rule 1. The system must maintain the *certified* relation, and ensure that only TPs certified to run on a CDI manipulates that CDI.

In a typical firm, the set of users who can use a TP is restricted. For example, in a bank, a teller cannot move millions of dollars from one bank to another; doing that requires a bank officer. The second enforcement rule ensures that only authorized users can run TPs on CDIs by defining a relation *allowed* that associates a user, a TP, and the set of CDIs that the TP can access on that user's behalf:

Enforcement Rule 2. The system must associate a user with each TP and set of CDIs. The TP may access those CDIs on behalf of the associated user. If a user is not associated with a particular TP and set of CDIs, then the TP cannot access those CDIs on behalf of that user.

This implies that the system can correctly identify users. The next rule enforces this:

Enforcement Rule 3. The system must authenticate each user attempting to execute a TP.

This ensures that the identity of a person trying to execute a TP is correctly bound to the corresponding user identity within the computer. The form of authentication is left up to the instantiation of the model, because differing environments suggest different authentication requirements. For example, a bank officer may use a biometric device and a password to authenticate herself to the computer that moves millions of dollars; a teller whose actions are restricted to smaller amounts of money may need only to supply a password.

Separation of duty, already discussed, is a key consideration in many commercial operations. The Clark-Wilson model captures it in the next rule:

Certification Rule 3. The *allowed* relation must meet the requirements imposed by separation of duty.

9 · 16 MATHEMATICAL MODELS OF COMPUTER SECURITY

A cardinal principle of commercial integrity is that the operations must be auditable. This requires logging of enough information to determine what the transaction did. The next rule captures this requirement:

Certification Rule 4. All TPs must append enough information to reconstruct the operation to an append-only CDI.

The append-only CDI is, of course, a log.

So far we have considered all inputs to TPs to be CDIs. Unfortunately, that is infeasible. In our bank example, the teller will enter account information and deposit and withdrawal figures; but those are not CDIs; the teller may mistype something. Before a TP can use that information, it must be vetted to ensure it will enable the TP to work correctly. The last certification rule captures this:

Certification Rule 5. A TP that takes a UDI as input must perform either a well-formed transaction or no transaction for any value of the UDI. Thus, it either rejects the UDI or transforms it into a CDI.

This also covers poorly crafted TPs; if the input can exploit vulnerabilities in the TP to cause it to act in unexpected ways, it cannot be certified under this rule.

Within the model lies a possible conflict. In the preceding rules, one user could certify a TP to operate on a CDI and then execute the TP on that CDI. The problem is that a malicious user may certify a TP that does not perform a well-formed transaction, causing the system to violate the integrity constraints. Clearly, an application of the principle of separation of duty would solve this problem, and indeed the last rule in the model does just that:

Enforcement Rule 4. Only the certifier of a TP may change the *certified* relation for that TP. Further, no certifier of a TP, or of any CDI associated with that TP, may execute the TP on the associated CDI.

This separates the ability to certify a TP from the ability to execute that TP and the ability to certify a CDI for a given TP from the ability to execute that TP on that CDI. This enforces the separation of duty requirement.

Now, revisit Lipner's requirements for commercial integrity models. The TPs correspond to Lipner's programs and the CDIs to the production data. To meet requirement 1, the Clark-Wilson certifiers need to be trusted, and ordinary users cannot certify either TPs or CDIs. Then Enforcement Rule 4 and Certification Rule 5 enforce this requirement. Requirement 2 is met by not certifying the development programs; as they are not TPs, they cannot be run on production data. The "special process" in requirement 3 is a TP. Certification Rule 4 describes a log; the special process in requirement 3 being a TP, it will append information to the log that can be audited. Further, the TP is by definition a controlled process, and Enforcement Rule 4 and Certification Rule 5 control its execution. Before the installation, the program being installed is a UDI; after it is installed, it is a CDI (and a TP). Thus, requirement 4 is satisfied. Finally, the Clark-Wilson model has a log that captures all aspects of what a TP does, and that is the log the managers and auditors will have access to. They also have access to the system state because they can run an IVP to check its integrity. Thus, Lipner's requirement 5 is met. So the Clark-Wilson model is indeed a satisfactory commercial integrity model.

This model is important for two reasons. First, it captures the way most commercial firms work, including applying separation of duty (something that Biba's model does not capture well). Second, it separates the notions of certification and enforcement. Enforcement typically can be done within the instantiation of the model. But the model cannot enforce how certification is done; it can only require that a certifier claim to have done it. This is true of all models, of course, but the Clark-Wilson model specifically states the assumptions it makes about certification.

CLASSIC MODELS 9 · 17

9.4.4 Chinese Wall Model. Sometimes called the *Brewer-Nash model*, the goal of the Chinese Wall model²⁷ is to prevent conflicts of interest. It does so by grouping objects that belong to the same company into *company data sets* and company data sets into *conflict-of-interest classes*. If two companies (represented by their associated company data sets) are in the same conflict-of-interest class, then a lawyer or stockbroker representing both would have a conflict of interest. The rules of the model ensure that a subject can read only one company data set in each conflict-of-interest class.

In general, objects are documents or resources that contain information that the company wishes to (or is required to) keep secret. There is, however, an exception. Companies release data publicly, in the form of annual reports; that information is carefully *sanitized* to remove all confidential content. To reflect business practice, the model must allow *all* subjects to see that data. The model therefore defines a conflict-of-interest class called the *sanitized class* that has one company data set holding *only* objects containing sanitized data.

Now consider a subject reading an object. If the subject has never read any object in the object's conflict-of-interest class, reading the object presents no conflict of interest. If the subject has read an object in the same company data set, then the only information that the subject has seen in that conflict-of-interest class is from the same company as the object it is trying to read, which is allowed. But if the subject has read an object in the same conflict-of-interest class but a *different* company data set, then were the new read request granted, the subject would have read information from two different companies for which there is a conflict of interest—exactly what the model is trying to prevent. So that is disallowed.

The next rule summarizes this:

Definition. The *CW-simple security property* says that a subject s can read an object o if and only if either:

1. s has not read any other object in o 's conflict-of-interest class; or
2. The only objects in o 's conflict-of-interest class that s has read are all in o 's company data set.

To see why this works, suppose all banks are in the same conflict-of-interest class. A stockbroker represents The Big Bank. She is approached to represent The Bigger Bank. If she agreed, she would need access to The Bigger Bank's information, specifically the objects in The Bigger Bank's company data set. But that would mean she could read objects from two company data sets in the same conflict-of-interest class, something the CW-simple security property forbids. The temporal element of the model is important; even if she resigned her representation of The Big Bank, she cannot represent The Bigger Bank because condition 2 of the CW-simple security property considers all objects she previously read. This makes sense, because she has had access to The Big Bank, and could unintentionally compromise the interests of her previous employer while representing The Bigger Bank.

The CW-simple security property implicitly says that s can read any sanitized object. To see this, note that if s has never read a sanitized object, condition 1 holds. If s has read a sanitized object, then condition 2 holds because all sanitized objects are in the same company data set.

Writing poses another problem. Suppose Barbara represents The Big Bank, and Percival works for The Bigger Bank. Both also represent The Biggest Toy Company, which—not being a financial institution—is in a different conflict-of-interest class from

9 · 18 MATHEMATICAL MODELS OF COMPUTER SECURITY

either bank. Thus, there is no conflict of interest in either Barbara’s or Percival’s representation of a bank and the toy company. But there is a path along which information can flow from Barbara to Percival, and vice versa, that enables a conflict of interest to occur. Barbara can read information from an object in The Big Bank’s company data set and write it to an object in The Biggest Toy Company’s company data set. Percival can read the information from the object in The Biggest Toy Company’s company data set, thereby effectively giving him access to The Big Bank’s information—which is a conflict of interest. That he needs Barbara’s help does not detract from the problem. The goal of the model requires that this conspiracy be prevented. The next rule does so:

Definition. The *CW^{*}-property* says that a subject s can write to an object o if and only if both of the following conditions are met:

1. The CW-simple security property allows s to read o ; and
2. All unsanitized objects that s can read are in the same company data set as o .

Now Barbara can read objects in both The Big Bank’s company data set and The Biggest Toy Company’s data set. But when she tries to write to The Biggest Toy Company’s data set, the CW^{*}-property prevents her from doing so as condition 2 is not met (because she can read an object in The Big Bank’s company data set).

This also accounts for sanitized objects. Suppose that Skyler represents The Biggest Toy Company and no other company. He can also read information from the sanitized class. When he tries to write to an object in The Biggest Toy Company’s company data set, he meets both conditions of the CW-simple security property (because he has only read objects in that company data set), and all unsanitized objects that he can read are in the same company data set as the object he can read. Thus, both conditions of the CW^{*}-property are met, so Skyler can write the object.

The conditions of the CW^{*}-property are very restrictive; effectively, a subject can write to an object only if it has access to the company data set containing that object, *and no other company data set except the company data set in the sanitized class*. But without this restriction, conflicts of interest are possible.

9.4.5 Summary. The four models discussed in this section have played critical roles in the development of our understanding of computer security. Although it is not the first model of confidentiality, the Bell-LaPadula model describes a widely used security scheme. The Biba model captured notions of “trust” and “trustworthiness” in an intuitive way, and recent advances in the analysis of programs for vulnerabilities have applied that model to great effect. The Clark-Wilson model moved the notion of commercial integrity models away from multilevel models to models that examine process integrity as well as data integrity. The Chinese Wall model explored conflict of interest, an area that often arises when one is performing confidential services for multiple companies or has access to confidential information from a number of companies. These models are considered classic because their structure and ideas underlie the rules and structures of many other models.

9.5 OTHER MODELS. Some models examine specific environments. The Clinical Information Systems Security model²⁸ considers the protection of health records, emphasizing accountability as well as confidentiality and integrity. *Traducement*²⁹

FURTHER READING 9 · 19

describes the process of real estate recordation, which requires a strict definition of integrity and accountability with little to no confidentiality.

Other models generalize the classic models. The best known are the models of *noninterference security* and *deducibility security*. Both are multilevel security models with two levels, HIGH and LOW. The noninterference model³⁰ defines security as the ability of a HIGH subject to interfere with what the LOW subject sees. For example, if a HIGH subject can prevent a LOW subject from acquiring a resource at a particular time, the HIGH subject can transmit information to the LOW subject. In essence, the interference is a form of writing, and must be prevented, just as the Bell-LaPadula model prevents a HIGH subject from writing to a LOW object. The deducibility model³¹ examines whether a LOW subject can infer anything about a HIGH subject's actions by examining only the LOW outputs. Both these models are useful in analyzing the security of systems³² and intrusion detection mechanisms,³³ and led to work that showed connecting two secure compute systems may produce a nonsecure system.³⁴ Further work is focusing on establishing conditions under which connecting two secure systems produces a secure system.³⁵

9.6 CONCLUSION. The efficacy of mathematical modeling depends on the application of those models. Typically, the models capture system-specific details and describe constraints ensuring the security of the system or the information on the system. If the model does not correctly capture the details of the *entire* system, the results may not be comprehensive, and the analysis may miss ways in which security could be compromised.

This is an important point. For example, the Bell-LaPadula model captures a notion of what the *system* must do to prevent a subject cleared for TOP SECRET leaking information to a subject cleared for CONFIDENTIAL. But if the system enforces that model, the TOP SECRET subject could still meet the CONFIDENTIAL subject and hand her a printed version of the TOP SECRET information. That is outside the system and so was not captured by the model. But if the model also embraces procedures, then a procedure is necessary to prevent this “writing down.” In that case, the flaw would be in the implementation of the procedure that failed to prevent the transfer of information—in other words, an incorrect instantiation of the model, exactly what Dorothy Denning’s comment in the introduction to this section referred to.

The models described in this section span the foundational (access-control matrix model) to the applied (Bell-LaPadula, Biba, Clark-Wilson, and Chinese Wall). All play a role in deepening our understanding of what security is and how to enforce it.

The area of mathematical modeling is a rich and important area. It provides a basis for demonstrating that the design of systems is secure, for specific definitions of *secure*. Without these models, our understanding of how to secure systems would be diminished.

9.7 FURTHER READING

- Anderson, R. “A Security Policy Model for Clinical Information Systems,” Proceedings of the 1996 IEEE Symposium on Security and Privacy (May 1996): 34–48.
- Bell, D., and LaPadula, L. “Secure Computer Systems: Unified Exposition and Multics Interpretation,” Technical Report MTR-2997 rev. 1, MITRE Corporation, Bedford, MA (March 1975).
- Biba, K. “Integrity Considerations for Secure Computer Systems,” Technical Report MTR-3153, MITRE Corporation, Bedford, MA (April 1977).

9 · 20 MATHEMATICAL MODELS OF COMPUTER SECURITY

- Bishop, M. *Computer Security: Art and Science*. Boston: Addison-Wesley Professional, 2002.
- Brewer, D., and M. Nash. "The Chinese Wall Security Policy," *Proceedings of the 1989 IEEE Symposium on Security and Privacy* (May 1989): 206–212.
- Clark, D., and D. Wilson. "A Comparison of Commercial and Military Security Policies," *Proceedings of the 1987 IEEE Symposium on Security and Privacy* (April 1987): 184–194.
- Cortier, V., and S. Kremer. *Formal Models and Techniques for Analyzing Security Protocols*—Volume 5 *Cryptology and Information Security Series*. IOS Press, 2011.
- Cremers, C., and S. Mauw. *Operational Semantics and Verification of Security Protocols*. Springer, 2012.
- Demillo, D., D. Dobkin, A. Jones, and R. Lipton, eds. *Foundations of Secure Computing*. New York: Academic Press, 1978.
- Denning, D. "The Limits of Formal Security Models," National Information Systems Security Conference, October 18, 1999; available at www.cs.georgetown.edu/~denning/infosec/award.html
- Denning, P. "Third Generation Computer Systems," *Computing Surveys* 3, No. 4 (December 1976): 175–216.
- Engeler, E. *Introduction to the Theory of Computation*. New York: Academic Press, 1973.
- Ferraiolo, D., J. Cugini, and D. Kuhn. "Role-Based Access Control (RBAC): Features and Motivations," *Proceedings of the Eleventh Annual Computer Security Applications Conference* (December 1995): 241–248.
- Gougen, J., and J. Meseguer. "Security Policies and Security Models," *Proceedings of the 1982 IEEE Symposium on Privacy and Security* (April 1982): 11–20.
- Graubert, R. "On the Need for a Third Form of Access Control," *Proceedings of the Twelfth National Computer Security Conference* (October 1989): 296–304.
- Haigh, J., R. Kemmerer, J. McHugh, and W. Young. "An Experience Using Two Covert Channel Analysis Techniques on a Real System Design," *IEEE Transactions in Software Engineering* 13, No. 2 (February 1987): 141–150.
- Harrison, M., and W. Ruzzo, "Monotonic Protection Systems." In D. Demillo et al., eds. *Foundations of Secure Computing*, pp. 337–363. New York: Academic Press, 1978.
- Harrison, M., W. Ruzzo, and J. Ullman. "Protection in Operating Systems," *Communications of the ACM* 19, No. 8 (August 1976): 461–471.
- Ko, C., and T. Redmond. "Noninterference and Intrusion Detection," *Proceedings of the 2002 IEEE Symposium on Security and Privacy* (May 2002): 177–187.
- Lampson, B. "Protection." *Proceedings of the Fifth Princeton Symposium of Information Science and Systems* (March 1971): 437–443.
- Lipner, S. "Non-Discretionary Controls for Commercial Applications," *Proceedings of the 1982 IEEE Symposium on Privacy and Security* (April 1982): 2–10.
- Mantel, H. "On the Composition of Secure Systems," *Proceedings of the 2002 IEEE Symposium on Security and Privacy* (May 2002): 88–101.
- McCullough, D. "Non-Interference and the Composability of Security Properties," *Proceedings of the 1987 IEEE Symposium on Privacy and Security* (April 1988): 177–186.
- Sandhu, R. "The Typed Access Matrix Model," *Proceedings of the 1992 IEEE Symposium on Security and Privacy* (April 1992): 122–136.

NOTES 9 · 21

- Seacord, R. *Secure Coding in C and C++*. Boston: Addison-Wesley, 2005.
- Walcott, T., and M. Bishop. "Traducement: A Model for Record Security," *ACM Transactions on Information Systems Security* 7, No. 4 (November 2004): 576–590.

9.8 NOTES

1. Recordation of real estate refers to recording deeds, mortgages, and other information about property with the county recorder. See <http://ag.ca.gov/erds1/index.php>
2. D. Denning, "The Limits of Formal Security Models," National Information Systems Security Conference, October 18, 1999; available at www.cs.georgetown.edu/~denning/infosec/award.html
3. B. Lampson, "Protection," *Proceedings of the Fifth Princeton Symposium of Information Science and Systems* (March 1971): 437–443; P. Denning, "Third Generation Computer Systems," *Computing Surveys* 3, No. 4 (December 1976): 175–216.
4. Harrison, 1976.
5. M. Harrison, W. Ruzzo, and J. Ullman, "Protection in Operating Systems," *Communications of the ACM* 19, No. 8 (August 1976): 461–471.
6. The halting problem is the question "Is there an algorithm to determine whether any arbitrary program halts?" The answer, "No," was proved by Alan Turing in 1936. See www.nist.gov/dads/HTML/haltingProblem.html. See also E. Engeler, *Introduction to the Theory of Computation* (New York: Academic Press, 1973).
7. E. Engeler, *Introduction to the Theory of Computation* (New York: Academic Press, 1973).
8. The interested reader is referred to Harrison et al., "Protection in Operating Systems," or to M. Bishop, *Computer Security: Art and Science* (Boston: Addison-Wesley Professional, 2002), p. 47 ff., for the proof.
9. Harrison et al., "Protection in Operating Systems."
10. M. Harrison and W. Ruzzo, "Monotonic Protection Systems," in D. Demillo et al., eds., *Foundations of Secure Computing*, pp. 337–363 (New York: Academic Press, 1978).
11. Harrison and Ruzzo, "Monotonic Protection Systems."
12. Harrison and Ruzzo, "Monotonic Protection Systems."
13. Harrison and Ruzzo, "Monotonic Protection Systems."
14. R. Sandhu, "The Typed Access Matrix Model," *Proceedings of the 1992 IEEE Symposium on Security and Privacy* (April 1992): 122–136.
15. Sandhu, "The Typed Access Matrix Model."
16. Also sometimes called *organization-controlled access control*, or ORGCON.
17. R. Graubert, "On the Need for a Third Form of Access Control," *Proceedings of the Twelfth National Computer Security Conference* (October 1989): 296–304.
18. D. Ferraiolo, J. Cugini, and D. Kuhn, "Role-Based Access Control (RBAC): Features and Motivations," *Proceedings of the Eleventh Annual Computer Security Applications Conference* (December 1995): 241–248.
19. As Benjamin Franklin once said, "Three can keep a secret if two of them are dead."
20. Saltzer and Schroeder 1975.
21. D. Bell and L. LaPadula, "Secure Computer Systems: Unified Exposition and Multics Interpretation," Technical Report MTR-2997 rev. 1, MITRE Corporation, Bedford, MA (March 1975).

9 · 22 MATHEMATICAL MODELS OF COMPUTER SECURITY

22. A macro virus can infect a document. See, for example, Bishop, M., *Computer Security: Art and Science*. (Boston, MA: Addison-Wesley Professional, 2002), section 22.3.8, and Chapter 16 in this *Handbook*.
23. Biba proposed three models: the Low-Water-Mark Policy model, the Ring Policy model, and the Strict Integrity Policy model. See Bishop, *Computer Security*, Section 6.2.
24. The explanation is too complex to go into here. The interested reader is referred to R. Seacord, *Secure Coding in C and C++* (Boston: Addison-Wesley, 2005), Chapter 6, for a discussion of this problem.
25. S. Lipner, “Non-Discretionary Controls for Commercial Applications,” *Proceedings of the 1982 IEEE Symposium on Privacy and Security* (April 1982): 2–10.
26. D. Clark and D. Wilson, “A Comparison of Commercial and Military Security Policies,” *Proceedings of the 1987 IEEE Symposium on Security and Privacy* (April 1987): 184–194.
27. D. Brewer and M. Nash, “The Chinese Wall Security Policy,” *Proceedings of the 1989 IEEE Symposium on Security and Privacy* (May 1989): 206–212.
28. R. Anderson, “A Security Policy Model for Clinical Information Systems,” *Proceedings of the 1996 IEEE Symposium on Security and Privacy* (May 1996): 34–48.
29. T. Walcott and M. Bishop, “Traducement: A Model for Record Security,” *ACM Transactions on Information Systems Security* 7, No. 4 (November 2004): 576–590.
30. J. Gougen and J. Meseguer, “Security Policies and Security Models,” *Proceedings of the 1982 IEEE Symposium on Privacy and Security* (April 1982): 11–20.
31. Gougen and Meseguer, “Security Policies and Security Models.”
32. J. Haigh, R. Kemmerer, J. McHugh, and W. Young, “An Experience Using Two Covert Channel Analysis Techniques on a Real System Design,” *IEEE Transactions in Software Engineering* 13, No. 2 (February 1987): 141–150.
33. C. Ko and T. Redmond, “Noninterference and Intrusion Detection.” *Proceedings of the 2002 IEEE Symposium on Security and Privacy* (May 2002): 177–187.
34. D. McCullough, “Non-Interference and the Composability of Security Properties,” *Proceedings of the 1987 IEEE Symposium on Privacy and Security* (April 1988): 177–186.
35. H. Mantel, “On the Composition of Secure Systems,” *Proceedings of the 2002 IEEE Symposium on Security and Privacy* (May 2002): 88–101.

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 10

UNDERSTANDING STUDIES AND SURVEYS OF COMPUTER CRIME

M. E. Kabay

10.1 INTRODUCTION	10·1	10.2.1 Some Fundamentals of Statistical Design and Analysis	10·3
10.1.1 Value of Statistical Knowledge Base	10·1	10.2.2 Research Methods Applicable to Computer Crime	10·9
10.1.2 Limitations on Our Knowledge of Computer Crime	10·2		
10.1.3 Limitations on the Applicability of Computer Crime Statistics	10·2	10.3 SUMMARY	10·11
		10.4 FURTHER READING	10·12
10.2 BASIC RESEARCH METHODOLOGY	10·3	10.5 NOTES	10·12

10.1 INTRODUCTION. This chapter provides guidance for critical reading of research results about computer crime. It will also alert designers of research instruments who may lack formal training in survey design and analysis to the need for professional support in developing questionnaires and analyzing results.

10.1.1 Value of Statistical Knowledge Base. Security specialists are often asked about computer crime; for example, customers want to know who is attacking which systems, how often, using what methods. These questions are perceived as important because they bear on the strategies of risk management; in theory, in order to estimate the appropriate level of investment in security, it would be helpful to have a sound grasp of the probability of different levels of damage. Ideally, one would want to evaluate an organization's level of risk by evaluating the experiences of other organizations with similar system and business characteristics. Such comparisons would be useful in competitive analysis and in litigation over standards of due care and diligence in protecting corporate assets.

10 · 2 UNDERSTANDING STUDIES AND SURVEYS OF COMPUTER CRIME

10.1.2 Limitations on Our Knowledge of Computer Crime. Unfortunately, in the current state of information security, no one can give reliable answers to such questions. There are two fundamental difficulties preventing us from developing accurate statistics of this kind. These difficulties are known as the problems of ascertainment.

10.1.2.1 Detection. The first problem is that an unknown number of crimes of all kinds are undetected. For example, even outside the computer crime field, we do not know how many financial frauds are being perpetrated. We do not know because some of them are not detected. How do we know they are not detected? Because some frauds are discovered long after they have occurred. Similarly, computer crimes may not be detected by their victims but may be reported by the perpetrators.

In a landmark series of tests at the Department of Defense (DoD), the Defense Information Systems Agency found that very few of the penetrations it engineered against unclassified systems within the DoD seem to have been detected by system managers. These studies were carried out from 1994 through 1996 and attacked 38,000 systems. About two-thirds of the attacks succeeded; however, only 4 percent of these attacks were detected.¹

A commonly held view within the information security community is that only one-tenth or so of all the crimes committed against and using computer systems are detected.

10.1.2.2 Reporting. The second problem of ascertainment is that even if attacks are detected, few seem to be reported in a way that allows systematic data collection. This commonly held belief is based in part on the unquantified experience of information security professionals who have conducted interviews of their clients; it turns out that only about 10 percent of the attacks against computer systems revealed in such interviews were ever reported to any kind of authority or to the public. The DoD studies mentioned earlier were consistent with this belief; of the few penetrations detected, only a fraction of 1 percent were reported to appropriate authorities.

Given these problems of ascertainment, computer crime statistics generally should be treated with skepticism.

10.1.3 Limitations on the Applicability of Computer Crime Statistics. Generalizations in this field are difficult to justify. Even if we knew more about types of criminals and the methods they use, it still would be difficult to have the kind of actuarial statistic that is commonplace in the insurance field. For example, the establishment of uniform building codes in the 1930s in the United States led to the growth in fire insurance as a viable business. With official records of fires in buildings that could be described using a standard typology, statistical information began to provide an actuarial basis for using probabilities of fires and associated costs to calculate reasonable insurance rates.

In contrast, even if we had access to accurate reports, it would be difficult to make meaningful generalizations about vulnerabilities and incidence of successful attacks for the information technology field. We use a bewildering variety and versions of processors, operating systems, firewalls, encryption, application software, backup methods and media, communications channels, identification, authentication, authorization, compartmentalization, and operations.

BASIC RESEARCH METHODOLOGY 10 · 3

How would we generalize from data about the risks at (say) a mainframe-based network running Multiple Virtual Systems (MVS) in a military installation to the kinds of risks faced by a UNIX-based intranet in an industrial corporation, or to a Windows New Technology (NT)-based Web server in a university setting? There are so many differences among systems that if we were to establish a multidimensional analytical table where every variable was an axis, many cells would likely contain no or only a few examples. Such sparse matrices are notoriously difficult to use in building statistical models for predictive purposes.

10.2 BASIC RESEARCH METHODOLOGY. This is not a chapter about social sciences research. However, many discussions of computer crime seem to take published reports as gospel, even though these studies may have no validity whatsoever. In this short section, we look at some fundamentals of research design so that readers will be able to judge how much faith to put in computer crime research results.

10.2.1 Some Fundamentals of Statistical Design and Analysis. The way in which a scientist or reporter represents data can make an enormous difference in the readers' impressions.

10.2.1.1 Descriptive Statistics. Suppose three companies reported these losses from penetration of their computer systems: \$1 million, \$2 million, and \$6 million. We can describe these results in many ways. For example, we can simply list the raw data; however, such lists could become unacceptably long as the number of reports increased, and it is hard to make sense of the raw data.

We could define classes such as “2 million or less” and “more than 2 million” and count how many occurrences there were in each class:

Class	Freq
$\leq \$2M$	2
$> \$2M$	1

Alternatively, we might define the classes with finer granularity as $< \$1M$, $\geq \$1M$ but $< \$2M$, and so on; such a table might look like this:

Class	Freq
$< \$1M$	0
$\geq \$1M \text{ & } < \$2M$	1
$\geq \$2M \text{ & } < \$3M$	1
$\geq \$3M \text{ & } < \$4M$	0
$\geq \$4M \text{ & } < \$5M$	0
$\geq \$5M \text{ & } < \$6M$	0
$\geq \$6$	1

10 · 4 UNDERSTANDING STUDIES AND SURVEYS OF COMPUTER CRIME

Notice how the definition of the classes affects perception of the results: The first table gives the impression that the results are clustered around \$2 million and gives no information about the upper or lower bounds.

10.2.1.1.1 Location. One of the most obvious ways we describe data is to say where they lie in a particular dimension. The *central tendency* of our three original data (\$1 million, \$2 million, and \$6 million) can be represented in various ways; for example, two popular measures are:

- Arithmetic mean or average = $(1+2+6)/3 = \$3M$
- Median (the middle of the sorted list of losses) = \$2M

Note that if we tried to compute the mean and the median from the first table (with its approximate classes), we would get the wrong value. Such statistics should be computed from the original data, not from summary tables.

10.2.1.1.2 Dispersion. Another aspect of our data that we frequently need is *dispersion*—that is, variability. The simplest measure of dispersion is the range: the difference between the smallest and the largest value we found; in our example, we could say that the range was from \$1 million to \$6 million or that it was \$5 million. Sometimes the range is expressed as a percentage of the mean; then we would say that the range was $5/3 = 1.6 \dots$ or ~ 167 percent.

The *variance* (σ^2) of these particular data is the average of the squared deviations from the arithmetic mean; the variance of the three numbers would be $\sigma^2 = (1-3)^2 + (2-3)^2 + (6-3)^2 / 3 = (4+1+9)/3 \approx 4.67$.

The square root of the variance (σ) is called the *standard deviation* and is often used to describe dispersion. In our example, $\sigma = \sqrt{4.67} \approx 2.16$.

Dispersion is particularly important when we compare estimates about information from different groups. The greater the variance of a measure, the more difficult it is to form reliable generalizations about an underlying phenomenon, as described in the next section.

10.2.1.2 Inference: Sample Statistics versus Population Statistics.

We can accurately describe any data using descriptive statistics; the question is what we then do with those measures.

Usually we expect to extend the findings in a *sample* or subset of a *population* to make generalizations about the population. For example, we might be trying to estimate the losses from computer crime in commercial organizations with offices in the United States and with more than 30,000 employees. Or perhaps our sample would represent commercial organizations with offices in the United States and with more than 30,000 employees and whose network security staff was willing to respond to a survey questionnaire.

In such cases, we try to infer the characteristics of the population from the characteristics of the sample. Statisticians say that we try to estimate the *parametric* statistics by using the *sample* statistics.

For example, we estimate the parametric (population) variance (usually designated σ^2) by multiplying the variance of the sample by $n/(n-1)$. Thus, we would say that the estimate of the parametric variance (s^2) in our sample would be $s^2 = 4.67 * 3/2 = 7$. The estimate of the parametric standard deviation (s) would be $s = \sqrt{7} \approx 2.65$.

BASIC RESEARCH METHODOLOGY 10 · 5

10.2.1.3 Hypothesis Testing. Another kind of inference that we try to make from data is *hypothesis testing*. For example, suppose we were interested in whether there was any association between the presence or absence of firewalls and the occurrence of system penetration. We can imagine collecting these data about penetrations into systems with or without firewalls:

Firewalls	Penetration		
	No	Yes	Totals
No	25	75	100
Yes	70	130	200
Totals	95	205	300

We would frame the hypothesis (the *null hypothesis*, sometimes represented as H_0) that there was *no* relationship between the two independent variables, penetration and firewalls, and *test* that hypothesis by performing a test of independence of these variables. In our example, a simple chi-square test of independence would give a *test statistic* of $\chi^2_{[1]} = 2.636$. If there really were no association between penetration and firewalls in the population of systems under examination, the parametric value of this statistic would be zero. In our imaginary example, we can show that such a large value (or larger) of $\chi^2_{[1]}$ would occur in only 10.4 percent of the samples taken from a population where firewalls had no effect on penetration. Put another way, if we took many samples from a population where the presence of firewalls was not associated with any change in the rate of penetration, we would see about 10.4 percent of those samples producing $\chi^2_{[1]}$ statistics as large as or larger than 2.636.

Statisticians have agreed on some conventions for deciding whether a test statistic deviates enough from the value expected under the null hypothesis to warrant inferring that the null hypothesis is wrong. Generally, we describe the likelihood that the null hypothesis is true—often shown as $p(H_0)$ —in this way:

- When $p(H_0) > 0.05$, we say the results are *not statistically significant* (often designated with the symbols *ns*);
- When $0.05 \geq p(H_0) > 0.01$, the results are described as statistically significant (often designated with the symbol *);
- When $0.01 \geq p(H_0) > 0.001$, the results are described as highly statistically significant (often designated with the symbols **);
- When $p(H_0) \leq 0.001$, the results are described as extremely statistically significant (often designated with the symbols ***).

10.2.1.4 Random Sampling, Bias, and Confounded Variables. The most important element of sampling is randomness. We say that a sample is *random* or *randomized* when every member of the population we are studying has an equal probability of being selected. When a population is defined one way but the sample is drawn nonrandomly, the sample is described as *biased*. For example, if the population we are studying was designed to be, say, all companies worldwide with more than 30,000 full-time employees, but we sampled mostly from such companies in the United States, the sample would be biased toward U.S. companies and their characteristics.

10 · 6 UNDERSTANDING STUDIES AND SURVEYS OF COMPUTER CRIME

Similarly, if we were supposed to be studying security in all companies in the United States with more than 30,000 full-time employees, but we sampled only from those companies that were willing to respond to a security survey, we would be at risk of having a biased sample.

In this last example, involving studying only those who respond to a survey, we say that we are potentially *confounding* variables: We are looking at people-who-respond-to-surveys and hoping they are representative of the larger population of people from all companies in the desired population. But what if the people who are *willing* to respond are those who have better security and those who do not respond have terrible security? Then *responding to the survey* is confounded with *quality of security*, and our biased sample could easily mislead us into overestimating the level of security in the desired population.

Another example of how variables can be confounded is comparisons of results from surveys carried out in different years. Unless exactly the same people are interviewed in both years, we may be confounding individual variations in responses with changes over time; unless exactly the same companies are represented, we may be confounding differences among companies with changes over time; if external events have led people to be more or less willing to respond truthfully to questions, we may be confounding willingness to respond with changes over time. If the surveys are carried out with different questions or used by different research groups, we may be confounding changes in methodology with changes over time.

10.2.1.5 Confidence Limits. Because random samples naturally vary around the parametric (population) statistics, it is not very helpful to report a *point estimate* of the parametric value. For example, if we read that the mean damage from computer crimes in a survey was \$180,000 per incident, what does that imply about the population mean?

To express our confidence in the sample statistic, we calculate the likelihood of being right if we give an *interval estimate* of the population value. For example, we might find that we would have a 95 percent likelihood of being right in asserting that the mean damage was between \$160,000 and \$200,000. In another sample, we might be able to narrow these *95 percent confidence limits* to \$175,000 and \$185,000.

In general, the larger the sample size, the narrower the confidence limits will be for particular statistics.

The calculation of confidence limits for statistics depends on some necessary *assumptions*:

- Random sampling
- A *known error distribution* (usually the *Normal* distribution—sometimes called a *Gaussian* distribution)
- *Equal variance* at all values of the measurements

If any of these assumptions is wrong, the calculated confidence limits for our estimates will be wrong; that is, they will be misleading. There are tests of these assumptions that analysts should carry out before reporting results; if the data do not follow Normal error distributions, sometimes one can apply *normalizing transformations*.

In particular, percentages do not follow a Normal distribution. Here is a reference table of confidence limits for various percentages in a few representative sample sizes.

BASIC RESEARCH METHODOLOGY 10 · 7**95 Percent Confidence Limits for Percentages**

Percentage	Sample size		
	100	500	1000
0	0–3.0%	0–0.6%	0–0.3%
10	4.9–17.6%	7.5–13.0%	8.2–12.0%
20	12.7–29.1%	16.6–23.8%	17.6–22.6%
50	40.0–60.1%	45.5–54.5%	46.9–53.1%
80	70.9–87.3%	76.2–83.4%	77.4–82.4%
90	82.4–95.1%	87.0–92.5%	88.0–91.8%
100	97.0–100%	99.4–100%	99.7–100%

10.2.1.6 Contingency Tables. One of the most frequent errors in reporting results of studies is to provide only part of the story. For example, one can read statements such as “Over 70 percent of the systems without firewalls were penetrated last year.” Such a statement may be true, but it cannot be interpreted correctly as meaning that systems with firewalls were necessarily more or less vulnerable to penetration than systems without firewalls. The statement is incomplete; to make sense of it, we need the other part of the implied *contingency table*—the percentage of systems *with* firewalls that were penetrated last year—before making any assertions about the relationship between firewalls and penetrations. Compare, for example, these two hypothetical tables:

	Without Firewalls	With Firewalls in Default Configuration		Without Firewalls	With Firewalls Properly Configured
Penetrated	70%	70%	Penetrated	70%	10%
Not Penetrated	30%	30%	Not Penetrated	30%	90%

In both cases, someone could say that “70 percent of the systems without firewalls were penetrated,” but the implications would be radically different in the two data sets. Without knowing the right-hand column, the original assertion would be meaningless.

10.2.1.7 Association versus Causality. Continuing our example with rates of penetration, another error that untrained people often make when studying statistical information is to mistake *association* for *causality*. Imagine that a study showed that a lower percentage of systems with fire extinguishers was penetrated than systems without fire extinguishers and that this difference was statistically highly significant. Would such a result necessarily mean that fire extinguishers *caused* the reduction in penetration? No. We know that it is far more reasonable to suppose that the fire extinguishers were installed in organizations whose security awareness and security policies were more highly developed than in the organizations where no fire extinguishers were installed. In this imaginary example, the fire extinguishers might actually have *no causal effect* whatever on resistance to penetration. This result would illustrate the effect of

10 · 8 UNDERSTANDING STUDIES AND SURVEYS OF COMPUTER CRIME

confounding variables: *presence of a fire extinguisher* with *state of security awareness and policies*.

10.2.1.8 Control Groups. Finally, to finish our penetration example, one way to distinguish between association and causality is to *control* for variables. For example, one could measure the state of security awareness and policy as well as the presence or absence of fire extinguishers and make comparisons only among groups with the same level of awareness and policy. There are also statistical techniques for mathematically controlling for differences in such *independent variables*.

10.2.1.9 A Priori versus a Posteriori Testing. Amateurs or beginners sometimes forget the principle of *random sampling* that underlies all statistical inference (see Section 10.2.1.4). None of the hypothesis tests or confidence limit calculations work if a sample is not random. For example, if someone is wandering through a supermarket and notices that Granny Smith apples seem to be bigger than Macintosh apples, selecting a sample—even a random sample—of the apples that specifically gave rise to the hypothesis will not allow reliable computations of probability that the apples have the same average weight. The problem is that those *particular* apples would not have been sampled at all had the observer not been moved to formulate the hypothesis. So even if a particular statistical comparison produces a sample statistic that appears to have a probability of, say, 0.001, it is not possible to know how much the sampling deviated from randomness.

Applying statistical tests to data *after* one notices an interesting descriptive value, comparison, or trend is known as *a posteriori testing*. Formulating a hypothesis, obtaining a random sample, and computing the statistics and probabilities in accordance with the assumptions of those statistics and probabilities is known as *a priori testing*.

A well-used example of the perils of a posteriori testing is the unfortunate habit of searching through sequences of results such as long strings of guesses collected in student tests of paranormal abilities and calculating statistical values on carefully selected subsets of the strings. These a posteriori tests are then presented as if they were a priori and cause great confusion and arguments, such as: “Look, even though the overall proportion of correct guesses was (say) 50.003 percent in this run of <some very large number> guesses, there was a run of <much smaller number> guesses that were correct <any value greater than 50 percent> of the time! The probability of such a result by chance is <very small number>. That proves that there was a real effect of <whatever the treatment was>.” Unfortunately, a long series of numbers can produce any desired nonrandom-looking string; there are even tests known as *runs tests* that can help a researcher evaluate the nonrandomness of such occurrences.

In practical terms, statisticians have established a convention for limiting the damaging effects of a posteriori testing: Use the 0.001 level of probability as the equivalent of the minimum probability of the null hypothesis. This custom makes it far less likely that an a posteriori comparison will trick the user into accepting what is in fact a random variation that caught someone’s eye.

The best solution to the bias implicit in a posteriori testing is to use a completely *new* sample for the comparison. In the apple example, one could ask the store manager for new, unobserved, and randomly selected batches of both types of apples. The comparison statistics would then be credible and could be expected to follow the parametric distribution underlying calculations of probability of the null hypothesis. The populations from which these apples were selected still would have to be carefully

BASIC RESEARCH METHODOLOGY 10 · 9

determined. Would the populations be apples at this particular store? For this particular chain? For this particular region of the country or of the world?

10.2.2 Research Methods Applicable to Computer Crime

10.2.2.1 Interviews. Interviewing individuals can be illuminating. In general, interviews provide a wealth of data that are unavailable through any other method. For example, one can learn details of computer crime cases or motivations and techniques used by computer criminals. Interviews can be structured (using precise lists of questions) or unstructured (allowing the interviewer to respond to new information by asking additional questions at will).

Interviewers can take notes or record the interviews for later word-for-word transcription. In unstructured interviewers, skilled interviewers can probe responses to elucidate nuances of meaning that might be lost using cruder techniques such as surveys. Techniques such as thematic analysis can reveal patterns of responses that can then be examined using exploratory data analysis.² Thematic analysis is a technique for organizing nonquantitative information without imposing a preexisting framework on the data; exploratory data analysis uses statistical techniques to identify possibly interesting relationships that can be tested with independently acquired data. Such exploratory techniques can correctly include a posteriori testing as described in Section 10.2.1.9, but the results are used to propose further studies that can use a priori tests for the best use of resources.

10.2.2.2 Focus Groups. Focus groups are like group interviews. Generally, the facilitator uses a list of predetermined questions and encourages the participants to respond freely and to interact with each other. Often, the proceedings are filmed from behind a one-way mirror for later detailed analysis. Such analysis can include nonverbal communications, such as facial expressions and other body language as the participants speak or listen to others speak about specific topics.

10.2.2.3 Surveys. Surveys consist of asking people to answer a fixed series of questions with lists of allowable answers. They can be carried out face to face or by distributing and retrieving questionnaires by telephone, mail, fax, and email. Some questionnaires have been posted on the Web.

The critical issue when considering the reliability of surveys is *self-selection bias*—the obvious problem that survey results include only the responses of people who agreed to participate. Before basing critical decisions on survey data, it is useful to find out what the response rate was; although there are no absolutes, in general, we tend to trust survey results more when the response rate is high. Unfortunately, response rates for telephone surveys are often less than 10 percent; response rates for mail and email surveys can be less than 1 percent. It is very difficult to make any case for random sampling under such circumstances, and all results from such low-response-rate surveys should be viewed as indicating the range of problems or experiences of the respondents rather than as indicators of population statistics.

Regarding Web-based surveys, there are two types from a statistical point of view: those that use strong identification and authentication and those that do not. Those that do not are vulnerable to fraud, such as repeated voting by the same individuals. Those that provide individual universal resource locators (URLs) to limit voting to one per

10 · 10 UNDERSTANDING STUDIES AND SURVEYS OF COMPUTER CRIME

person nonetheless suffer from the same problems of self-selection bias as any other survey.

10.2.2.4 Instrument Validation. Interviews and other social sciences research methodologies can suffer from a systematic tendency for respondents to shape their answers to please the interviewer or to express opinions that may be closer to the norm in whatever group they see themselves. Thus, if it is well known that every organization ought to have a business continuity plan, some respondents may misrepresent the state of their business continuity planning to look better than they really are.

In addition, survey instruments may distort responses by phrasing questions in a biased way; for example, the question “Does your business have a completed business continuity plan?” may have a more accurate response rate than the question “Does your business comply with industry standards for having a completed business continuity plan?” The latter question is not neutral and is likely to increase the proportion of “yes” answers.

The sequence of answers may bias responses; exposure to the first possible answers can inadvertently establish a baseline for the respondent. For example, a question about the magnitude of virus infections might ask:

In the last 12 months, has your organization experienced total losses from virus infections of

- (a) \$1 million or greater;
- (b) less than \$1 million but greater than or equal to \$100,000;
- (c) less than \$100,000;
- (d) none at all?

To test for bias, the designer can create versions of the instrument in which the same information is obtained using the opposite sequence of answers:

In the last 12 months, has your organization experienced total losses from virus infections of

- (a) none at all;
- (b) less than \$100,000;
- (c) less than \$1 million but greater than or equal to \$100,000;
- (d) \$1 million or greater?

The sequence of questions can bias responses; having provided a particular response to a question, the respondent will tend to make answers to subsequent questions about the same topic conform to the first answer in the series. To test for this kind of bias, the designer can create versions of the instrument with questions in different sequences.

Another instrument validation technique inserts questions with no valid answers or with meaningless jargon to see if respondents are thinking critically about each question or merely providing any answer that pops into their heads. For example, one might insert the nonsensical question, “Does your company use steady-state quantum interference methodologies for intrusion detection?” into a questionnaire about security and invalidate the results of respondents who answer yes to this and other diagnostic questions.

Finally, independent verification of answers provides strong evidence of whether respondents are answering truthfully. However, such intrusive investigations are rare.

SUMMARY 10 · 11

10.2.2.5 Meta-analysis. Sometimes it is useful to evaluate a hypothesis based on several studies. We can combine probabilities P of the *same* null hypothesis from k trials using the formula

$$X^2 = -2\sum \ln P$$

where X^2 is distributed as $\chi^2_{[2k]}$ if the null hypothesis is true in all the trials.

For example, suppose a forensic specialist is evaluating the possibility that log files have been tampered with in a specific period.

- She runs a test of the frequency distribution of individual digits in the data for the suspect period to evaluate the likelihood that the distribution is consistent with the null hypothesis of randomness. The P for the two-tailed chi-square test is 0.072ns.
- She looks at the average number of disk I/Os per second in the records for the suspect period and compares the data with the same statistic in the control period using ANOVA; the P for the two-tailed test is 0.096ns.

In this case, $X^2 = -2 \sum \ln P$ $X^2 = -2 * (\ln 0.072 + \ln 0.096) = -2 * (-2.63109 - 2.34341) = 9.948992$ with 4 degrees of freedom. The P for the null hypothesis is 0.0413*. In other words, the chances of observing the results of both tests by chance alone if the suspect data were consistent with the raw data from comparison data is statistically significant. There is reason to reject the null hypothesis: Someone may very well have tampered with the log files for the suspect period.

This technique is subject to constraints. Most meta-analyses require the probabilities of the null hypothesis to be computed for a single tail in the same direction; it doesn't make sense to combine probabilities from conflicting hypotheses using two-tailed probabilities. However, as in the example above, if the null hypotheses are consistent, even two-tailed probabilities may be usefully combined.

Another problem is more systemic. There is considerable reason to be concerned that investigators and publishers sometimes suppress experimental results that do not conform to their expectations or desires. Such suppression biases the published results to *appear* to support the desired result.

Identifying such suppression is difficult. One approach is to examine the distributions of the published data and look for indications of data exclusion. For example, if the frequency distribution for raw data in a published report shows an abrupt disappearance of data in one direction (e.g., if the frequency distribution looks like a normal curve except that the left side suddenly drops to zero at a certain point) then the publication may be using a truncated data set. Meta-analysis based on data of dubious validity will itself be dubious. Garbage in, garbage out.

10.3 SUMMARY. In summary, all studies about computer crime should be studied carefully before we place reliance on their results. Some basic take-home questions about such research:

- What is the population we are sampling?
- Keeping in mind the self-selection bias, how representative of the wider population are the respondents who agreed to participate in the study or survey?
- How large is the sample?
- Are the authors testing for the assumptions of randomness, normality, and equality of variance before reporting statistical measures?

10 · 12 UNDERSTANDING STUDIES AND SURVEYS OF COMPUTER CRIME

- What are the confidence intervals for the statistics being reported?
- Are comparisons confounding variables?
- Are correlations being misinterpreted as causal relations?
- Were the test instruments validated?

10.4 FURTHER READING

Textbooks

If you are interested in learning more about survey design and statistical methods, you can study any elementary textbook on the social sciences statistics. Here are some sample titles:

- Babbie, E. R., F. S. Halley, and J. Zaino. *Adventures in Social Research: Data Analysis Using SPSS 11.0/11.5 for Windows*, 5th ed. Pine Forge Press, 2003.
- Bachman, R., and R. K. Schutt. *The Practice of Research in Criminology and Criminal Justice*, 3rd ed. Sage Publications, 2007.
- Carlberg, C. *Statistical Analysis: Microsoft Excel 2010*. Que, 2011.
- Chambliss, D. F., and R. K. Schutt. *Making Sense of the Social World: Methods of Investigation*, 2nd ed. Pine Forge Press, 2006.
- Cox, D. R., and C. A. Donnelly. *Principles of Applied Statistics*. Cambridge University Press (ISBN 978-1107644458), 2011. 212 pp.
- Kabay, M. E. *Statistics in Business, Finance, Management, and Information Technology: A Layered Introduction with Excel*. Free textbook (PDF), 2013. 205 pp.
www.mekabay.com/courses/academic/norwich/qm213/statistics.text.pdf
- Sirkin, R. M. *Statistics for the Social Sciences*, 3rd ed. Sage Publications, 2005.
- Warner, R. M. *Applied Statistics: From Bivariate through Multivariate Techniques*, 2nd ed. SAGE Publications (ISBN 978-1412991346), 2012. 1208 pp.

Websites

- Education Insider (2011) “Explore Statistics in the Blogosphere: Top 10 Statistics Blogs.” http://education-portal.com/articles/Explore_Statistics_in_the_Blogosphere_Top_10_Statistics_Blogs.html
- StatPac, “Survey & Questionnaire Design,” www.statpac.com/surveys

10.5 NOTES

1. GAO (1996) “Computer Attacks at Department of Defense Pose Increasing Risks.” General Accounting Office Report to Congressional Requesters GAO/AIMD-98-84 (May 1996), p. 19
2. M. E. Kabay, “CATA: Computer-Aided Thematic Analysis,” 2006. Available at www.mekabay.com with narrated lectures at www.mekabay.com

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 11

FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

William A. Zucker and Scott J. Nathan

11.1	INTRODUCTION	11·2		
11.2	THE MOST FUNDAMENTAL BUSINESS TOOL FOR PROTECTION OF TECHNOLOGY IS THE CONTRACT	11·3		
11.2.1	Prevention Begins at Home—Employee and Fiduciary Duties	11·4	11.4.4	Fair Use Exception
11.2.2	Employment Contract, Manual, and Handbook	11·4	11.4.5	Formulas Cannot Be Copyrighted
11.2.3	Technology Rights and Access in Contracts with Vendors and Users	11·4	11.4.6	Copyright Does Not Protect the “Look and Feel” for Software Products
11.3	PROPRIETARY RIGHTS AND TRADE SECRETS	11·5	11.4.7	Reverse Engineering as a Copyright Exception
11.3.1	Remedies for Trade Secret		11.4.8	Interfaces
11.3.2	Misappropriation Vigilance Is a Best Practice	11·6 11·8	11.4.9	Transformative Uses
			11.4.10	Derivative Works
			11.4.11	Semiconductor Chip Protection Act of 1984
			11.4.12	Direct, Contributory, or Vicarious Infringement
			11.4.13	Civil and Criminal Remedies
				DIGITAL MILLENNIUM COPYRIGHT ACT
				11·14
11.4	COPYRIGHT LAW AND SOFTWARE	11·8	11.6	CIRCUMVENTING TECHNOLOGY MEASURES
11.4.1	Works for Hire and Copyright Ownership	11·9	11.6.1	Exceptions to the Prohibitions on Technology Circumvention
11.4.2	Copyright Rights Adhere from the Creation of the Work	11·9		11·16
11.4.3	First Sale Limitation	11·9	11.7	PATENT PROTECTION
			11.7.1	Patent Protection Requires Disclosure
				11·18

11 · 2 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

11.7.2	Patent Protection in Other Jurisdictions	11·19	11.10.3	Other Open Source Licenses	11·34
11.7.3	Patent Infringement	11·19	11.10.4	Business Policies with Respect to Open Source Licenses	11·34
11.8	PIRACY AND OTHER INTRUSIONS	11·20	11.11	APPLICATION INTERNATIONALLY	11·34
11.8.1	Marketplace	11·20	11.11.1	Agreement on Trade-Related Aspects of Intellectual Property Rights	11·35
11.8.2	Database Protection	11·20	11.11.2	TRIPS and Trade Secrets	11·36
11.8.3	Applications of Transformative and Fair Use	11·21	11.11.3	TRIPS and Copyright	11·37
11.8.4	Internet Hosting and File Distribution	11·21	11.11.4	TRIPS and Patents	11·37
11.8.5	Web Crawlers and Fair Use	11·23	11.11.5	TRIPS and Anticompetitive Restrictions	11·38
11.8.6	HyperLinking	11·23	11.11.6	Remedies and Enforcement Mechanisms	11·38
11.8.7	File Sharing	11·23			
11.9	OTHER TOOLS TO PREVENT UNAUTHORIZED INTRUSIONS	11·24	11.12	RECENT DEVELOPMENTS IN INTELLECTUAL PROPERTY LAW	11·39
11.9.1	Trespass	11·24	11.12.1	AIA	11·39
11.9.2	Terms of Use	11·25	11.12.2	The PROTECT IP Act (PIPA)	11·39
11.9.3	Computer Fraud and Abuse Act	11·26	11.12.3	The Stop Online Piracy Act (SOPA)	11·41
11.9.4	Electronic Communications and Privacy	11·29	11.12.4	Patent Trolls	11·43
11.9.5	Stored Communications Act	11·31			
11.10	OPEN SOURCE	11·33	11.13	CONCLUDING REMARKS	11·44
11.10.1	Open Source Licenses	11·33	11.14	FURTHER READING	11·44
11.10.2	GPL	11·33	11.15	NOTES	11·44

11.1 INTRODUCTION. This chapter is not for lawyers or law students. Rather, it is written for computer professionals who might find it useful to understand how their concerns at work fit into a legal framework, and how that framework shapes strategies that they might employ in their work. It is not intended to be definitive but to help readers spot issues when they arise and to impart an understanding that is the first part of a fully integrated computer security program.

Cyberlaw is a compendium of traditional law that has been updated and applied to new technologies. When gaps have developed or traditional law is inadequate, particular statutes have been enacted. It is a little like the old story of the three blind men and the elephant: One of the blind men touching the elephant's leg believes he is touching a tree; the other touching its ear believes it is a wing, and the third, touching the tail, thinks it is a snake. Issues of cyberspace, electronic data, networks, global transmissions, and positioning have neither simple unitary solutions nor a simple body of law to consult.

THE MOST FUNDAMENTAL BUSINESS TOOL FOR PROTECTION 11 · 3

In thinking about the application of law to computer security, it is helpful to think about the problems as issues in which the computer is

- The target of the activity
- The tool used for the activity
- Incidental to the activity itself

For example, “hacking” into a computer can be analogized to the tort¹ of trespass (i.e., entering the property of another without permission), and “cracking” can be viewed as conversion of someone else’s property. Similarly, using the computer to make illegal copies is a violation of copyright law in its most basic sense. Although trademark law has very little to do with computers, using trade names as part of keywords for search engines, or domain names to misdirect Internet traffic to a competitive Web site can be a violation of a trademark. While touching on some of the more traditional tort remedies, this chapter focuses on the property rights being invaded by such activities and the remedies that exist in the context of a business operation.

Recognizing that the body of law which touches on these problems is as global as the Internet itself, this chapter is intended to help readers actually see the elephant in the room. In selecting what legal issues to highlight, we have tried to consider the routine needs of the computer professional. We have focused largely on the law of the United States, recognizing that these problems and subject matters often transcend national boundaries. There is a very simple reason for this. Most often, the impact of the computer security attack, denial of service, decryption, or theft of computer materials will have occurred here, or have a direct impact here, no matter where it originates. Imagine for a second a gunman—standing in Canada—who takes aim at someone in the United States, pulls the trigger, and hits his target. Since there is purposeful conduct aimed at this country, in the ordinary instance the U.S. judiciary will not only assert jurisdiction over the gunman but also apply its laws. There may be other problems, such as actually catching the gunman, but the example underlines the importance of the law of the United States for entities located here. For orientation purposes, we have also included a section at the end of this chapter that discusses some international issues.

One other introductory note: We use the phrase “security program” in this chapter with some frequency. Understanding that this phrase can mean one thing to a lawyer or risk manager and another thing to a computer security professional, we intend it as a shorthand reference to the generic and systemic effort to secure information stored on computers and not solely to the applications that may be employed as part of that effort.

11.2 THE MOST FUNDAMENTAL BUSINESS TOOL FOR PROTECTION OF TECHNOLOGY IS THE CONTRACT. The computer security professional’s job is to understand, anticipate, and then worry about risk: risks that are beyond control and risks that can be controlled. The most fundamental tool for controlling risk, whether predictable or unforeseeable, is the contract. Unlike other forms of risk control, a contract need not be static; it can be adaptable. We can limit use; we can limit distribution; we can impose conditions and confidentiality; we can specify rights as well as provide for certain remedies through contract. Contracts actually can take many forms: the traditional signed agreement; an email exchange; Web site or product terms of use; employment agreements; workplace manuals and policies; and so-called shrink-wrap or click-wrap agreements. We sell or license products. Where we can contract, we can also define and limit risk.

11 · 4 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

11.2.1 Prevention Begins at Home—Employee and Fiduciary Duties.

There is an old hoary concept in the law that employees owe to their employers the fiduciary duty of utmost loyalty. The scope and extent of that fiduciary duty is a matter of common law that varies in each state. Generally, employees' fiduciary duty prohibits them from using any property that belongs to the employer in competition with the employer or for personal gain. Employees, however, are entitled to retain and use for whatever purpose their own skill and knowledge, which arguably could include contacts that they develop over the course of their employment unless those contacts are trade secrets. What comes or does not come within the ambit of fiduciary duty has spawned endless arguments and lawsuits. There is a simple remedy to this problem: the contract that covers technology issues and ownership as well as it covers pay and other benefits.

11.2.2 Employment Contract, Manual, and Handbook. Whatever policy the security professional develops should be implemented through the organization's employment contract, manual, and handbook. Many contractual provisions can be applied, such as: nondisclosure agreements; definition of proprietary policy; restrictive covenants; concessions of ownership regarding discoveries, know-how, improvements, inventions, and the like during the term of employment; email policies; terms of use regarding computer systems; and statements of authorized and unauthorized activity. The point is that employment contracts and handbooks should be the starting point for computer security.

11.2.3 Technology Rights and Access in Contracts with Vendors and Users.

Security. Security protection necessarily includes vigilance about all contracts and licenses with vendors and users. This may not be sexy, but it is blocking and tackling. Vendors can be subject to many of the same limitations and nondisclosure agreements as employees. Rights of access to intranets and data should be controlled and privileges specified. Careful consideration should be given to what rights a user will have, the rules surrounding user access, and enforcement of those rules. Is this a sale or a license? There are many virtues to controlling technology through licenses (as opposed to sales), including imposing limits on rights of use, and specifying remedies for breaches of the license, or for unauthorized activity that involves the licensed product.

“Shrink-wrap” or “click-wrap” licenses have become common parlance. They are now accepted tools for licensing and controlling software distribution so long as: (a) they are business to business and thus between parties of roughly equal bargaining position; (b) their terms for other users or consumers are not unconscionable; and (c) they do not violate public policy. Concerns over whether contractual terms are unconscionable or the contracts are ones of adhesion arise because the licenses are not products of negotiation but of fiat, which users accept when they open the shrink-wrapped package or through an online click. These concerns have been addressed through requirements that users have been provided with adequate notice of the terms, an opportunity to reject, and conduct that sufficiently manifests consent. For shrink-wrap agreements, the opening of the product, its installation, and retention have been deemed sufficient acts to show consent to the terms of the license, noting that if the consumer does not wish to consent, the product could be returned.² Thus, it is not necessary for the prospective user to be aware of all of the terms of a license before purchase if the remedy includes return after purchase. The license can impose restrictions on use, limit the number of machines on which the product can be installed, copying, and even available remedies.³

PROPRIETARY RIGHTS AND TRADE SECRETS 11 · 5

The issues of notice, actual or constructive, an opportunity to accept or reject, and manifestation of consent have led to general acceptance of online agreements such as the presentation of licensing terms followed by an active need to check, accept, or reject by clicking on the appropriate box.⁴ The same analysis applies to terms of use especially for intranet or network use.⁵ In *Register.com v. Verio, Inc.*,⁶ downloading data from a WHOIS database, having knowledge of the terms of use, was acceptance of those terms even if there was no click-through. These examples show that terms of use, properly positioned, can be binding on the user.

An active security program begins with a review of the contracts, licenses, and terms of use in all relationships with your organization. Just because a contractual arrangement has not existed does not mean that you cannot create one through proper notice of the terms of the contract and conduct that shows assent to those terms. Such contracts are the security professional's first line of defense. They give you the ability to limit risk with an organization's employees, contractors, vendors, and affiliates. With that in mind, this chapter addresses issues that arise largely outside of the terms of contractual protections and also suggests additional potential self-help remedies.

11.3 PROPRIETARY RIGHTS AND TRADE SECRETS. For many years, unless an idea was patentable, the primary protection for internal business data, confidential or proprietary information, and computer code was through the common law doctrine of trade secrets.⁷ Generally, a trade secret might be considered any internal, nonpublished manufacturing know-how, drawings, formulas, or sales information used in a trade or business that has commercial applicability and that provides a business with some strategic advantage.⁸ Such information, so long as it was (a) not published or disseminated to others who were not obligated to maintain its confidentiality,⁹ and (b) maintained in confidence with the protecting organization, could be protected as a trade secret.

The law of trade secret thus recognized a business's ownership or proprietary interest in such information, data, or processes. There are, however, important practical limitations on the application of trade secret protection. First and foremost, for any product sold in the market, the law does not protect against a competitor seeing the product and then using it to figure out how to manufacture like or similar items. Competitors are therefore free to reverse engineer a product so long as the reverse engineering is done wholly independently.

The second caveat is that an organization has to prove not only that the information qualifies for trade secret protection, but also that it protected the secrecy of the information as required by the law of the applicable jurisdiction. This means that ownership will be a matter not of record but of case-by-case proof, making enforcement of trade secret protection time consuming and expensive. Generally, the required proof consists of a showing that there was an active security program in place that was sufficient to protect the information as confidential. Various programs may be deemed adequate, depending on the circumstances, but usually such programs have five principles in common:

1. An inventory of trade secret information that is periodically updated
2. A security program to protect the technology at issue, often on a need-to-know basis with clear marking of information as "confidential, access restricted"
3. A written description of the security program that is provided to all employees

11 · 6 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

4. An enforcement officer or oversight procedure
5. An enforcement program, including litigation, if necessary, to enjoin unauthorized access or distribution

In the field of computing, these principles often mean that source code or other readable formats should be secured in a locked file and marked CONFIDENTIAL. All representations of the code as stored on magnetic or other media should be marked CONFIDENTIAL and secured. Computerized information should be password protected with restrictions on circulation of the password and periodic password changes. A notice of confidentiality should be displayed as soon as access to the program is obtained, with appropriate warnings on limitation of use. Levels of access should be controlled so that privileges to copy, read, and write are appropriately restricted. Surveillance of entries and logon should be routinely conducted to verify that there has been no unauthorized entry. Finally, periodic audits should be conducted to test and substantiate the security procedures.

For many years, each state developed its own brand of trade secret protection through evolving judicial decisions that establish something in this country called the common law, as distinguished from legislative enactments of a statute addressing the same issue. In 1985, the Uniform Trade Secrets Act (UTSA) was promulgated by the National Conference of Commissioners on Uniform State Laws, with one of its purposes to make uniform the rights and remedies available to a holder of a trade secret. This model law, however, needed to be adopted by each state before it became the law of the state. As of this writing, it has been adopted to some degree in 46 states with the exception of Massachusetts, New Jersey, New York, and Texas.

The UTSA defines a trade secret as information, including a formula, pattern, compilation, program device, method, technique, or process, that: (a) derives independent economic value, actual or potential, from *not being generally known to, and not being readily ascertainable by proper means by*, other persons who can obtain economic value from its disclosure or use; and (b) is the subject of efforts that are reasonable under the circumstances to maintain its secrecy. It also defines the unlawful taking of a trade secret, or misappropriation, as the wrongful use of a trade secret, including (a) knowingly acquiring the secret through improper means or (b) disclosing the secret without consent.

11.3.1 Remedies for Trade Secret Misappropriation. Misappropriation of a trade secret is the unauthorized use or disclosure of the trade secret. In simple parlance, it is a taking or theft. The taking can be by one who owes a fiduciary duty of confidentiality, such as an employee; it can be in breach of an agreement of confidentiality; or the taking can occur through improper access or means. The misappropriation can be treated under common law as the tort of conversion, trespass, unfair competition, or interference with contractual relations. As discussed, there are now specific statutory provisions under the UTSA for trade secret misappropriation. The UTSA grants the wronged party certain remedies that include enjoining the use of the misappropriated property, damages, and attorney's fees. When the misappropriation is of a physical item, such as a disk drive, the owner may ask the court to order seizure and return of its property.¹⁰ In addition, where the misappropriation also violates other laws protecting intellectual property, such as where the taking infringes a copyright, the property owner may be entitled to additional relief.

Exactly what remedies are available will vary among the states. Interestingly, the very uniformity that the UTSA was intended to create has led to different treatment of

PROPRIETARY RIGHTS AND TRADE SECRETS 11 · 7

available claims and remedies. For example, before the UTSA, an employee's theft of the employer's confidential customer lists triggered a common law claim for breach of the implied fiduciary obligation owed by an employee to the employer as well as a claim for misappropriation of trade secrets. The UTSA provides that its remedies preempt other common law remedies; in other words, a claim under the UTSA trumps the claim for breach of fiduciary duty as well as the claim for misappropriation of trade secrets. There is a split in the courts as to whether the UTSA replaces only common law causes of action for misappropriation of trade secrets or extends to any tortious claims for relief that arise out of the misappropriation no matter how stated. The broader reach of the UTSA appears to be favored by the growing majority of courts that have considered this issue to date. The takeaway from this uncertainty is that computer security professionals should protect trade secrets, confidential information, and other valuable data through contractual terms with, among others, employees, vendors, and users to minimize the reliance on the UTSA.

In the event of a misappropriation, in addition to civil remedies, often separate state statutes treat the taking as a theft and a criminal act. Such statutes are generally state specific. Prior to 1996, the Trade Secrets Act (TSA) was the only federal statute prohibiting trade secret misappropriation. The TSA, however, was of limited utility because it did not apply to private sector employees and provided only limited criminal sanctions.¹¹ To combat an increase in computer crimes, Congress enacted the Economic Espionage Act of 1996 (EEA), which provided greater protection for the proprietary and economic information of both corporate and governmental entities against foreign and domestic theft.¹²

The EEA criminalizes two principal categories of corporate espionage: economic espionage and theft of trade secrets.¹³ Section 1831 punishes those who steal trade secrets "to benefit a foreign government, foreign instrumentality, or foreign agent." Section 1832 is the general criminal trade secret provision.¹⁴ The EEA criminalizes stealing, concealing, destruction, sketching, copying, transmitting, or receiving trade secrets without authorization, or with knowledge that the trade secrets have been misappropriated. It also criminalizes attempting to and conspiring to do any of these acts.¹⁵ The EEA penalizes parties responsible for a taking that is intended to benefit a foreign government with fines up to \$250,000 and imprisonment up to 10 years.¹⁶

The EEA explicitly defines a trade secret to include information stored in electronic media and includes "programs or codes, whether tangible or intangible" so long as:

- (a) the owner thereof has taken reasonable measures to keep such information secret; and
- (b) the information derives independent economic value, actual or potential, from not being generally known to, and not being readily ascertainable through proper means by, the public.¹⁷

Although one might assume that this definition is relatively straightforward, not everything is as it appears. In a case of domestic trade secret theft, the Court of Appeals for the Seventh Circuit examined what the EEA means when it says that the data or material "derives independent economic value, actual or potential, from not being generally known to, and not being readily ascertainable through proper means by, the public."¹⁸ Noting that others had assumed that the word "public" meant the general public, the court in *Lange* astutely observed that this was not, in fact, the case. Moreover, the standard for measuring the persons who might readily ascertain the economic value of (in this case) the design and composition of airplane brake assemblies is not the average person in the street, for this assumes (as the court mentions) that any person can understand and apply something as arcane as Avogadro's number. Instead, the definition

11 · 8 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

of the term “the public” should take into account the segment of the population that would be interested in and understand the nature of that which has allegedly been misappropriated.

The international reach of the act is limited, extending outside of the United States only if: “(1) the offender is a natural person who is a citizen or permanent resident alien of the United States, or an organization organized under the laws of the United States or a State or political subdivision … or (2) an act in furtherance of the offense was committed in the United States.”¹⁹ Few defendants have been charged under the act since its passage in 1996, so the precise reach has yet to be tested. However, the language of the EEA applies its provisions to corporations with headquarters or operations subject to U.S. jurisdiction that could be prosecuted under the act. Finally, the remedies under the EEA can be invoked only by the United States. There is no private right of action under the act.

11.3.2 Vigilance Is a Best Practice. The key points of practice to remember are: Security and trade secret law are forever linked together. A trade secret cannot exist without such security. The maxim “Eternal vigilance is the price of liberty,” often attributed to Thomas Jefferson, should in the context of business information protection be restated as “Eternal vigilance is the price of trade secret protection.” It is not as catchy a phrase, but it is the price each business must pay if it relies in whole or in part on trade secret law for protection. In such situations, the greatest assurance of protection can be obtained through rigorous contractual terms and strenuous enforcement.

11.4 COPYRIGHT LAW AND SOFTWARE. Because of anxiety over the true extent of protection afforded software under patent and copyright law, software programs initially were protected as trade secrets. Such protection became increasingly problematic in today’s society, where information technology and pressure for the free flow of information makes confidentiality controls more difficult to police. Copyright law now has evolved to include computer programs.

Since 1964, the United States Copyright Office has permitted registration of computer programs, although judicial decisions were divided on the applicability of the Copyright Act. In 1976, Congress passed the Copyright Act of 1976, which did little to resolve the ambiguity. Clarification finally was obtained in the Computer Software Copyright Act of 1980, which explicitly extended the protection of the copyright laws to software.²⁰ Any type of work that can be fixed in any tangible medium can be protected by copyright as literary works based on the authorship of the source and object code²¹ even if the work can only be machine reproduced.

Copyright protection, however, does not protect “ideas.”²² Rather, it protects the particular expression of the idea. As can be seen by the parallel proliferation of spreadsheet programs, the idea for the spreadsheet program cannot be protected, but the particular code that produces the spreadsheet can be. In order to qualify for copyright protection, the work must be (a) original, (b) fixed in a tangible medium, and (c) not just the embodiment of an idea. Once obtained, copyright protection grants to the copyright owner the exclusive right to reproduce, to publish, to prepare derivative works, to distribute, to display, and to perform the copyrighted work. In 1990, Congress passed the Computer Software Rental Amendments Act,²³ which added to the list of copyright infringements the distribution of a computer program for commercial advantage. Materials copyrighted after 1978 are protected for the lesser of 75 years from the date of first publication or 100 years from the date of creation.

COPYRIGHT LAW AND SOFTWARE 11 · 9

11.4.1 Works for Hire and Copyright Ownership. The copyright for a work does not always belong to the person who creates it. The most frequent exceptions are works that fall under the concept of a “work for hire.” A work for hire is not owned by the creator but by the persons who hired the creator to create the work. Most often the concept applies to employees who have created a work *within the scope* of their employment. The key concept is the scope of employment. Even though a work is created outside of the office and normal working hours, it still will be a work for hire if it is within the scope of employment. However, a work that falls outside the scope of employment and that is created outside the office is likely not to be deemed a work for hire. Because of such issues, it is better practice when dealing with employees or independent contractors to provide specificity in an agreement as to what is a work and when the creation of a work will be governed by the doctrine of work for hire.

11.4.2 Copyright Rights Adhere from the Creation of the Work. Everyone who has looked at a copyrighted work is probably familiar with the symbol © affixed to any published copyrighted work, together with the name of the copyright holder and the year of creation or publication of the work. For many years, such notice was a *fortiori* necessary for copyright protection. Today, however, the copyright arises from the creation of a copyrighted work itself. It is still good practice to advise the world of potential infringement by inserting the formalities of a copyright on the work itself. In addition, one should register the work in the United States Copyright Office, which is currently developing a process for online registration. Registration of the copyright also permits one to claim statutory damages ranging from \$500 to \$20,000 for each violation, which often is useful to prevent additional infringements when no actual damages can be demonstrated. Moreover, in some jurisdictions, it may be necessary to register the copyright with the copyright office before one can actually sue to protect the copyright.

The change in copyright protection has interesting applications when applied to electronic works. The creation of the work in some permanent form is sufficient to trigger copyright protection. Thus, the creation of an electronic copy is sufficient permanency. What that means is that any electronic data are already conceivably subject to copyright protection at the time that they are viewed or received. Thus, in using any such information or “work,” care must be taken that one does not infringe on a potential copyright without a license.

11.4.3 First Sale Limitation. The holder of the copyright has the right to sell or license the work. If the work is sold, the holder essentially loses all rights to control the resale of the work. This is known as the first sale doctrine. Once the item is placed in commerce, subsequent transfers cannot be restricted. The doctrine applies only to the copy that has actually been sold. It does not create a license to copy the item itself.

To avoid what sometimes can be a problem if the program winds up in the hands of a competitor, companies often prefer to license the item instead of selling it outright. If the work is licensed, only those rights that are contained in the license are transferred. All other rights of ownership remain with the licensor. Thus, a breach of the license gives the licensor of the copyrighted work the right to reclaim the work or prevent its further use or publication. However, if the license has all the basic indicia of a sale, it will be treated as one, notwithstanding the label.

One interesting intersection of these two principles is the requirement when upgrading software that the old version be present. As a condition of making the upgrade

11 · 10 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

available at a reduced rate, the seller normally requires that the older version be authenticated before the newer version can be installed. Such requirements are legal, as the owner of the earlier version could choose to sell it, but then would have to pay a higher price for the newer version and work to restrain subsequent sales of software from a user who expects to upgrade in the future.

11.4.4 Fair Use Exception. All copyright protection is subject to the doctrine of fair use.²⁴ Fair use permits the use of a work without authorization for a limited purpose. But what use constitutes fair use? The Copyright Act of 1976 suggested four, nonexclusive factors, for a court to consider:

1. What is the purpose and character of the use?
2. What is the nature of the copyrighted work?
3. How much of the copyrighted work is used?
4. What is the effect on the potential market for the work?

Despite its codification in the Copyright Act of 1976, fair use remains a nebulous doctrine—an equitable rule of reason, with each case to be decided on its own facts.²⁵ It is often misquoted and misapplied. The essential concept behind the doctrine of fair use is to permit public discussion, review, and debate of a copyrighted work without violating the copyright. Thus, the Copyright Act of 1976 gives as examples of fair use, situations of “criticism, comment, news reporting, teaching (including multiple copies for classroom use), scholarship, or research.”

Fair use is not an antidote for failing to license a work. It should be invoked with care—understanding that the more material that is used and the more commercial the purpose, the less likely a court will find it applicable. Indeed, sometimes the only way to harmonize cases on whether a use is a fair use is to decide whether the court ultimately viewed the user as a “good” or “bad” guy.

11.4.5 Formulas Cannot Be Copyrighted. There are limitations on what expressions can be protected by copyright law. A frequent source of argument is whether, since one cannot protect the idea, the expression is directly driven by its content (i.e., the expression is simply a function of the idea). For that reason, formulas cannot be copyrighted.²⁶ This means that when formulas are part of a computer program, other modes of defense need to be considered, such as trade secret or possibly patent protection. If one were to disclose the formula through copyright publication, one would lose the ability to protect that information.

11.4.6 Copyright Does Not Protect the “Look and Feel” for Software Products. Copyright protection ordinarily extends to the physical manifestation of the computer program in the source code and object code. The operation of that code, as it translates to what the human mind perceives, has been described as the “look and feel” of the program. In attempting to quantify the concept of “look and feel,” courts have considered whether the organization, structure, and sequence of the program can be protected. In the United States, *Whelan Associates, Inc. v. Jaslow Dental Lab, Inc.*,²⁷ gave the greatest extension to protecting look and feel. In that case, none of the code had been copied and the program operated on a different platform. Nonetheless, copyright infringement was found because the organization, structure, and sequence of the program had been copied. The court recognized that the structure and logic of

COPYRIGHT LAW AND SOFTWARE 11 · 11

the program are the most difficult to create and that the idea could be protected as it was embodied in the program structure since, given the variety that was possible, the structure was not necessarily just an extension of the idea. Since *Whelan*, courts in the United States have retreated from such broad protection. In 1992, *Computer Associates, Inc. v. Altai, Inc.*²⁸ developed the so-called abstraction-filtration test. The results of that test define as unprotectable: (a) program structures that are dictated by operating efficiency or functional demands of the program and therefore deemed part of the idea and (b) all tools and subroutines that may be deemed included in the public domain. Only what remains is to be compared for possible copyright infringement.

While protection of look and feel may vary among the different federal circuits, in general, the courts are swinging away from broader protection. However, this may not necessarily be true internationally; English law appears to grant the broader protections afforded by the *Whelan* decision.

11.4.7 Reverse Engineering as a Copyright Exception. Within the field of computer software, cases have considered whether “dissection” in order to reverse engineer the program is a violation of the copyright. To those involved in protecting software programs, as well as those involved with competing products, the answer appears to be that reverse engineering does not constitute an infringement, even though the disassembly of the program falls squarely within the category of acts prohibited by the Copyright Act because of the doctrine of fair use. The Ninth Circuit in *Sega Enterprises Ltd. v. Accolade, Inc.*²⁹ found as a matter of law that:

... where disassembly is the only way to gain access to the ideas and functional elements embodied in a copyrighted computer program and where there is a legitimate reason for seeking such access, disassembly is a fair use of the copyrighted work.³⁰

The Ninth Circuit is not the only circuit that has upheld reverse engineering against a copyright claim. The Federal Circuit reached a similar conclusion regarding reverse engineering of object code to discern the “ideas” behind the program in *Atari Games Corp. v. Nintendo of America, Inc.*³¹ The fair use rationale of *Sega* was also adopted by the Eleventh Circuit in *Bateman v. Mnemonics, Inc.*³² on the grounds that it advanced the sciences. In addition, in *Assessment Techs. of WI, LLC, v. WIREDATA, Inc.*, the Seventh Circuit relied on *Sega* and determined that WIREDATA, Inc. could extract uncopyrighted data from a copyrighted computer program, noting that the purpose of the extraction was to get the raw data, not compete with Assessment Technologies by selling copies of the program itself.³³ In *Evolution, Inc. v. SunTrust Bank*, the Tenth Circuit relied on both *Sega* and WIREDATA when it allowed the defendant to copy part of plaintiff’s source code to extract uncopyrighted data from plaintiff’s copyrighted computer program.³⁴ Thus, unless careful thought is given to the application of copyright protection, merely copyrighting the software will not necessarily protect against imitation.

11.4.8 Interfaces. There is an open issue as to whether copyright protects the format for interfacing between application and data. Competitors, particularly in the area of gaming, look to reverse engineer the interface format to make new modules compatible with existing hardware. Such reverse engineering has been held not to violate the copyright laws, so long as the new product does not display copyrighted images or other copyrightable expressions.³⁵ Thus, the nonprotectable interface may be protected if such copyrighted images or expressions are embedded in the display.

11 · 12 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

11.4.9 Transformative Uses. One of the factors that the doctrine of fair use considers is the “amount and substantiality of the portion used in relation to the copyrighted work as a whole.”³⁶ In practical terms, this means that courts look at how much was taken and for what purpose. One could take a little but still take the essence of the program. One could also take a little that did not attempt to duplicate but rather used the copyrighted material as a springboard for a new creation. Out of this qualitative and quantitative investigation comes the notion of transformative use, which became the coin of analysis in the Supreme Court’s 1994 decision in *Campbell v. Acuff-Rose Music, Inc.*³⁷ *Campbell* addressed the concept in terms of a claim of copyright infringement involving a rap parody of a popular song. There, taking its clues from the opening language of Section 107 codifying fair use, the Supreme Court asked whether the “new” work “adds something new, with a further purpose or different character, altering the first with new expression, meaning or message; it asks, in other words, whether and to what extent the new work is transformative.”³⁸ The Court then laid down the test to be applied.

Although such transformative use is not absolutely necessary for a finding of fair use, … the goal of copyright, to promote science and the arts, is generally furthered by the creation of transformative works. Such works thus lie at the heart of the fair use doctrine’s guarantee of breathing space within the confines of copyright, … and the more transformative the new work, the less will be the significance of other factors, like commercialism, that may weigh against a finding of fair use.³⁹

Thus, a transformative use may play off of a prior copyright and still not be deemed an infringement so long as the resulting new work is just that—new.

11.4.10 Derivative Works. Under Section 106(2) of the Copyright Act of 1976, the copyright owner has the exclusive right “to prepare derivative works based upon the copyrighted work.” The act defines a “derivative work” as:

… a work based upon one or more pre-existing works, such as a translation, musical arrangement, dramatization, fictionalization, motion picture version, sound recording, art reproduction, abridgement, condensation, or any other form in which a work may be recast, transformed, or adapted. A work consisting of editorial revisions, annotations, elaborations, or other modifications which, as a whole, represent an original work of authorship, is a “derivative work.”

A “derivative work” is thus defined as an original work that is independently copyable. To infringe the exclusive right to prepare a derivative work granted by the Copyright Act to the copyright owner, the infringer need not actually have copied the original work or even have fixed in a tangible medium of expression the allegedly infringing work.⁴⁰ The right, therefore, to create the derivative work can be a useful tool in counterbalancing attempts to pirate computer programs and the issue of fair use.

The Copyright Act creates an exemption for a lawful owner of a purchased license for a computer program to adapt the copyrighted program if the actual adaptation “is created as an essential step in the utilization of the computer program in conjunction with a machine and it is used in no other manner.”⁴¹ The adaptation cannot be transferred to a third party. The right to adapt is, in essence, the right to modify or, in the language of the act, to create a derivative work. Such changes can be made even without the consent of the software owner so long as such modifications are used only internally and are necessary to the continuing use of the software.⁴²

COPYRIGHT LAW AND SOFTWARE 11 · 13

11.4.11 Semiconductor Chip Protection Act of 1984. The Semiconductor Chip Protection Act of 1984 (SCPA) protects as part of the Copyright Act “mask works fixed in a semiconductor product.”⁴³ The SCPA protects not the product itself but the copying of the circuit design or blueprint. Because of reverse engineering, the protections afforded by SCPA are limited in practice.

11.4.12 Direct, Contributory, or Vicarious Infringement. Copyright infringement generally requires a showing of substantial similarity between allegedly offending use and the protected expression contained in a work. Infringement can occur through the simple act of printing (without permission), by posting on the Web or other form of unauthorized distribution, by creating a derivative work, or by another act that interferes with the copyright holder’s rights.

A copyright can be infringed directly, contributorily, or vicariously. Direct infringement is the term ascribed to the actor who violates the copyright. Contributory infringement involves knowingly providing the means for the violation to occur. Liability for contributory infringement may be predicated on actively encouraging (or inducing) infringement through specific acts, or on distributing a product that distributees use to infringe copyrights, if the product is not capable of “substantial” or “commercially significant” noninfringing uses.⁴⁴ But secondary liability for copyright infringement does not exist in the absence of direct infringement by a third party. Vicarious infringement occurs when one is responsible for or controls the actions of another who violates the infringement. The usual situation is that of an employer’s responsibility for the acts of an employee.

Not all situations admit themselves of simple answers, as when a person commits direct infringement by actually photocopying a work. New technologies constantly pose issues as to whether infringement has occurred and whether the infringement violates the public interest. In general, when faced with an issue of potential copyright infringement, the questions to ask are:

- Can the product or service be used to infringe a copyright, or is the product capable of substantial noninfringing uses?
- If so, did the owner of the product or service encourage the user to use it for infringement?
- Alternatively, did the owner of the product or service have knowledge of the specific infringing use and have the ability to prevent it?

Today, we take Internet service providers (ISPs) for granted. But application of these questions initially led courts to conclude that ISPs were liable for contributory infringement. For example, a Web site that encouraged and facilitated the uploading of copyrighted materials was found to be a direct infringer of the copyright of the owner even though the provider did not actually do the uploading.⁴⁵ Similarly, an ISP that was notified of a copyright violation that was posted on its server and failed to correct it could be found to have contributory liability for the infringement.⁴⁶ In its wisdom, however, Congress, in the Digital Millennium Copyright Act (DMCA), created a safe harbor for ISPs so that, as a matter of public policy, an ISP does not have to monitor each and every transmission for potential copyright infringement.

11.4.13 Civil and Criminal Remedies. The Copyright Act contains several sections that specifically address the penalties and remedies for infringement. They

11 · 14 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

include injunctive relief (i.e., a court order terminating the infringing conduct),⁴⁷ impounding and disposing of infringing articles,⁴⁸ damages,⁴⁹ litigation costs and attorneys' fees,⁵⁰ and criminal penalties.⁵¹ Although this chapter cannot address all of the permutations of remedies and penalties available, a few are worth mentioning.

Generically, a copyright owner must choose between its actual losses (i.e., what it actually lost and any profits realized by the infringer) and statutory damages.⁵² Actual damages imply economic losses actually suffered as a result of the infringement. The kinds of actual damages that have been awarded include development costs of the software,⁵³ the economic consequences of lost customers,⁵⁴ lost future sales,⁵⁵ the value of the infringer's licensing fees where the licensor is precluded from market sales,⁵⁶ lost market value of the material infringed,⁵⁷ and lost royalty payments.⁵⁸ An award of actual damages is not automatic; the license holder has the burden of proving that the infringing activity and the economic loss are causally connected, at which point the infringing party must show that the license holder would have incurred the loss anyway.⁵⁹

A copyright owner may elect to receive statutory damages rather than actual damages and the infringer's profits.⁶⁰ Making the election is mandatory, and it must be done before final judgment is entered. Once the election is made, it is final. The statutory damages generally range from \$500 to \$20,000 "for all infringements involved in the action, with respect to any one work, for which any two or more infringers are liable jointly and severally. . . . For purposes of this section, all the parts of a compilation or derivative work constitute one work."⁶¹ This amount may be increased to \$100,000 if the court finds that the infringement was willful and reduced to \$200 if the court finds that the infringer "was not aware and had no reason to believe" that the act was an infringement.⁶²

Statutory damages theoretically⁶³ are intended to approximate the actual damages suffered, and were crafted as an alternative compensation scheme for copyright owners, when actual damages are difficult to calculate. In determining whether to elect actual or statutory damages, a copyright owner ought to perform a careful analysis to determine how many separate infringements occurred that justify, under the statute, separate awards. Although posting different copyrighted computer software programs on a bulletin board for downloading constitutes multiple infringements,⁶⁴ making multiple copies of the same cartoon character in different poses constitutes a single infringement because only one work was copied.⁶⁵

As mentioned, this is one of the statutory schemes that discourage frivolous litigation by imposing the cost of litigating on the losing party. The statute permits the substantially prevailing party to recover its reasonable attorneys' fees and costs from the losing party. Who is the substantially prevailing party and what constitutes reasonable attorneys' fees are separate and distinct issues that will be decided by the courts.

Copyright violations also can be criminally prosecuted, and generally require demonstration of *mens rea*, or intent. One or more infringements having a total retail value of more than \$1,000 within a 180-day period or "for purposes of commercial advantage or private financial gain" can be punished by one to five years of imprisonment and fines. Even without demonstration of a motive of financial gain, 10 or more infringements having a value in excess of \$2,500 can result in up to three years in jail and fines. Repeated violations carry stiffer penalties. Finally, one who knowingly aids or abets a copyright infringement is also subject to criminal prosecution.

11.5 DIGITAL MILLENNIUM COPYRIGHT ACT.

In 1998, Congress passed the Digital Millennium Copyright Act (DMCA) to address concerns raised by the

CIRCUMVENTING TECHNOLOGY MEASURES 11 · 15

Internet and copyright issues in the context of our increasingly technological society. The DMCA creates a civil remedy for its violation as well as criminal penalties starting after October 2000. One of the purposes of the DMCA is to protect the integrity of copyright information. Removal of a copyright notice, or distribution knowing that such copyright has been removed, is now actionable.⁶⁶

11.6 CIRCUMVENTING TECHNOLOGY MEASURES. Article 11 of the World Intellectual Property Organization Copyright Treaty required all signatory countries to provide adequate legal protection and remedies against the circumvention of technical measures intended to secure copyrights. In response, Congress adopted Section 1201 of the DMCA, which generally prohibits the act of circumventing, and trafficking in the technology that enables circumvention of, protection measures designed to control access to copyrighted work.⁶⁷ Both civil and criminal remedies also now exist under the DMCA if one circumvents “a technological measure that effectively controls access to a work protected” by the Copyright Act.⁶⁸ It is a civil violation and a crime to “manufacture, import, offer to the public, provide or otherwise traffic in any technology, product, service, device, component, or part thereof” that “is primarily designed or produced for the purpose of circumventing a technological measure that effectively controls access to a work protected” under the Copyright Act.⁶⁹ A technological measure effectively controls access to a work if the measure, “in the ordinary course of its operation, requires the application of information or a process or a treatment, with the authority of the copyright owner, to gain access to the work.”⁷⁰ One circumvents such technology measure if one uses a means “to descramble a scrambled work, to decrypt an encrypted work, or otherwise to avoid, bypass, remove, deactivate, or impair a technological measure,” without the authority of the copyright owner.⁷¹

In *RealNetworks, Inc. v. Streambox, Inc.*,⁷² Streambox distributed software that enabled users to bypass the authentication process employed by RealNetworks, which distributes audio and video content over the Internet. Thus, Streambox users could get the benefit of the RealNetworks streaming audio and video content without compensating the copyright owners. The United States District Court in Washington State found that the Streambox software was a technological measure that was designed to circumvent the access and copy control measures intended to protect the copyright owners.⁷³

In a case involving digital video disc (DVD) encryption, a U.S. District Court in New York enjoined posting links to sites where visitors may download the decryption program as trafficking in circumvention technology and a violation of the DMCA.⁷⁴ In *Universal City Studios, Inc. v. Reimerdes*, the court rejected an argument that the use of the decryption software constituted free expression protected by the First Amendment of the U.S. Constitution. On appeal, the appellant argued that the injunction violated the First Amendment because computer code was speech, was entitled to full protection, and was unable to survive the strict scrutiny given to protected speech.⁷⁵ The appellate court found that the computer code used in the program was protected speech:

Communication does not lose constitutional protection as “speech” simply because it is expressed in the language of computer code. Mathematical formulae and musical scores are written in “code,” i.e., symbolic notations not comprehensible to the uninitiated, and yet both are covered by the First Amendment. If someone chose to write a novel entirely in computer object code by using strings of 1’s and 0’s for each letter of each word, the resulting work would be no different for constitutional purposes than if it had been written in English. The “object code” version would be incomprehensible to readers outside the programming community (and tedious to read even for most within the community), but it would be no more incomprehensible

11 · 16 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

than a work written in Sanskrit for those unversed in that language. The undisputed evidence reveals that even pure object code can be, and often is, read and understood by experienced programmers. And source code (in any of its various levels of complexity) can be read by many more. See *Universal I*, 111 F. Supp. 2d at 326. Ultimately, however, the ease with which a work is comprehended is irrelevant to the constitutional inquiry. If computer code is distinguishable from conventional speech for First Amendment purposes, it is not because it is written in an obscure language.⁷⁶

The court then analyzed the type of scrutiny that should be applied where the restriction is content neutral:

Having concluded that computer code conveying information is “speech” within the meaning of the First Amendment, we next consider, to a limited extent, the scope of the protection that code enjoys. As the District Court recognized, *Universal I*, 111 F. Supp. 2d at 327, the scope of protection for speech generally depends on whether the restriction is imposed because of the content of the speech. Content-based restrictions are permissible only if they serve compelling state interests and do so by the least restrictive means available. See *Sable Communications of California, Inc. v. FCC*, 492 U.S. 115, 126, 106 L. Ed. 2d 93, 109 S. Ct. 2829 (1989). A content-neutral restriction is permissible if it serves a substantial governmental interest, the interest is unrelated to the suppression of free expression, and the regulation is narrowly tailored, which “in this context requires … that the means chosen do not ‘burden substantially more speech than is necessary to further the government’s legitimate interests.’” *Turner Broadcasting System, Inc. v. FCC*, 512 U.S. 622, 662, 129 L. Ed. 2d 497, 114 S. Ct. 2445 (1994) (quoting *Ward v. Rock Against Racism*, 491 U.S. 781, 799, 105 L. Ed. 2d 661, 109 S. Ct. 2746 (1989)).⁷⁷

Finding that the government’s interest in preventing unauthorized access to encrypted copyrighted material is unquestionably substantial, and that the regulation of decryption programs served that interest, the appellate court upheld the prohibitions against both posting of, and linking to, the decryption program.

Not all efforts to “circumvent” restrictions, however, come within the prohibitions of the DCMA. In *I.M.S. Inquiry Mgmt. Sys. v. Berkshire Info. Sys.*,⁷⁸ the defendant had used a valid password provided to plaintiff’s own customers and user identification to view plaintiff’s e-Basket system exactly as the customer itself might have done. The court concluded that although this might be viewed as a technology measure, it was not circumvention of a digital wall within the meaning of the DCMA.

11.6.1 Exceptions to the Prohibitions on Technology Circumvention.

The DMCA, however, explicitly carves out all defenses to copyright infringement, including the doctrine of fair use, as being unaffected by the passage of the DMCA. In some circumstances fair use can include reverse engineering.

11.6.1.1 Fair Use and Reverse Engineering. Thus, one can spy through reverse engineering still without running afoul of copyright protection or the DMCA.

However, in *Bowers v. Baystate Technologies, Inc.*,⁷⁹ a split Federal Circuit Court of Appeals found that a shrink-wrap license prohibiting reverse engineering was enforceable against the licensee who had reverse engineered Bowers’s CAD Designer’s Toolkit to develop a competing product. The *Bowers* court found that the contractual language trumped the “fair use” permitted under the Copyright Act. The Fifth Circuit has reached the opposite result in the earlier decision of *Vault Corp. v. Quaid Software, Ltd.*,⁸⁰ specifically finding that the Copyright Act preempts state law that attempts to prohibit disassembly, and holding a mass distribution license agreement unenforceable.

CIRCUMVENTING TECHNOLOGY MEASURES 11 · 17

Thus, the extent to which *Bowers* may be followed is still unclear, but it appears to be questioned in subsequent decisions.⁸¹ *Bowers* suggests a course that businesses can attempt to follow to curtail reverse engineering, which is to limit that right by contract. If *Bowers* becomes widely accepted, the United States will be in conflict with the European Union on this issue. In its 1991 Software Directive, the European Union set forth a right to reverse engineer that is consonant with “fair use” under the Copyright Act. The Software Directive also provided that the right cannot be waived by contract. So, until *Bowers* is settled, if a shrink-wrap license prohibits reverse engineering, it would be best to consider engaging in such activity abroad.

11.6.1.2 Other Exceptions. The DMCA also creates an important exception that recognizes the right to reverse engineer if (a) the person has lawfully obtained the right to use a copy of a computer program, and (b) the sole purpose of circumventing the technology measure is to identify and analyze “those elements of the program that are necessary to achieve interoperability of an independently created computer program with other programs.”⁸² The DMCA creates a similar exemption for circumvention for the purpose of “enabling the interoperability of an independently created computer program with other programs, if such means are necessary to achieve such interoperability.”⁸³ The term “interoperability” is defined to encompass the “ability of computer programs to exchange information and of such programs mutually to use the information which has been exchanged.”⁸⁴ The information acquired through these permitted acts of circumvention may also be provided to third parties so long as it is solely used for the same purposes.⁸⁵

Circumvention is permissible under these exemptions, however, “only to the extent [that it] does not constitute copyright infringement.” Two cases, *Chamberlain Group, Inc. v. Skylink Techs., Inc.*,⁸⁶ and *Lexmark Int'l, Inc. v. Static Control Components, Inc.*,⁸⁷ are particularly instructive. In both cases, the courts permitted a competitor’s access and reverse engineering under this exemption. In contrast, in *Storage Tech Corp. v. Custom Hardware Engineering Consulting, Inc.* (D. Mass. 2004), the defendant bypassed a protective access key to activate the diagnostics program by copying the code into the random access memory (RAM) of the defendant’s access device. The District Court found that this copying constituted infringement. The result was reversed in a 2 to 1 decision in the United States Federal Circuit⁸⁸ based on a reading of sections 117(a) and (c) of the DMCA, which permits copying for maintenance purposes. This string of decisions has led to recommendations that access be controlled by a method that would cause copyright infringement and that access protect not just the copyrighted program but copyrighted data so as to exclude the rationale of the Federal Circuit. Suggestions have been made that certain parts of copyrighted executable code be encrypted and that a decryption key be required that will create a copy of the code and protected data as part of the process so as to create an argument of copyright infringement. These types of recommendations remain untested, and the simpler course may be control through terms inserted into the licensing agreement.

Exempt from the DMCA, as well, are “good faith” acts of circumvention where the purpose is encryption research. A permissible act of encryption research requires that (a) the person lawfully have obtained a copy, (b) the act is necessary to the research, (c) there was a good faith effort to obtain authorization before the circumvention, and (d) such act does not constitute an infringement under a different section of the Copyright Act or under the Computer Fraud and Abuse Act of 1986. With the caveat that it must be an act of good faith encryption research, the technological means for circumvention can be provided to others who are working collaboratively on such

11 · 18 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

research. The issue of good faith encryption research looks to what happened to the information derived from the research. If it was disseminated in a manner that was likely to assist infringement, as opposed to reasonably calculated to advance the development of encryption technology, then the act still falls outside of the exemption. Other factors that go into the determination of good faith are whether the person conducting the research is trained, experienced, or engaged in the field of encryption research and whether the researcher provides the copyright owner with a copy of the findings.

The DMCA also has a bias against the collection or dissemination of personally identifying information. Thus, it is not a violation of the DMCA to circumvent a technology measure that essentially protects, collects, or disseminates personally identifying information, provided that the circumvention has no other effect, and provided that the program itself does not contain a conspicuous notice warning against the collection of such information and a means to prevent or restrict such collection.⁸⁹ In short, one can disable cookies if the program does not itself permit a user to do so.

Finally, insofar as relevant to this chapter, the DMCA also excludes from its scope “security testing.” The DMCA grants permission to engage in security testing that, but for that permission, would violate the terms of the DMCA. If the security testing, for some reason, violated some other provision of the Copyright Act or the Computer Fraud and Abuse Act of 1986, then it is still an act of infringement. The DMCA, in part, considers whether a violation occurred, and by whom the information was used. The factors to be considered include if the information was used to promote the security of the owner or operator of the computer network or system, if it was shared with the developer, and if it was used in a manner that would not facilitate infringement.⁹⁰ For purposes of the DMCA, security testing means accessing either an individual computer or network for the purpose of “good faith testing, investigating, or correcting, a security flaw or vulnerability, with the authorization of the owner or operator.”⁹¹

11.6.1.3 Remedies. The criminal penalties for violation of the DMCA can be quite severe. If the violation is willful for commercial gain, the first offense bears a fine of up to \$500,000 or 5 years’ imprisonment. Subsequent violations bear fines of up to \$1 million dollars or 10 years imprisonment. Civil remedies include an order to restrain the violation, damages for lost profits, damages for recovery of the infringer’s profits, or statutory damages for each violation. Depending on the section of the DMCA at issue, each violation can generate fines of up to \$2,500 or \$25,000. Since each act of infringement can constitute a violation, the statutory fines can become quite substantial.

11.7 PATENT PROTECTION. Ideas, which are not protected by copyright, can be protected through a patent. In general, the patent laws protect the functionality of a product or process.

11.7.1 Patent Protection Requires Disclosure. A patent can be properly obtained if the invention is new, useful, nonobvious, and disclosed. The patent exchanges a grant of an exclusive monopoly over the invention in return for disclosure. Disclosure is the trigger point for patentability. The disclosure supports the claims of patentability (i.e., it sets up the claim that the invention is both new and nonobvious) and also the scope of what can be protected. Thus, 35 U.S.C. section 112 provides:

The specification shall contain a *written description* of the invention, and of the manner and process of making and using it, in such full, clear, concise and exact terms as to *enable any person skilled in the art* to which it pertains, or with which it is most nearly connected, to make

PATENT PROTECTION 11 · 19

and use the same, and shall set forth the *best mode* contemplated by the inventor of carrying out his invention.

The specification shall conclude with one or more claims particularly pointing out and *distinctly claiming* the subject matter which the applicant regards as his invention. [Emphasis added.]

A patent therefore must disclose the best mode for implementing the invention, a clear written description of the invention, sufficient detail so that a practitioner can understand and make use of the description, and distinct claims, in order for a patent to issue.⁹² Through adequate disclosure of the invention, the application gives notice of the technology involved in the patent so as to put the public on fair notice of what would constitute an infringement. From a public policy perspective, the disclosure enlarges the public knowledge. From the inventor's perspective, the trade-off is disclosure for exclusivity. Depending on how the invention is to be used and the areas in which protection will be necessary, disclosure may not be the best means of protecting the invention. This is particularly true if the inventor is not convinced it will be deemed nonobvious from prior art, in which case it will be subject to challenge, or if, after disclosure, other companies may legally use the disclosed information for competitive advantage. The effects of disclosure should be carefully considered before applying for patent protection.

11.7.2 Patent Protection in Other Jurisdictions. Patent protection is jurisdictional. What that means, in general, is that a patent has legal meaning in the country that granted it. The United States is a signatory to the Paris Convention for the Protection of Industrial Properties, which has roughly 160 signatories. The Paris Convention essentially grants a one-year grace period for filing national patent applications in each selected signatory, to obtain the benefit of the original filing date in the United States. An alternative, open to members of the Paris Convention, is the Patent Cooperation Act. This permits the filing of an international patent that basically gives the patentee an 8- to 18-month window to test feasibility, and which simplifies the national application process.

11.7.3 Patent Infringement. Like the remedies for copyright infringement, the remedies for patent infringement include injunctive relief and damages that, by statute, are not less than a reasonable royalty for the infringing use.⁹³ If the infringement is willful, the damages can be trebled. Attorneys' fees can be awarded, but only in exceptional cases.

In the area of exported computer software, an issue of note has arisen under 35 U.S.C. section 271(f). Section 271(f) was added in 1984 to the patent law to prevent infringers from avoiding liability by finishing goods outside of the United States. An infringer will be liable if its intent is to manufacture or supply a component from the United States to be combined elsewhere, if it would be an infringement had it occurred within the United States. Exported software may be considered a "component" under section 271(f). In *Microsoft Corp. v. AT&T Corp.*,⁹⁴ the issue was whether a master disk supplied by Microsoft abroad for duplication and installation abroad of its Windows program ran afoul of AT&T's patent. In overruling the Federal Circuit, the Supreme Court concluded that it did not.

Section 271(f) prohibits the supply of components "from the United States ... in such manner as to actively induce the combination of such components." § 271(f)(1). Under this formulation, the very components supplied from the United States, and not copies thereof, trigger § 271(f).

11 · 20 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

liability when combined abroad to form the patented invention at issue. Here, as we have noted, the copies of Windows actually installed on the foreign computers were not themselves supplied from the United States. Indeed, those copies did not exist until they were generated by third parties outside the United States. Copying software abroad, all might agree, is indeed easy and inexpensive. But the same could be said of other items: “Keys or machine parts might be copied from a master; chemical or biological substances might be created by reproduction; and paper products might be made by electronic copying and printing.”... The absence of anything addressing copying in the statutory text weighs against a judicial determination that replication abroad of a master dispatched from the United States “supplies” the foreign-made copies from the United States within the intendment of § 271(f).

Unless section 271(f) is amended, it may have profound implications for subverting the ability of a U.S. company to control patent infringement where software is a component of a patented invention.

11.8 PIRACY AND OTHER INTRUSIONS. For as long as ideas and innovation have been a source of commercial or social value, the terms on which these ideas and innovations have been available for use and exchange by others has been the subject of significant tension. Although inventors and creators of commercially viable products and processes want to maximize the return on their investment, marketplace pressure for cost efficiency (often motivated by human and corporate greed) fuels a constant drive to remove the inventors’ and creators’ royalties from the cost of production. Thus, the ancient notion of piracy, the unauthorized boarding of a ship to commit theft, and the unauthorized use of another’s invention or production⁹⁵ remains alive and well. The piracy we speak of is not simply the unauthorized copying of millions of compact discs (CDs); increasingly it includes the unauthorized scraping of data from Web sites, abuse of authorized Internet use, theft of employee data, and similar activities.

11.8.1 Marketplace. The demand for unlicensed access to and use of software and entertainment media increases annually. In its 2007 survey regarding computer security among corporate and governmental institutions, the Computer Security Institute and the U.S. Federal Bureau of Investigation found that 59 percent of all respondents discovered employees who abused Internet privileges for a variety of unauthorized purposes.⁹⁶ A 2007 study by the Software & Industry Information Association reported worldwide revenue loss from the piracy (unlawful copying and distribution) of software exceeding \$28.8 billion in 2007.⁹⁷ In countries such as China, despite recent overtures to the contrary, piracy is not merely sanctioned, it constitutes an investment by government agencies.⁹⁸

In recent years, in large part due to the saturation of Internet access, there has been a tremendous proliferation of technologies designed to access and distribute (without authorization) protected software applications and entertainment media. This has posed a tremendous challenge for license holders, legislators, and law enforcement authorities. The results have included attempts to punish both unauthorized access and use of protected material. In the process, there has been a transformation in the definition of what is protected and some confusion about the extent of that protection when the Internet is involved.

11.8.2 Database Protection. Databases, the organized compilation of information in an electronic format, are prominent elements of any discussion concerning copyright protection. Compilations of information, data, and works are protectable under the Copyright Act.⁹⁹ To secure copyright protection for a compilation, a party

PIRACY AND OTHER INTRUSIONS 11 · 21

must demonstrate that (1) it owned a valid copyright in the compilation; (2) the alleged infringer copied at least a portion of the compilation; and (3) the portion so copied was protected under the Copyright Act.¹⁰⁰ In this context, the Copyright Act protects the “original” selection, coordination, or arrangement of the data contained in the compilation.¹⁰¹

To the extent that compilations contain purely factual information (e.g., existing prices of products and services), there is no protection because the facts themselves lack originality.¹⁰² It does not matter that the author “created” the facts of the prices being charged for the product or service.¹⁰³ To sustain a claim of copyright protection for compilations of fact, the author must demonstrate creativity in the arrangement of the data. Standard or routine arrangements are likewise beyond the act’s umbrella.¹⁰⁴ This is in contrast to the European Union’s Database Directive, which does not require creativity as an element for the protection of a database. Rather it protects investment in databases under copyright protection subject, however, to fair use qualifications.

The United States Supreme Court has held that the compilation into a database of original works by contributing authors to newspapers and magazines violates the copyrights of the individual authors when the database does not reproduce the authors’ articles as part of the original collective work to which the articles were contributed. In *New York Times Co., Inc. v. Tasini*,¹⁰⁵ authors who contributed articles and other works to the *New York Times*, *Time* magazine, and *Newsday* sued when they learned that the articles that they sold to the publishers for use in the respective publications were being reproduced and made available online, through LEXIS/NEXIS, an online database, and on CD-ROM. In most instances, the reproductions were of the individual articles outside of the newspaper or magazine context, in a collection of works separately protected by the Copyright Act. The Supreme Court held that, because the publishers of the new collective works made no original or creative contribution to the individual authors’ original works, they could not reproduce and distribute those works outside of the format that each publisher created for the original collections of works, without permission from, or payments to, each author.¹⁰⁶

11.8.3 Applications of Transformative and Fair Use. The concepts of transformative use and fair use (to the extent that they are separable) discussed earlier in this chapter have played a substantial role in recent decisions involving the authorized use of electronic media and the Internet. The starting point for this application of the doctrine is the U.S. Supreme Court’s decision in *Sony Corporation v. Universal City Studios, Inc.*,¹⁰⁷ the famous battle over Betamax initiated by the movie industry. At issue was whether electronic recording machines could record television programs to permit individuals to “time-shift” television programs (i.e., to record programs for viewing at a time other than the time of airing). In its decision, the *Sony* Court found that time shifting was a productive use of the television programs for a purpose other than the original commercial broadcast, and was not an attempt either to duplicate the original purpose or to impact the commercial market for these programs. The Court emphasized the noncommercial element inherent in time shifting.¹⁰⁸

11.8.4 Internet Hosting and File Distribution. The growth of the breadth and scope of the Internet has been accompanied by increasing questions about the extent to which the distribution of otherwise protected expressions change their form when converted into an electronic format. These questions arise for ISPs, which provide the pathway for distributing protected material, and for end users who post such materials

11 · 22 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

on their Web sites and bulletin boards. For ISPs, the DMCA provides some initial comfort.

Title II of the DMCA, designated the “Online Copyright Infringement Limitation Act,” establishes several infringement liability safe harbors for service providers. The “Information residing on systems or networks at direction of users”¹⁰⁹ safe harbor is available to any provider of “online services or network access, or the operator of facilities thereof, . . .” including “digital online communications, between or among points specified by user, of material of the user’s choosing, without modification to the content of the material as sent or received”¹¹⁰ that “has adopted and reasonably implemented, and informs subscribers and account holders of the service provider’s system or network of, a policy that provides for the termination in appropriate circumstances of subscribers and account holders of the service provider’s system or network who are repeat infringers” and “accommodates and does not interfere with standard technical measures.”¹¹¹ To qualify for the safe harbor, the service provider must demonstrate that:

1. It has no actual or constructive knowledge that information on its system is infringing, it is not aware of circumstances from which infringement is apparent or, upon obtaining such knowledge or awareness, it acts expeditiously to remove those materials;
2. It receives no financial benefit directly attributable to the infringing activity, *and*
3. Upon receipt of a notice of infringing material on its system, responds expeditiously to remove, or disable access to, the material.¹¹²

Assuming that the safe harbor does not apply (as, for instance, because the ISP failed to act on a notice of infringing activity), many service providers may nonetheless escape liability. In the first, and seminal, case on this topic, *Religious Technology Center v. Netcom On-Line Communication Services, Inc.*,¹¹³ an ISP hosted a bulletin board service on which Church of Scientology publications were posted by a former minister. The District Court held that the ISP must demonstrate that its use was of public benefit (facilitating dissemination of creative works including, but not limited to, the infringing work); that its financial gain was unrelated to the infringing activity (e.g., subscription fees from providing email systems rather than fees from the display or sale of the infringing work); that its use was unrelated to the use of the owner of the work; that the ISP copied only what was necessary to provide its service; and that its use of the material had no demonstrable effect on the potential market for the work.¹¹⁴ In *CoStar Group, Inc. v. LoopNet, Inc.*, the Fourth Circuit relied on *Netcom*, its codification in the DMCA, and the fact that the DMCA does not limit the application of other infringement defenses, and held that “the automatic copying, storage, and transmission of copyrighted materials, when instigated by others, does not render an ISP strictly liable for copyright infringement under §§501 and 106 of the Copyright Act.”¹¹⁵

For Web site owners and users who post allegedly infringing material, the courts have had much less difficulty discarding the transformative fair use arguments. This has been particularly true in the purely commercial setting, as where the infringing party gains direct financial benefit from the infringing material,¹¹⁶ and where the posted material is an exact copy of the protected work without any transformation to something creative or original.¹¹⁷ In a case that goes to the heart of the open-access nature of the Internet, one court recently held that a copyright owner who posts its work on the Internet for

PIRACY AND OTHER INTRUSIONS 11 · 23

free distribution as shareware may defeat a transformative fair use defense by also posting an express reservation of distribution rights.¹¹⁸

11.8.5 Web Crawlers and Fair Use. The Internet, premised on open exchange of data and economic efficiency, has spawned a spate of data search and aggregation software tools that scan the Web looking for information requested by the user. The process used by these search engines¹¹⁹ includes identifying data on the Web that conforms to the search parameters and then downloading that data. Since the copying usually occurs without the express permission of the copyright owner, some have argued that such copying constitutes an infringement. Although there is very little precedent concerning the application of transformative fair use to automated data retrieval systems, at least one court has upheld the use of the defense to an infringement claim.¹²⁰

11.8.6 HyperLinking. In *Perfect 10 v. Google, Inc.*,¹²¹ affirmed in part and remanded in part, *Perfect 10, Inc. v. Amazon.com, Inc.*,¹²² Perfect 10 (P10) claimed that Google was infringing its ownership of copyrights in certain images and thumbnails hosted by third-party and P10's Web sites when Google's image search picked them up for display as framed full-size images and as thumbnails on computers and cell phones. The court concluded that hyperlinking did not constitute display for purposes of direct copyright infringement. On appeal, the case was remanded for further consideration as to whether the conduct fell within the general rule for contributory liability. To appreciate the context in which the courts are wrestling with these issues in the light of new technology, a review of the District Court's analysis should be studied.

11.8.7 File Sharing. Transformative fair use will not protect the verbatim retransmission of protected work in a different medium when there is a substantial and detrimental impact on the market for the protected work. In *A&M Records, Inc. v. Napster, Inc.*,¹²³ Napster enabled users to share music files over the Internet by downloading the file-sharing software to their hard drive, using the software to search for MP3 music files stored on other computers, and transferring copies of MP3 files from other computers. The court of appeals held that Napster users were merely retransmitting original works in a different medium and that this did not constitute a transformation of the original work. The court also found that sharing of music files over the Internet had, and would have, a significant and detrimental impact on the existing and potential market for CDs and digital downloads of the copyright owners' works. Picking up on the *Sony* decision's emphasis on the distinction between commercial and personal use, the Court of Appeals found that Napster's Web site effectively made the works available for use by the general public and not simply for the personal use of individual users.¹²⁴

Napster's demise, however, did not end the controversy over file sharing. Trying to avoid Napster's method of directly enabling file sharing, entities such as Grokster and StreamCast developed software creating peer-to-peer networks through which individual computers communicate to exchange files without the necessity of a central server.¹²⁵ The Supreme Court recently revisited copyright infringement and file sharing specifically with respect to these peer-to-peer networks and applied the "inducement rule" to file-sharing services. Evidence demonstrated that 90 percent of the files available to download from Grokster and StreamCast were copyrighted works, and Grokster and StreamCast conceded that most users were downloading copyrighted material. There was also an abundance of evidence that through their respective software

11 · 24 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

applications and advertisements, both entities marketed themselves as the alternative to Napster, and their business models demonstrated “that their principal object[ive] was [the] use of their software to download copyrighted works.”¹²⁶ The Court vacated the court of appeals’ affirmation of summary judgment for Grokster and StreamCast, and rejected the court of appeals’ broad interpretation of *Sony Corp. v. Universal City Studios*, but declined to further discuss the balance between protecting copyrighted works and promoting commerce in the context of how much noninfringing use each service was capable of providing, and did not at all discuss the issue of fair use. Instead, the Court noted that *Sony* did not preclude other forms of infringement liability and, focusing on the intent of the defendants in their inducement of file sharing, held that “one who distributes a device with the object of promoting its use to infringe copyright, as shown by clear expression or other affirmative steps taken to foster infringement, is liable for the resulting acts of infringement by third parties.”¹²⁷ Citing *Sony*, the Court further opined that mere knowledge of potential or actual infringement are not sufficient bases for liability, but that “the inducement rule … premises liability on purposeful, culpable expression and conduct, and thus does nothing to compromise legitimate commerce or discourage innovation having a lawful purpose.”¹²⁸

Since the service and software in *Grokster* had other lawful purposes, the Supreme Court’s decision underscores the importance of proving an intent to infringe or cause infringement. Thus, when asking a court to look behind stratagems and disclaimers that hide unlawful purposes, the copyright holder should consider what other evidence exists or is likely to exist of product design, advertising, marketing, external and internal communications, revenue plans, and other factors that would prove unlawful intent. In addition, for copyright holders, the problem remains that many providers of file-sharing software may not be subject to the jurisdiction of U.S. courts and that file-sharing software, such as “Darknet,” provides anonymity to users illegally downloading copyrighted materials. As will be discussed, many countries are signatories to TRIPS (see Section 11.11.1 of this chapter) and subscribe to international copyright protection. Following *Grokster*, the maker of KaZaa file-sharing software was enjoined in Australia from using its software to commit copyright infringement. The remedy required alteration of the software so that it would not duplicate copyrighted works.

11.9 OTHER TOOLS TO PREVENT UNAUTHORIZED INTRUSIONS. Several legal principles and laws support the right to prevent and prosecute unauthorized intrusions. These include the definition of trespass, terms of use, and several critically important and widely used laws explicitly addressing the issues.

11.9.1 Trespass. Trespass is a common law concept that we are all familiar with when applied to land. We have all seen and probably at some point in our youth violated the no-trespassing signs that are posted on an unfriendly neighbor’s property. Trespass is also a concept that can apply to computers and informational databases. Courts have been taking older concepts and reapplying them to new situations.

In *eBay, Inc. v. Bidder’s Edge, Inc.*,¹²⁹ the Federal District Court granted eBay an injunction forbidding Bidder’s Edge from using a software robot to scrape information from eBay’s Web site. The court based the injunction on its finding that accessing the Web site in a manner that was beyond eBay’s posted notice (there were actual letters of objection) constituted a trespass. The court reasoned that the “electronic signals

OTHER TOOLS TO PREVENT UNAUTHORIZED INTRUSIONS 11 · 25

sent by Bidder's Edge to retrieve information from eBay's computer system [were] sufficiently tangible to support a trespass cause of action." The court further viewed the ongoing violation of eBay's fundamental right to exclude others from its computer system as creating sufficient irreparable harm to warrant an injunction. Thus, it was not necessary that eBay prove that the access actually interfered with the operation of the Web site. Rather, proof of the "intermeddling with or use of another's personal property" was sufficient to establish the cause of action for trespass. What is significant here is that eBay did permit others to access its Web site under license, and the court viewed conduct that exceeded the licensed use, upon notice to the violator, to be a trespass.

However, the applicability of trespass to unauthorized computer activity is not settled. Where trespass involves an object, rather than land, there must not only be improper use but also some harm to the physical condition or value of the object, or the misuse must deprive the rightful owner of the use of the object for a substantial period of time. The two must be causally related. In *Intel v. Hamidi*,¹³⁰ the California Supreme Court reversed a lower court's banning a former employee from sending unsolicited emails on the grounds of trespass. The court thought that the reach of the doctrine had been extended too far, concluding that bad analogies (i.e., viewing servers as houses and electronic waves as intrusions) create bad law. The court declined to view computers as real property. Rather, finding that they were like other personal property, the court found that this communication was no different from a letter delivered by mail or a telephone call. In short, the court declined to find a trespass because there was an "unwelcome communication, electronic or otherwise" that had fictitiously caused an "injury to a communication system." Here there was no injury to the computer system although Intel claimed injury to its business.

Intel v. Hamidi simply warns against overbreadth of application of the concept of trespass. If injury to the computer system can be demonstrated, then the concept of trespass does lie as a tool in the arsenal of remedies assuming that the trespasser can be identified.

11.9.2 Terms of Use. Terms of use can constitute a contract with respect to Web site usage. Thus, in any situation where electronic access is requested or permitted, the terms and conditions of use, together with an acknowledgment that such terms have been seen and consented to, can be enforced as restricting usage. In *Register.com, Inc. v. Verio, Inc.*,¹³¹ the Second Circuit upheld an order enjoining Web site access primarily on the issue of contract. There, as described by the Second Circuit, the defendant Verio, against whom the preliminary injunction was issued, was engaged in the business of selling a variety of Web site design, development, and operation services. In the sale of such services, Verio competed with Register's Web site development business. To facilitate its pursuit of customers, Verio undertook to obtain daily updates of the WHOIS information relating to newly registered domain names. To achieve this, Verio devised an automated software program, or robot, which each day would submit multiple successive WHOIS queries through the port 43 accesses of various registrars. Upon acquiring the WHOIS information of new registrants, Verio would send them marketing solicitations by email, telemarketing, and direct mail. To the extent that Verio's solicitations were sent by email, the practice was inconsistent with the terms of the restrictive legend Register attached to its responses to Verio's queries.

Register at first complained to Verio about this use and then adopted a new restrictive legend on its Web site that undertook to bar mass solicitation "via direct mail, electronic

11 · 26 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

mail, or by telephone.” The court concluded that Verio’s conduct formed a contract, like buying an apple at a roadside fruit stand, which Verio breached:

We recognize that contract offers on the Internet often require the offeree to click on an “I agree” icon. And no doubt, in many circumstances, such a statement of agreement by the offeree is essential to the formation of a contract. But not in all circumstances. While new commerce on the Internet has exposed courts to many new situations, it has not fundamentally changed the principles of contract. It is standard contract doctrine that when a benefit is offered subject to stated conditions, and the offeree makes a decision to take the benefit with knowledge of the terms of the offer, the taking constitutes an acceptance of the terms, which accordingly become binding on the offeree. See, e.g., *Restatement (Second) of Contracts* § 69 (1)(a) (1981) (“Silence and inaction operate as an acceptance … where an offeree takes the benefit of offered services with reasonable opportunity to reject them and reason to know that they were offered with the expectation of compensation.”)

* * *

Returning to the apple stand, the visitor, who sees apples offered for 50 cents apiece and takes an apple, owes 50 cents, regardless whether he did or did not say, “I agree.” The choice offered in such circumstances is to take the apple on the known terms of the offer or not to take the apple. As we see it, the defendant in Ticketmaster and Verio in this case had a similar choice. Each was offered access to information subject to terms of which they were well aware. Their choice was either to accept the offer of contract, taking the information subject to the terms of the offer, or, if the terms were not acceptable, to decline to take the benefits.

Id., at 403; and was also a trespass because:

The district court found that Verio’s use [**31] of search robots, consisting of software programs performing multiple automated successive queries, consumed a significant portion of the capacity of Register’s computer systems. While Verio’s robots alone would not incapacitate Register’s systems, the court found that if Verio were permitted to continue to access Register’s computers through such robots, it was “highly probable” that other Internet service providers would devise similar programs to access Register’s data, and that the system would be overtaxed and would crash. We cannot say these findings were unreasonable.

Id., at 405.

Similarly, although in a different setting, in *ProCD v. Zeidenberg*,¹³² where ProCD sold a CD with noncopyrightable data. Access to the data, however, was controlled by a license agreement; if there was no acceptance, there was also no access. The license agreement prohibited the use of the data for any commercial use. Zeidenberg took the data and posted it on a Web site, which he used commercially to sell advertising. Thus, the data were being used to attract visitors. The court found the license limitation on use enforceable.

The importance of this decision is that so long as the owner prominently specifies the limitations, the restrictions can become a contract that is accepted by accepting the benefits of access and can be one safeguard against misuse of the access.

11.9.3 Computer Fraud and Abuse Act¹³³

11.9.3.1 Prohibited Behavior and Damages. In 1984, Congress passed the original version of the Computer Fraud and Abuse Act (CFAA).¹³⁴ The general purpose was to protect “Federal interest computers” by criminalizing intentional and unauthorized access to those computers that resulted in damage to the computers or the data stored on them. The statute was substantially amended in 1986¹³⁵ and again in 1996¹³⁶ and now contains both criminal and private civil enforcement provisions.

OTHER TOOLS TO PREVENT UNAUTHORIZED INTRUSIONS 11 · 27

The statute proscribes these activities:

- ... knowingly accessing a computer without authority or in excess of authority, thereafter obtaining U.S. government data to which access is restricted and delivering, or attempting to deliver, the data to someone not entitled to receive it¹³⁷;
- intentionally accessing a computer without authority or in excess of authority and thereby obtaining protected consumer financial data¹³⁸;
- intentional and unauthorized access of a U.S. government computer that affects the use of the computer by or for the U.S. government¹³⁹;
- accessing a computer used in interstate commerce knowingly and with the intent to defraud and, as a result of the access, fraudulently obtaining something valued in excess of \$5,000¹⁴⁰;
- causing damage to computers used in interstate commerce by (i) knowingly transmitting a program, code, etc. that intentionally causes such damage, or (ii) intentionally accessing the computer without authority and causing such damage¹⁴¹;
- knowingly, and with the intent to defraud, trafficking in computer passwords for computers used in interstate commerce or by the U.S. government¹⁴²; and
- transmitting threats to cause damage to a protected computer with the intent to extort money or anything of value.¹⁴³

The linchpin among the relevant decisions concerning access to data under the CFAA is whether the access is “without authority” or “in excess of authority.” The factors considered by the courts include the steps taken by the owner of the information to protect against disclosure or use, the extent of the defendants’ knowledge regarding their authority to access or use the data, and the use(s) made of the data after gaining access. The legislative history indicates that the statute was intended to “punish those who illegally use computers for commercial advantage.”¹⁴⁴

Broadly speaking, there are two sets of circumstances to consider. In the first instance, is the actual access authorized, either expressly or impliedly? In the Internet context, where there is a presumption of open access, the site or data owners must show that they took steps to protect the contents of their site and to limit access to the data at issue.¹⁴⁵ Once those steps are taken, the protection constitutes a wall through which even automated search retrieval systems may not go without express permission.¹⁴⁶ Without the wall, there must be some evidence of an intent to access for an impermissible purpose, as when Intuit inserted cookies into the hard drives of home computers.¹⁴⁷

Second, has the authorized access been improperly exceeded? Generally speaking, those who use their permitted access for an unauthorized purpose to the detriment of the site or data owner have violated the CFAA. Examples include employees who obtain trade secret information and transmit it via the employer’s email system to a competitor for which the employee is about to begin work¹⁴⁸; using an ISP subscription membership to gain access to and harvest email addresses of other subscribers in order to transmit unsolicited bulk emails¹⁴⁹; and using access to an employer’s email system to alter and delete company files.¹⁵⁰

The criminal penalties range from fines to imprisonment for up to 20 years for multiple offenses. As discussed in Section 11.9, the CFAA has become a prominent element of claims by the U.S. government and private parties seeking to protect data that are not always protected by other statutory schemes.

11.9.3.2 Its Application to Web Crawling and Bots. Web robots, or “bots,” have become widespread to scrape data from Web sites. All of that data generally are available to the public. That is, any individual can access the same information,

11 · 28 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

but not with the speed or accuracy of a Webspider. But when does such “scraping” run afoul of the CFAA? To what extent does the law protect site operators or company data from penetration by an outside third party?

The key to the analysis under the CFAA is to ask whether the data are in fact publicly available. Are there technical barriers, such as passwords or codes, that have to be circumvented? Do the terms of use prohibit access or use other than by an individual consumer? These questions are critical to determining whether the access either exceeds authority or is without authority under the CFAA.

If the answer to either one of these questions (or similar questions) is yes, one needs to consider access carefully since such access and downloading of data is likely to violate the CFAA. In *EF Cultural Travel v. Zefer Corporation*, Zefer designed a Web bot to scrape travel trip and pricing information from the Web site of EF Cultural Travel (EF) for use by a competitive travel Web site. The bot, designed by Zefer, downloaded the information by calling URLs on which each separate trip and pricing information was stored, reading the source code for the key features, and storing the information on a spreadsheet. The bot did so in a fashion not to burden or interfere with EF’s Web site. Once gathered, the information was turned over to a competitor, who used the information to adjust price and trip information that it offered. Zefer’s scraping did not occur continuously, but only on two dedicated occasions. EF sued, claiming that a violation of the CFAA had occurred. The First Circuit Court of Appeals disagreed, refusing to read into what is or is not authorized some “reasonable expectations” standard, instead requiring that the Web site operator expressly state any limitations on access in its terms and conditions. On remand to the Federal District Court, the court, following the First Circuit, granted summary judgment for Zefer.

11.9.3.3 Simple Preventive Measures. Not surprisingly, there are several methods for preventing unauthorized access in the first instance and, if unsuccessful, in prevailing in any subsequent claim arising under the CFAA. Perhaps the most obvious measure, and one that the First Circuit Court of Appeals underscored, is to make sure that each visitor to a Web site is adequately notified that the owner of the site intends only limited use or access to the data on the site. The notice can take many forms.

For example, a detectable message easily identifiable on a home page warning visitors that the posted information is available only for viewing and not for use in any manner adverse to the host’s interests would be sufficient. Understandably, most Web hosts are reluctant to post such a blatant limitation—it is not necessarily “good for business.” For those interested in an equally effective but less direct message, an increasingly common practice is to compel site visitors to register before gaining access to links and other pages available through the home page. The more difficult the registration process, the greater the host’s apparent intent to restrict access to, and use of, the information that will be accessible after registration is completed.

Those hosts that require the payment of money, some kind of membership, or an access agreement before providing access establish what, for purposes of statutes like the CFAA that criminalize unauthorized access, will most often be seen as providing sufficient notice of the limits of authorized access. In the case of membership sites, the presumption is that each registrant is prequalified and therefore authorized to view and use the more restricted data, at least for purposes consistent with the terms of access. Enforceable click-wrap access agreements establish not only notice of access limitations; they also secure each visitor’s agreement to use the Web site and the data therein within the stated limitations.

OTHER TOOLS TO PREVENT UNAUTHORIZED INTRUSIONS 11 · 29

Securing Web-based data against unauthorized use or users is, in some ways, antithetical to the information-sharing intent and purpose of the Web. In this regard, however, the question arising when we post information on the Web differs little from the question posed over the centuries regarding the extent to which each of us wants our competitors or adversaries to use our proprietary work against our interests. The greater the concern, the more likely that each host will have to limit the data posted on the Web, or else increase each visitor's awareness of the rules of access.

11.9.4 Electronic Communications and Privacy. Electronic privacy is becoming the issue in our society of databases and networking. Most of the U.S. "privacy" statutes are subject matter specific: the Telephone Consumer Protection Act of 1991 (do not call, for telemarketers); Health Insurance Portability Accountability Act of 1996 (privacy with respect to uses and disclosure of medical information); Children's Online Privacy Protection Act of 1998 (regulating collection of information from children under the age of 13 by Web sites directed to children); Gramm-Leach-Bliley Act of 1999 (regulating sharing of customer data by financial institutions); Controlling the Assault of Non-Solicited Pornography and Marketing Act of 2003 (restricting spammers and requiring an ability to opt out); the Fair and Accurate Credit Transaction Act of 2003 (providing very limited assistance with respect to identity theft such as the obligation to provide a yearly credit report). These laws do not provide assurances of privacy in the same way that the European Union did in its 1996 Data Protection Directive.¹⁵¹ The EU Data Directive establishes protections against release of personal data, including emails, within the European Union, and restrictions on the transmission of such data outside the EU to countries or companies that do not have equivalent protections in place.¹⁵²

In 2005, ChoicePoint, a large data broker, admitted that it had sold personal data on over 160,000 people to phony companies established by identity thieves. Since then, other companies have announced data break-ins and data leaks. As a result of such data security breaches, approximately half of the states have passed laws that require disclosure of unauthorized access to personal data.¹⁵³

In the United States, the primary protection for privacy remains a lawsuit for tortious invasion of one's privacy. Because those rights are defined state by state, a review is beyond the scope of this chapter. However, most states recognize some form of the tort of invasion of privacy, and the tort has been recognized in the *Restatement (Second) of Torts* § 652, which courts reference as an authoritative source of the law. In general, the Restatement makes actionable (a) intentional intrusion, that is highly offensive to a reasonable man, into the seclusion of another's private affairs, (b) the public disclosure of private facts if such disclosure is highly offensive to a reasonable person, and is not a legitimate public concern, and (c) the appropriation for his own use or benefit of the name or likeness of another.

This chapter has already discussed the fiduciary obligation owed by employees to their employers with respect to confidential information. The development of the tort of privacy suggests that companies owe a similar obligation to their employees. Although slightly different in scope, but foreshadowing the growing body of law in this area, in *Remsburg v. Docusearch, Inc.*,¹⁵⁴ the New Hampshire Supreme Court was faced with a database company that had supplied information to a client that included a woman's personal information. The client used it to confront her and kill her. The New Hampshire Supreme Court held that the company had to act with "reasonable care in

11 · 30 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

disclosing a third person's personal information to a client." This decision is as yet an unanswered invitation to other courts.

On the federal level, the CFAA, of course, does address "unauthorized" access to computerized information. In addition, Congress has enacted some statutory regulations that specifically address electronic communications and privacy.

11.9.4.1 Wiretap Act and Electronic Communications Privacy Act.

The Omnibus Crime Control and Safe Streets Act of 1968, generally referred to as the Federal Wiretap Act,¹⁵⁵ established the general parameters for permitted interception of communications by law enforcement. As originally crafted, the Wiretap Act covered only "wire and oral communications." In 1986, Congress enacted the Electronic Communications Privacy Act (ECPA),¹⁵⁶ which amended the Wiretap Act and created the Stored Wire and Electronic Communications and Transactional Records Act (Stored Communications Act or SCA) to "update and clarify federal privacy protections and standards in light of changes in computers and telecommunication technologies."¹⁵⁷ The SCA makes it unlawful to knowingly access a prohibited electronic communications service facility without authority, or in excess of authority, and for such public service provider to disclose information contained in such facilities. The ECPA allows a private plaintiff to bring a claim for knowing or intentional violation of the statute to recover actual damages or the statutory minimum of \$1,000.

The 1986 amendment extended the Wiretap Act's coverage to include "electronic communications," which is defined as "any transfer of signs, signals, writing, images, sounds, data, or intelligence of any nature transmitted in whole or in part by a wire, radio, electromagnetic, photo-electronic or photo-optical system."¹⁵⁸ "Intercept" is defined as "the aural or other acquisition of the contents of any wire, electronic, or oral communication through the use of any electronic, mechanical, or other device."¹⁵⁹ Consequently, the Wiretap Act now makes it an offense to "intentionally *intercept* ... any wire, oral, or *electronic communication*."¹⁶⁰ Thus, the definitions in the act now cover Internet transmissions such as emails or file transfers.

There is an important exception to this prohibition. Under the "consent of a party" exception, it is permissible to intercept communications where "one of the parties to the communication has given prior consent to such interception."¹⁶¹ The requisite consent may be express or implied from the surrounding circumstances.¹⁶² Furthermore, an employer may obtain consent by informing the employee of the monitoring practices in an employment contract or in an employee handbook.¹⁶³

Under the "provider exception," a provider of electronic communication services "whose facilities are used in the transmission of a wire or electronic communication, [may] intercept, disclose or use that communication in the normal course of his employment while engaged in any activity which is a necessary incident ... to the protection of the rights or property of the provider of that service."¹⁶⁴ This exception may allow an employer to lawfully intercept communications to detect an employee's unauthorized disclosure of trade secrets to third parties.¹⁶⁵

11.9.4.2 Contemporaneous Transmission Requirement.

The Wiretap Act only prohibits *interceptions* of electronic communications,¹⁶⁶ a term that has been more narrowly defined by the courts than the definition in the act might suggest. The definition of interception provides that an individual "intercepts" a wire, oral, or electronic communication "merely by *acquiring* its contents, regardless of when or under what circumstances the acquisition occurs."¹⁶⁷ In the context of this section, a serious question arises about the legality of intercepting electronic communications as they

OTHER TOOLS TO PREVENT UNAUTHORIZED INTRUSIONS 11 · 31

were being transmitted and once they were stored, either temporarily or permanently. Although Congress intended to liberalize one's ability to monitor "wire communications" while it sought to make the monitoring of "electronic communications" more difficult,¹⁶⁸ courts have held that Congress intended to make acquisitions of electronic communications unlawful under the Wiretap Act "*only if* they occur *contemporaneously with* their transmissions"¹⁶⁹ and before they actually cross the finish line and become stored.¹⁷⁰ This is, of course, an interesting fiction when applied to Internet transmissions, which consist of packages that are broken up and passed from router to router as well as from temporary storage to temporary storage. It is a far cry from the interception of a telephone call. It may simply be that in applying the language of the statute, the courts are faced with applying it to a technology that was not really in existence when the statute was amended in 1986.

In recent years, the courts have attempted to apply the contemporaneous transmission requirement to various situations. For example, cookies used to recover personal data from visitors to a Web site constitute an interception of a contemporaneous electronic communication and a violation of the Wiretap Act.¹⁷¹ Noting that electronic communications are generally in transit and in storage simultaneously, the court reasoned that users communicated simultaneously with the pharmaceutical client's Web server and with the software company's Web server and, thus, the information was acquired contemporaneously with its transmission.¹⁷²

Where electronic transmissions are found in RAM or on the hard drive, they are stored communications and can be retrieved because they are outside of the Wiretap Act.¹⁷³ Similarly, an email that is recovered after it has been sent and received does not satisfy the contemporaneous transmission requirement and therefore has not been intercepted under the Wiretap Act.¹⁷⁴ Perhaps in response to these and other decisions, in 2001 Congress amended the Wiretap Act to apply the contemporary transmission requirement to wire communications that could not be retrieved, thereby permitting the recovery of stored wire communications.¹⁷⁵

11.9.4.3 Konop v. Hawaiian Airlines, Inc. The *Konop* decision appears to be the most oft-cited case on the issue of "interception" under the Wiretap Act. Konop, the plaintiff, was an airline pilot who created and maintained a Web site where he posted bulletins critical of his employer, Hawaiian Airlines, Inc., and the airline union. Konop controlled access to his Web site by requiring visitors to log in with a user name and password and by creating a list of authorized users.

An officer of Hawaiian Airlines asked one such authorized user for permission to use his name to access the Web site. The officer logged on several times, and another officer, using the same technique, also logged on to view the information posted on Konop's bulletin. Konop eventually filed suit against Hawaiian Airlines, alleging that it violated the Wiretap Act when its officer gained unauthorized access to Konop's Web site.

The court first reiterated that the act only prohibits *interceptions* of electronic communications.¹⁷⁶ "Interception," the court held, requires that the party acquire the information contemporaneous with its transmission, and not while it is in electronic storage. In this case, the court concluded that the employer did not violate the Wiretap Act because the officers accessed an electronic communication located on an idle Web site, which did not satisfy the contemporaneous transmission requirement.¹⁷⁷

11.9.5 Stored Communications Act. Unlike the Wiretap Act, the Stored Communications Act (SCA),¹⁷⁸ as its name suggests, establishes the limitations of access to stored communications (i.e., communications accessed *after* their

11 · 32 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

transmission).¹⁷⁹ Specifically, the SCA makes it unlawful to “intentionally [access] without authorization a facility through which an electronic communication service is provided … and thereby [obtain], [alter], or [prevent], authorized access to a wire or electronic communication while it is in electronic storage.”¹⁸⁰ The SCA defines “electronic storage” as “(A) any temporary, intermediate storage of a wire or electronic communication incidental to the electronic transmission thereof; and (B) any storage of such communication by an electronic communication service provider for purposes of backup protection of such communication.”¹⁸¹ The SCA exempts from liability conduct “authorized … by the person or entity providing a wire or electronic communications service”¹⁸² or “by a user of that service with respect to a communication of or intended for that user.”¹⁸³

11.9.5.1 Electronic Storage: Backup Files. The essential element that separates the SCA from the Wiretap Act is that the accessed communications reside in electronic storage. Therefore, the first question is what constitutes electronic storage. In *Theofel v. Fary-Jones*,¹⁸⁴ the United States Court of Appeals for the Ninth Circuit attempted to answer this question.

In *Theofel*, overzealous lawyers for Fary-Jones secured, through a subpoena issued to an ISP, emails sent and received by their opponents in the lawsuit, a company called Integrated Capital Associates (ICA). The subpoena requested from the ISP virtually every email ever sent or received by ICA and its employees. In response, the ISP posted a smattering of the emails on a Web site accessible to Fary-Jones and its lawyers. When ICA learned of these activities, it sued Fary-Jones for, among other things, violation of the SCA.

According to the court in *Theofel*, Congress recognized that users of ISPs have a legitimate interest in protecting the confidentiality of communications in electronic storage at a communications facility. Moreover, this legitimate interest cannot be overcome by fraud or by someone who knowingly exploits a mistake that permits access to what is otherwise protected. The court found that the use of the subpoena to access ICA’s emails when it was reasonably plain, at least to counsel, that the subpoena was invalid, negated any apparent authority that Fary-Jones and its lawyers may have had to view ICA’s emails.

Fary-Jones claimed that the ICA emails were not in “electronic storage” and therefore no violation of the SCA occurred. The court disagreed. As stated earlier, electronic storage exists when messages are stored on a temporary, intermediate basis as part of the process of transmitting the message to the recipient, and when messages are stored as part of a backup process. In this instance, the court found that the emails, which had apparently been delivered to their recipients, were stored by the ISP as part of its backup process for retrieval after initial receipt. Access to those emails was therefore protected by the SCA, which Fary-Jones and its lawyers violated.

11.9.5.2 Electronic Storage: Temporarily Stored Communications.

Recent cases interpreting the meaning of “temporary, intermediate storage … incidental to” transmission of the communication have adhered to the letter of the law more than its spirit. In two cases involving the installation of cookies that were subsequently accessed by software companies for commercial gain, the courts have held that cookies are permanently (or at least indefinitely) installed in the consumer’s hard drive and therefore cannot be considered “temporary, intermediate storage.”¹⁸⁵ The *DoubleClick* decision also emphasized that the “temporary, intermediate storage”

OPEN SOURCE 11 · 33

element of the SCA means what it says, that is, the prohibited conduct involves only the unauthorized access to communications while they are being temporarily stored by an intermediate and does not include access to stored messages after they have been received.¹⁸⁶ In the context of an employer’s right to examine an employee’s emails, the employee will have no claim that an employer has violated the SCA when the employer opens emails sent or received by the employee once the email has been either received or discarded.¹⁸⁷

11.10 OPEN SOURCE. With the continued proliferation of the Internet and computer software, the licensing, distribution, and use of open source code has gained publicity and added importance in the practice of intellectual property and computer security. “Open source” describes the distribution of computer code that is available (i.e., open) to all others and therefore allows computer programmers to read, apply, and modify the code, and also redistribute any changes.¹⁸⁸ The open source movement began with Richard Stallman’s development of Gnu’s Not UNIX (GNU), a freeware form of UNIX that was meant to be free software (free as in the freedom to use, modify, and distribute the software).¹⁸⁹ GNU’s development created the first open source license, the General Public License (GPL). Linux, an open source–based operating system and an alternative to Microsoft Windows, experienced tremendous growth through its use of the GPL.¹⁹⁰ The prevalence of open source issues is evidenced by the 1998 formation of the Open Source Initiative (OSI), which not only promotes open source development and encourages its use by business¹⁹¹ but also offers links to and information about most of the available open source licenses.

11.10.1 Open Source Licenses. The author of an open source code holds a copyright that operates as other copyrights do, but the code is released under a certain license on a nonproprietary basis. There are various types of open source licenses. The first open source license was the GPL, as described. It offers the broadest application of free software. In contrast, other licenses do not seek to perpetuate the free nature of a particular program. According to the Open Source Initiative, there are nearly 60 open source licenses now available for authors of source code,¹⁹² all of which assert certain requirements of the software user.

11.10.2 GPL. Licensing under the GPL is premised on Stallman’s idea of “copyleft,” which basically uses copyright as a tool to ensure the continued free distribution of source code.¹⁹³ In other words, the GPL affords application, modification, and distribution rights to the copyrighted source code only if the user agrees that the distribution terms remain the same. This creates an endless chain of GPLs attached to future distributions of either the original or derived versions, regardless of their form.¹⁹⁴ This endless chain often is referred to as the GPL’s “viral effect,” as GPL-protected code multiplies from any modifications of original GPL-protected code.¹⁹⁵ The GPL applies not just to an originally protected software program but also to what it broadly defines as the “Program”:

[A]ny such program or work, and a “work based on the Program” means either the Program or any derivative work under copyright law: that is to say, a work containing the Program or a portion of it, either verbatim or with modifications and/or translated into another language.¹⁹⁶

Moreover, although the GPL also states that independent and separate sections of a derivative work are not subject to the GPL’s terms when they are distributed as separate

11 · 34 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

works, the GPL does apply when the user distributes those same independent and separate “sections as part of a whole which is based on the Program. . . .”¹⁹⁷ The broad application given to the program under the GPL further enhances the viral effect of the license.

Other provisions of the GPL require users who distribute verbatim copies of the source code to publish copyright notices, disclaim warranties, and provide copies of the GPL. In addition, the modifier/user must attach to any modifications a notice that the software was changed, must distribute or license the software free of charge to third parties, and must provide appropriate copyright notices, warranty disclaimers, and GPL terms and conditions. In sum, the GPL’s sweeping terms not only seek to achieve the free software goals of the FSF but also to impact whether authors chose the GPL, and whether businesses utilize software subject to the GPL.

11.10.3 Other Open Source Licenses. The Berkeley Software Distribution (BSD) License and the Massachusetts Institute of Technology (MIT) License are very similar in that they both require copyright notices, disclaimers of warranties, and liability limitations. The BSD further prohibits contributors or similar organizations from endorsing the program and also requires a copy of the BSD’s terms to be distributed with the software.

11.10.4 Business Policies with Respect to Open Source Licenses. The issue of whether distribution of a proprietary work that incorporates a small portion of GPL-protected code subjects that proprietary work to the terms of the GPL has never been litigated.¹⁹⁸ This is one risk of using open source software. Another risk is that failure to comply with the GPL’s terms could lead to litigation.¹⁹⁹ For instance, MySQL sought to enjoin Progress Software Corporation from distributing MySQL’s Gemini program without a GPL-compliant agreement.²⁰⁰ Because there was a factual dispute as to whether Gemini was a derivative work or an independent work under the GPL, and because Progress stipulated that it disclosed Gemini’s source code and would withdraw the end user license for commercial users, the court did not grant the injunction as to the GPL.²⁰¹

Given the expanding use of open source, businesses need to develop comprehensive policies addressing their use of open source to avoid liability and publicly releasing their own proprietary technology.²⁰² Concerns generally involve license requirements regarding the distribution of the software and its modifications,²⁰³ since those activities usually require the company to release the source code for any distributed modification, and modifications often terminate vendors’ support agreements.²⁰⁴ In addition, distributing unmodified open source as part of a proprietary program may require the company to release its own proprietary open source code.²⁰⁵ It is more likely, however, that the company would be enjoined from distributing the open source or would have to pay damages.²⁰⁶ These considerations should be addressed not only through company policy but also by choosing the best source code to use in programming, given the company’s internal and external needs and the specific licensing requirements of that source code.

11.11 APPLICATION INTERNATIONALLY. Because the laws of the United States are the laws of just one nation among many, the enforcement of U.S. law and the protection of intellectual property rights in large part depend on international treaties. To the extent that the infringing acts or acts of piracy may be deemed to occur in the United States, or the infringers can be found in the United States, then the United States

APPLICATION INTERNATIONALLY 11 · 35

has sufficient jurisdiction over these acts to enforce its laws. In other words, such actors can be sued directly in the courts of the United States for violation of the laws of the United States.

Apart from direct enforcement, international protection is usually a vehicle of bilateral agreements between the United States and individual countries or a function of international protocols or treaties to which the United States is a signatory. Thus, for example, the Paris Convention for the Protection of Industrial Property²⁰⁷ establishes a system for recognizing priority of invention, but only among member countries. In addition, there is the Patent Cooperation Treaty (PCT), a multilateral treaty with more than 50 signatories. The PCT permits the filing of an international application that simplifies the filing process when a patent is sought in more than one nation. For copyright protection, there is also a series of international treaties and agreements that include the Berne Convention,²⁰⁸ the Universal Copyright Convention, and the World Trade Organization (WTO) Agreement.²⁰⁹ Canada, Mexico, and the United States also signed the North American Free Trade Agreement (NAFTA) in December 1992. NAFTA addresses intellectual property and requires that member states afford the same protections to intellectual property as members of the General Agreement on Tariffs and Trade (GATT). At a minimum, members of GATT must adopt four international conventions, including the Paris Convention and the Berne Convention.

These agreements, conventions, and treaties in large part do not attempt to reconcile the differences in the national laws of intellectual property. The particular national rules and nuances are simply too complicated, and there are too many differences of opinion to expect that these differences could be internally reconciled. Rather, in large measure, these international accords attempt to codify comity between the member nations so that each will recognize the legitimacy of the intellectual property rights in the other.

11.11.1 Agreement on Trade-Related Aspects of Intellectual Property Rights.

On December 8, 1994, the Agreement on Trade-Related Aspects of Intellectual Property Rights (TRIPS) was signed into law in this country. The signing of TRIPS required changes to be made in United States statutes and regulations to bring them into conformity with international norms. TRIPS, however, was a product of the United States and other industrial countries pressing for stronger, more uniform standards for international treaties concerning intellectual property. The basic structure of TRIPS is to set the minimum standard of protection for intellectual property with each member nation free to adopt more stringent standards. Under the rubric used in the United States, TRIPS applies to copyrights, patents, trademarks, service marks, mask works (integrated circuit designs), and trade secrets. It also covers geographical indications²¹⁰ and industrial designs.²¹¹ Not addressed by TRIPS, although part of the international jargon for intellectual property, are breeder's rights²¹² and utility models.²¹³ Thus, TRIPS establishes no standards as applied to these concepts, leaving each nation to set the parameters of protection unimpeded by TRIPS.

It is not by accident that TRIPS was negotiated within the context of GATT, which had set the international standards for trade tariffs and had provided remedies of trade retaliation if such standards were not adhered to. The structure of GATT provided the means under which developing countries agreed to reduce their trade tariffs in exchange for the right to export innovative products under an exclusive monopoly conveyed by intellectual property rights. The second benefit to the GATT format was to provide a means for trade retaliation if, under the dispute resolution provisions of TRIPS, the WTO determines that there is noncompliance. In reality, it is obvious that TRIPS

11 · 36 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

benefits those industrial nations that are more likely to be at the forefront of innovation and more concerned with the protection of their citizens' intellectual property.²¹⁴ The major concession wrung by the developing countries under TRIPS was obtaining a period of 4 to 11 years to implement TRIPS and to bring their national laws into conformity.

TRIPS generally reflects the U.S. view that focuses on the economic underpinnings for intellectual property rights as serving the greater societal interests. There is thus a shift from "societal" interests to "enterprise" interests. In particular, TRIPS adopts high minimum standards for patents, which will require significant legislative changes in developing countries. The copyright section, however, affords less protection than may be afforded by European nations, but it is in line with treatment in the United States. In short, TRIPS responds to the concern of enterprises in the United States that too loose a system of international protection has enabled imitation of U.S. innovations through copying and outright piracy.

11.11.2 TRIPS and Trade Secrets. Under its category for "Protection of Undisclosed Information," TRIPS provides protection for the type of information routinely referred to as trade secrets in the United States. Member nations are required to implement laws that safeguard lawfully possessed information from being disclosed to, acquired by, or used by others without consent and contrary to "honest commercial practices" if such information is (a) a secret in that it is not in the public domain, (b) has commercial value because it is a secret, and (c) has been subject to reasonable steps to maintain its secrecy.

Because discussions that led to TRIPS are not institutionally preserved, unlike the United States Congressional Record, there is no negotiating history to be consulted to flesh out the meaning of the spare paragraphs instituting trade secret protection. There do, however, appear to be differences from the total panoply of protections afforded in the United States. The concept of public domain articulated by TRIPS is information that is "not, as a body or in the precise configuration and assembly of its components, generally known among or readily accessible to persons within the circles that normally deal with the kind of information in question." This articulation appears to be addressing technological formulations of information, as opposed to general commercial information, such as financial information, that is generally considered proprietary and confidential in the United States. The focus on a technology formulation for protected information is bolstered by the TRIPS requirement that the information have commercial value. Thus, other types of information that are not part of a traded article may be deemed to have no commercial value and therefore to fall outside of the scope of protection. Depending on the particular jurisdiction in the United States, there is a distinction between confidential information and trade secrets based on the requirement that a trade secret must have commercial value. This, in turn, has been held to mean that information that is not exploited commercially is unprotectable under the law of trade secret. For example, the results of failed experiments that never resulted in a commercial product lack commercial value, even though such experiments are certainly helpful in the next round of exploration, in that they are signposts of what not to do.

The lesson to be drawn is that one should not assume symmetry of protections just because of the TRIPS provision. Instead, as part of the reasonable steps to maintain secrecy, enterprises need to consider carefully thought out and structured contractual provisions as well as a system of data caching that leaves truly confidential data in the United States, even if access is permitted outside. Improper takings of such data are,

APPLICATION INTERNATIONALLY 11 · 37

arguably, acts that occur in the United States, and such acts are subject to enforcement and punishment under the laws of the United States.

11.11.3 TRIPS and Copyright. TRIPS embraces the U.S. general model for copyright protection in its opening statement that “[c]opyright protection shall extend to expressions and not to ideas, procedures, methods of operation or mathematical concepts as such.” All member nations agree that, as to the protection of copyrights, the Berne Convention will apply. Under the Berne Convention, the duration of a copyright is the life of the author plus 50 years. If the life of a natural person is not involved, then it is ordinarily 50 years from publication. In addition, computer programs, whether in source or object code, are to be protected as literary works under the Berne Convention. TRIPS also recognizes that compilations of data can be protected as creative works. Article 10, ¶ 2 explicitly provides:

Compilations of data or other material, whether in machine readable or other form, which by reason of the *selection* or *arrangement* of their contents constitute *intellectual creations* shall be protected as such. Such protection, which shall not extend to the data or material itself, shall be without prejudice to any copyright subsisting in the data or material itself. (Emphasis added.)

TRIPS, therefore, does establish some minimum standard in the growing debate over what protections will be afforded a database. In the United States, the clear demarcation point for unprotected information is compilations that represent no more than “sweat-of-the-brow” efforts. Such compilations cannot be copyrighted.²¹⁵ The classic example of a sweat-of-the-brow effort is the copying and alphabetical organizing of names, addresses, and telephone numbers that are in telephone books. In the United States, the key for copyright protection is the creator’s original contribution of selection and arrangement. Thus, arguably, the TRIPS provision mimics the law of the United States.

The European Union (EU) has taken a more protective path. In its 1996 European DataBase Directive, the EU granted databases *sui generis* protection as their own unique form of intellectual property. Under the EU Directive, a database is “a collection of independent works, data or other materials arranged in a systematic or methodical way and individually accessible by electronic or other means.” A database may be protected either because it represents a work of “intellectual creation” or because it was compiled through “substantial investment.” The EU Directive protects such databases from unauthorized extraction or use for a period of 15 years, with the ability to extend the period for an additional 15 years if there was a “substantial new investment” in the database. Such protection extends to databases of EU members and to databases of nationals of other countries that offer protections similar to the EU.

The United States, despite a number of legislative proposals, has not adopted a concomitant rule. The result, at least for multinationals, is that entities that rely on databases should consider “locating” such databases within an EU member to take advantage of the EU’s database protections.

11.11.4 TRIPS and Patents. TRIPS requires that all members recognize the right to patent products or processes in all fields of technology. A patentable invention must be new, inventive, and have an industrial application. The patent application must fully and clearly disclose the invention so that a person skilled in the art could carry out the invention. The best mode for carrying out the invention as of the filing date must also be disclosed. Patent rights are to be enforced without discrimination as to

11 · 38 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

place of invention or whether the product is imported or produced locally. The patent of a product conveys the exclusive right to prevent, without consent of the inventor, the making, using, offering for sale, selling, or importing of the product. The patent of a process conveys the exclusive right to prevent all of the above for products that result from the process as well as the use of the process itself. The holder of a patent also has the rights to assign, transfer, or license the patent. The minimum period for a protecting a patent is 20 years from filing.

TRIPS gives each member state the right to carve out from patentability certain subject matters that have as their purpose the protection of human, animal, or plant life, or to avoid serious prejudice to the environment. In addition, TRIPS permits a member state to allow other use without authorization from the patent holder. The section defining when such use is permissible is the most detailed section among the patent provisions of TRIPS. In general, it permits such use only (a) after an effort to obtain a license from the patent holder on reasonable commercial terms and conditions, (b) with adequate remuneration to the patent holder, (c) if such use is limited predominantly to the domestic market of the member nation, and (d) if there is a review of the decision to permit, as well as the compensation, by a “higher authority in that Member.”

One of the circumstances envisioned by TRIPS is the grant of a second patent that cannot be exploited without infringing an earlier (first) patent. In such cases, a member nation may grant authority if the invention embodied in the second patent represents an “important technical advance of considerable economic significance” with respect to the first patent’s invention and a cross-license on reasonable terms is granted to the holder of the first patent to use the second patent. For process patents, TRIPS creates a limited burden on the alleged infringer to prove that the identical product was produced using a different process. In particular, a member state can create a presumption that the process patent was violated in circumstances where the product is new, or where the patent holder is unable to demonstrate what process was actually used.

11.11.5 TRIPS and Anticompetitive Restrictions. TRIPS acknowledges that some licensing practices or other conditions with respect to intellectual property rights may restrain competition, adversely affect trade, and impede the transfer and dissemination of technology. Accordingly, TRIPS permits member nations to specify practices that constitute an abuse of intellectual property rights and to adopt measures to control or limit such practices, so long as the regulation is consistent with other provisions of TRIPS. In the event that a national of a member nation violates another member’s laws and regulations regarding anticompetitive activity, TRIPS provides for the right of the involved nations to exchange information confidentially regarding the nationals and their activities.

11.11.6 Remedies and Enforcement Mechanisms. Each member nation is expected to provide an enforcement mechanism under its national laws to permit effective action against any act of infringement. Such procedures are to include remedies to prevent acts of infringement as well as to deter future acts. TRIPS imposes the obligation that all such procedures be “fair and equitable” and not be “unnecessarily complicated or costly” or involve “unwarranted delays.”²¹⁶ In general, these remedies mean access to civil judicial procedures with evidentiary standards that shift the burden of going forward to the claimed infringer, once the rights holder has presented reasonably available evidence to support its claim. Damages may be awarded sufficient to compensate the rights holder for the infringement if the “infringer knew or had reasonable grounds to know that he was engaging in infringing activity.” This means

RECENT DEVELOPMENTS IN INTELLECTUAL PROPERTY LAW 11 · 39

that vigilance and notice are essential to have meaningful protection for intellectual property rights, since notice is the best means for setting up a damage claim. TRIPS permits its members to allow the recovery of lost profits or predetermined (statutory) damages even when the infringer did not know that it was engaged in infringing behavior. Although injunctive relief is to be provided for, remedies may be limited in circumstances involving patent holders, as discussed, where adequate compensation is paid, and the alleged infringer has otherwise complied with the provisions of its national law permitting such use upon payment of reasonable compensation. In order to deter further infringement, infringing materials may be ordered destroyed or noncommercially disposed of.

In addition to civil remedies, TRIPS requires criminal penalties in cases of “willful trademark counterfeiting or copyright piracy on a commercial scale.”²¹⁷

11.12 RECENT DEVELOPMENTS IN INTELLECTUAL PROPERTY LAW²¹⁸

11.12.1 AIA. Peter E. Heuser of Schwabe, Williamson, & Wyatt summarized the Leahy-Smith America Invents Act (AIA) of 2011²¹⁹ as follows: “The AIA is the most important legislative patent reform in over 50 years. The AIA will change how patents are granted, how patent litigation will proceed and what kinds of inventions are eligible for patents, among other things.”²²⁰ The author summarized the main features of the AIA in detailed discussions of the following areas:

- First-to-file will now establish priority of invention
- Prior commercial user defense is established
- New post-grant proceedings for patent validity challenges
- The Patent and Trademark Office (PTO) will no longer grant patents on tax strategy
- Special transitional review for certain patents related to financial products and services
- Most PTO fees will increase by 15 percent
- Limited prioritized examination will be available
- New rules will affect litigation by nonpracticing entities
- False patent marking claims are curbed
- Other provisions will make it more difficult to attack patent validity

Complete information about the legislation is available through the Library of Congress THOMAS database.²²¹

11.12.2 The PROTECT IP Act (PIPA). The PROTECT IP Act (Preventing Real Online Threats to Economic Creativity and Theft of Intellectual Property Act)²²² or PIPA, was introduced in the U.S. Senate in May 2011 but failed to make it to the floor of the Senate.²²³ After extensive public opposition, including a worldwide temporary blackout of thousands of Web sites in protest of PIPA and the Stop Online Piracy Act (SOPA, below),²²⁴ the bill was suspended in January 2012 pending further analysis.²²⁵

PIPA’s main points include the following (quoting several sections from the THOMAS database):

11 · 40 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

- Preventing Real Online Threats to Economic Creativity and Theft of Intellectual Property Act of 2011 or the PROTECT IP Act of 2011—(Sec. 3) Authorizes the Attorney General (AG) to commence: (1) an in personam action against a registrant of a nondomestic domain name (NDN) used by an Internet site dedicated to infringing activities (ISDIA) or an owner or operator of an ISDIA accessed through an NDN; or (2) if such individuals are unable to be found by the AG or have no address within a U.S. judicial district, an in rem action (against a domain name itself, in lieu of such individuals) against the NDN used by an ISDIA.
- Defines ISDIA as a site that: (1) has no significant use other than engaging in or facilitating copyright infringement, circumventing technology controlling access to copyrighted works, or selling or promoting counterfeit goods or services; or (2) is designed, operated, or marketed and used to engage in such activities.
- Defines NDN as a domain name for which the registry that issued the domain name and operates the relevant top level domain, and the registrar for the domain name, are located outside the United States.
- Allows the court, upon application by the AG after an NDN-related in personam or in rem action is commenced under this section, to issue a temporary restraining order or an injunction against the NDN, registrant, owner, or operator to cease and desist further ISDIA activity if the NDN is used within the United States to access an ISDIA directing business to U.S. residents and harming U.S. intellectual property right holders.
- Directs the AG to identify and provide advance notice to operators of nonauthoritative domain name system servers (NDNSSs), financial transaction providers (FTPs), Internet advertising services (IASs), and providers of information location tools (ILTs), including search engines, online directories, and other indexes with hypertext links or referrals to online locations, whose action may be required to prevent such NDN-related ISDIA activity.
- Sets forth the preventative measures required to be taken by NDNSSs, FTPs, IASs, and ILTs upon being served with a court order in such an NDN-related action commenced by the AG.
- (Sec. 4) Authorizes the AG or an intellectual property right owner harmed by an ISDIA to commence: (1) an in personam action against a registrant of an ISDIA's domain name or an owner or operator of an ISDIA accessed through a domain name; or (2) if such individuals are unable to be found or have no address within a U.S. judicial district, an in rem action against a domain name used by an ISDIA.
- Allows the court, upon application by the relevant plaintiff after an in personam or in rem action concerning a domain name is commenced under this section, to issue a temporary restraining order or injunction against a domain name, registrant, owner, or operator to cease and desist further ISDIA activity if the domain name is: (1) registered or assigned by a domain name registrar or registry located or doing business in the United States, or (2) used within the United States to access an ISDIA directing business to U.S. residents and harming U.S. intellectual property right holders.
- Directs the relevant plaintiff to identify and provide advance notice to FTPs and IASs whose action may be required to prevent such ISDIA activity.

RECENT DEVELOPMENTS IN INTELLECTUAL PROPERTY LAW 11 · 41

- Requires, upon being served with a court order after such an in personam or in rem action concerning a domain name is commenced by the AG or a private right owner under this section: (1) FTPs to take reasonable specified preventative measures, and (2) IASs to take technically feasible and reasonable measures.
- Sets forth provisions regarding the entities that may be required to take certain preventative measures in actions concerning both domain names and NDNs: (1) granting immunity to such entities for actions complying with a court order, (2) authorizing the relevant plaintiff to bring an action for injunction relief against a served entity that knowingly and willfully fails to comply with a court order, and (3) permitting such entities to intervene in commenced actions and request modifications, suspensions, or terminations of related court orders.
- (Sec. 5) Provides immunity from liability for: (1) FTPs or IASs that, in good faith, voluntarily take certain preventative actions against ISDIAs, and (2) domain name registries and registrars, FTPs, ILTs, or IASs that, in good faith, withhold services from infringing sites that endanger public health by distributing prescription medication that is counterfeit, adulterated, misbranded, or without a valid prescription. ...

11.12.3 The Stop Online Piracy Act (SOPA). The Stop Online Piracy Act (SOPA), H.R. 3261,²²⁶ is summarized in the THOMAS database as follows:

- ... Authorizes the Attorney General (AG) to seek a court order against a U.S.-directed foreign Internet site committing or facilitating online piracy to require the owner, operator, or domain name registrant, or the site or domain name itself if such persons are unable to be found, to cease and desist further activities constituting specified intellectual property offenses under the federal criminal code including criminal copyright infringement, unauthorized fixation and trafficking of sound recordings or videos of live musical performances, the recording of exhibited motion pictures, or trafficking in counterfeit labels, goods, or services.
- Sets forth an additional two-step process that allows an intellectual property right holder harmed by a U.S.-directed site dedicated to infringement, or a site promoted or used for infringement under certain circumstances, to first provide a written notification identifying the site to related payment network providers and Internet advertising services requiring such entities to forward the notification and suspend their services to such an identified site unless the site's owner, operator, or domain name registrant, upon receiving the forwarded notification, provides a counter notification explaining that it is not dedicated to engaging in specified violations. Authorizes the right holder to then commence an action for limited injunctive relief against the owner, operator, or domain name registrant, or against the site or domain name itself if such persons are unable to be found, if: (1) such a counter notification is provided (and, if it is a foreign site, includes consent to U.S. jurisdiction to adjudicate whether the site is dedicated to such violations), or (2) a payment network provider or Internet advertising service fails to suspend its services in the absence of such a counter notification.
- Requires online service providers, Internet search engines, payment network providers, and Internet advertising services, upon receiving a copy of a court order relating to an AG action, to carry out certain preventative measures including withholding services from an infringing site or preventing users located in

11 · 42 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

the United States from accessing the infringing site. Requires payment network providers and Internet advertising services, upon receiving a copy of such an order relating to a right holder's action, to carry out similar preventative measures.

- Provides immunity from liability for service providers, payment network providers, Internet advertising services, advertisers, Internet search engines, domain name registries, or domain name registrars that take actions required by this Act or otherwise voluntarily block access to or end financial affiliation with such sites.
- Permits such entities to stop or refuse services to certain sites that endanger public health by distributing prescription medication that is adulterated, misbranded, or without a valid prescription.
- Expands the offense of criminal copyright infringement to include public performances of: (1) copyrighted work by digital transmission, and (2) work intended for commercial dissemination by making it available on a computer network. Expands the criminal offenses of trafficking in inherently dangerous goods or services to include: (1) counterfeit drugs; and (2) goods or services falsely identified as meeting military standards or intended for use in a national security, law enforcement, or critical infrastructure application.
- Increases the penalties for: (1) specified trade secret offenses intended to benefit a foreign government, instrumentality, or agent; and (2) various other intellectual property offenses as amended by this Act.
- Directs the U.S. Sentencing Commission to review, and if appropriate, amend related Federal Sentencing Guidelines.
- Requires the Secretary of State and Secretary of Commerce to appoint at least one intellectual property attaché to be assigned to the U.S. embassy or diplomatic mission in a country in each geographic region covered by a Department of State regional bureau.

Critics of the legislation include the American Civil Liberties Association, some educators, some law professors, and the United States Student Association.²²⁷ Arguments included the following:

- The bill would lead to removal of much noninfringing content from the Web, resulting in infringement of free speech.
- Eliminating the focus articulated in PIPA about concentrating on sites dedicated to infringing activity would waste government resources on an enormous range of sites.
- ISPs, search engine providers, payment network providers, and advertising services would all have to obey the Attorney General's orders to block all access to sites with infringing content, thus blocking access to all the sites' noninfringing content as well.
- Educational uses could be severely constrained if a single infringing document led to the shutdown of an entire site.
- Sites with a single link to infringing content could be classified as "facilitating" infringement and thus be shut down.

RECENT DEVELOPMENTS IN INTELLECTUAL PROPERTY LAW 11 · 43

- The bill would violate standards of due process by allowing administrative shutdown without providing an opportunity for the owners of the accused sites a chance to defend themselves.
- SOPA's potential barriers to access could severely affect the worldwide movement to pressure dictatorial regimes such as that of the People's Republic of China in their consistent suppression of free access to information.
- Librarians, educators, and students could be subject to administrative shutdown even for what could be justified as fair use of copyright materials.

The proposed bill was dropped at the same time as PIPA (above).

11.12.4 Patent Trolls. Groups aggressively targeting users of little-known patents, often purchased from inventors who have never exercised their rights before, are known as *nonpracticing entities* or *patent trolls*. Some of these companies devote their entire business to suing or threatening to sue on the basis of their acquired patents.²²⁸

In one notorious case, a company bought

... the Canadian patent known as "Automatic Information, Goods, and Services Dispensing System (Canada '216)" whose complete text is available at <http://patents1.ic.gc.ca/> details?patent_number=1236216&language=EN_CA [and] specifically addresses "a system for automatically dispensing information, goods and services to a customer on a self-service basis including a central data processing centre in which information on services offered by various institutions in a particular industry is stored. One or more self-service information and sales terminals are remotely linked to the central data processing centre and are programmed to gather information from prospective customers on goods and services desired, to transmit to customers information on the desired goods or services from the central data processing centre, to take orders for goods or services from customers and transmit them for processing to the central data processing centre, to accept payment, and to deliver goods or services in the form of documents to the customer when orders are completed. The central data processing centre is also remotely linked to terminals of the various institutions serviced by the system, so that each institution can be kept up-dated on completed sales of services offered by that institution." [Note that Canadian spelling is used above.] Think about this patent. Does it not remind you unavoidably of what you did the last time you ordered a book or bought something online? Or performed any other commercial transaction on the Web?²²⁹

A study published by the Boston University School of Law²³⁰ found that patent trolls "... cost U.S. software and hardware companies US\$29 billion in 2011...."²³¹

In the House of Representatives, Peter DeFazio (D-OR) introduced HR.6245, Saving High-Tech Innovators from Egregious Legal Disputes Act of 2012, in August 2012.²³² It would "[Amend] federal patent law to allow a court, upon finding that a party does not have a reasonable likelihood of succeeding in an action disputing the validity or alleging infringement of a computer hardware or software patent, to award the recovery of full litigation costs to the prevailing party, including reasonable attorney's fees...." At the time of writing (May 2013), the bill was still in the hands of the Subcommittee on Intellectual Property, Competition and the Internet of the House Committee on the Judiciary.

In May 2013, Senator Charles Schumer (D-NY) introduced S.866, the Patent Quality Improvement Act, an amendment to the AIA to extend its provisions for challenging patents on business methods.²³³

The Library of Congress THOMAS database describes the substance of the proposal as follows:

11.44 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

- Amends the Leahy-Smith America Invents Act to remove the eight-year sunset provision with respect to the transitional post-grant review program available to review the validity of covered business method patents, thereby making the program permanent.
- Expands the term “covered business method patent” to include a patent that claims a method or corresponding apparatus for performing data processing or other operations used in the practice, administration, or management of any enterprise, product, or service, except technological inventions. (Current law limits the program to financial products or services.)²³⁴

11.13 CONCLUDING REMARKS. Data security ultimately involves the protection of proprietary or personal data and intellectual property. The competition to acquire and retain intellectual property legally is invariably met by unethical and illegal efforts to deprive legitimate owners of their rights. It is necessary, therefore, to be fully aware of the mechanisms and procedures required to protect these rights as part of any computer security program.

This chapter has attempted to delineate the most important aspects of the problem. However, many facets of the legal questions remain unanswered or have been answered generally rather than in the context of a particular problem. Prudent guardians of intellectual property should monitor relevant judicial determinations continuously and be certain to integrate them into a planned approach to protect these most valuable assets.

11.14 FURTHER READING

- Bently, L., and B. Sherman. *Intellectual Property Law*, 3rd ed. Oxford, UK: Oxford University Press, 2008.
- Bloomberg. “Intellectual Property News” Web site. <http://topics.bloomberg.com/intellectual-property>
- Bouchoux, D. E. *Intellectual Property: The Law of Trademarks, Copyrights, Patents, and Trade Secrets*, 4th ed. Cengage Learning, 2012.
- Guardian Newspaper. *Intellectual property* archive of more than 500 recent articles. www.guardian.co.uk/law/intellectual-property
- McJohn, S. *Intellectual Property: Examples and Explanations*, 4th ed. New York, NY: Aspen Publishers, 2012.
- Nard, C. A., D. W. Barnes, and M. J. Madison. *The Law of Intellectual Property*, 3rd ed. New York, NY: Aspen Publishers, 2011.
- Poltorak, A. I., and P. J. Lerner. *Essentials of Intellectual Property*, 2nd ed. Hoboken, NJ: John Wiley & Sons, 2011.
- Stim, R. *Patent, Copyright & Trademark: An Intellectual Property Desk Reference*, 12th ed. Berkeley, CA: Nolo Press, 2012.

11.15 NOTES

1. For the uninitiated, a tort is a civil wrong (i.e., an act or failure to act that violates common law rules of civil society, and is distinguished from criminal wrongdoing).
2. See *ProCD, Inc. v. Zeidenberg*, 86 F.3d 1447 (7th Cir. 1996) (product could be returned if shrink-wrap terms were unacceptable).

NOTES 11 · 45

3. See *Information Handling Services, Inc. v. LRP Publications, Inc.*, 2000 U.S. Dist. LEXIS 14531 (E.D. Pa., Sept. 20, 2000) (limit on unauthorized copies); *Hughes v. America Online, Inc.*, 204 F. Supp. 2d 178 (D. Ma. 2002) (enforcing forum selection clause).
4. See *LLAN Systems, Inc. v. Netscout Service Level Corp.*, 183 F. Supp. 328 (D. Mass. 2002) (click-wrap software agreement enforceable under Uniform Commercial Code as acceptance of an offer).
5. See *Motise v. America Online, Inc.*, 346 F. Supp. 2d 563 (S.D. N.Y. 2004) (user who logged on through another's account is bound by the terms of use even though not read).
6. 356 F.3d 393 (2d Cir. 2004).
7. See *Kewanee Oil Co. v. Bicron Corp.*, 416 U.S. 470, 473, 94 S. Ct. 1879, 40 L. Ed. 2d 315 (1974).
8. It is easy to confuse the notion of common law trade secret law with protection of confidential information. There is a distinction, however. At its core, trade secret law requires commercial application and utility, which is not true of confidential information that is generally protected as a matter of contract. For example, a failed experiment has no commercial utility and is not generally considered a trade secret, although it easily could be deemed confidential information.
9. The need to protect the information from general dissemination is what, in part, has given rise to the practice of Non Disclosure Agreements.
10. UTSA, 14 U.L.A. § 2(a).
11. See Trade Secrets Act, 18 U.S.C. § 1905; see also J. Michael Chamblee, J. D., *Validity, Construction, and Application of Title I of Economic Espionage Act of 1996*, 177 A.L.R. Fed. 609, *2 (2003) (hereinafter "Chamblee at __"). Other federal statutes, such as the National Stolen Property Act, 18 U.S.C. § 2314, were likewise of marginal utility in combating the rising problem of economic espionage. See Chamblee at *2.
12. Craig L. Uhrich, Article: *The Economic Espionage Act—Reverse Engineering and the Intellectual Property Public Policy*, 7 Mich. Telecomm. Tech. L. Rev. 147148-49 (2000/2001) (hereinafter "Uhrich at __"). Uhrich observes that the FBI investigated over 200% more economic espionage cases in 1996 than it had in 1994. See Uhrich at 151.
13. 18 U.S.C. §§ 1831, 1832.
14. *Id.*
15. 18 U.S.C. §§ 1831, 1832.
16. 18 U.S.C. §§ 1832 and 3571.
17. 18 U.S.C. § 1839 (3).
18. *United States v. Lange*, 312 F.3d 263 (7th Cir. 2002) (emphasis added).
19. 18 U.S.C. § 1839.
20. The 1980 Computer Software Copyright Act carved out for owners of computer programs a right to adapt, and for that purpose to copy, the program so that it functions on the actual computer in which it is installed. See discussion under the subheading "Derivative Works."
21. See, e.g., *Computer Management Assistance Co. v. Robert F. DeCastro, Inc.* 220 F.3d 396 (5th Cir. 2000) and *Engineering Dynamics, Inc. v. Structural Software, Inc.*, 26 F.3d 1335 (5th Cir. 1994).

11 · 46 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

22. Ideas, if protectable at all, are protected by patent.
23. The Copyright Act, 17 U.S.C. § 109(b).
24. The Copyright Act itself in sections 108 through 121 provides detailed limitations on the copyright owner's exclusive rights. These limitations are simply a matter of statutory construction. In addition, courts developed the doctrine of fair use in an effort to balance the rights of copyright owner and the public interest. That doctrine is now codified as part of the copyright statute in 17 U.S.C. § 107.
25. See the House Report No. 94-1476, 94th Cong., 2d Sess. 62 (1976) on the 1976 Act.
26. The Copyright Act, 17 U.S.C. §102(b).
27. 797 F.2d 1222 (3rd Cir. 1986).
28. 982 F.2d 693 (2d Cir. 1992).
29. 977 F.2d 1510 (9th Cir. 1992), amended, *Sega Enterprises Ltd. v. Accolade, Inc.*, 1993 U.S. App. Lexis 78.
30. 977 F.2d at 1527–1528.
31. 975 F.2d 832 (Fed. Cir. 1992), *petition for rehearing denied*, 1992 U.S. App. Lexis 30957 (1992).
32. 79 F.3d 1532 (11th Cir. 1996).
33. 350 F.3d 640, 645 (7th Cir. 2003).
34. *Evolution, Inc. v. Suntrust Bank*, 342 F. Supp. 2d 943, 956 (D. Kan. 2004).
35. Compare *Micro Star v. Formgen, Inc.*, 154 F.3d 1107 (9th Cir. 1998) (infringement found because copyrighted images displayed) with *Lewis Galoob Toys, Inc. v. Nintendo of America, Inc.*, 964 F.2d 965 (9th Cir. 1992) (no infringement although product compatible with Nintendo product).
36. 17 U.S.C. § 107.
37. 510 U.S. 569 (1994).
38. *Id.* at 577.
39. *Id.* at 580.
40. See the House Report No. 94-1476, 94th Cong., 2d Sess. 62 (1976) on the 1976 Copyright Act.
41. 17 U.S.C. § 117.
42. *Aymes v. Bonelli*, 47 F.3d 23 (2d Cir. 1995).
43. 17 U.S.C. § 901(a).
44. *MGM Studios Inc. v. Grokster, Ltd.*, 545 U.S. 913, 930 (2005).
45. *Playboy Enterprises v. Frena*, 839 F. Supp. 1552 (M.D. Fla. 1993); see also *Sega Enterprises v. MAPHIA*, 857 F. Supp. 679 (N.D. Cal. 1994), and 948 F. Supp. 923 (N.D. Cal. 1996) (providing site for and encouraging uploading of copyrighted games was copyright infringement).
46. *Religious Technology Center v. Netcom On-line Communication Services, Inc.*, 90 F. Supp. 1361 (N.D. Cal. 1995).
47. 17 U.S.C. § 502.
48. 17 U.S.C. § 503.
49. 17 U.S.C. § 504.
50. 17 U.S.C. § 505.
51. 17 U.S.C. § 506.

NOTES 11 · 47

52. 17 U.S.C. § 504(a).
53. *See Harris Market Research v. Marshall Marketing and Communications, Inc.*, 948 F.2d 1518 (10th Cir. 1991).
54. *See Regents of the University of Minnesota v. Applied Innovations, Inc.*, 685 F. Supp. 698, aff'd, 876 F.2d 626 (8th Cir. 1987) 698.
55. *Id.*
56. *See Cream Records, Inc. v. Jos. Schlitz Brewing Co.*, 754 F.2d 826 (9th Cir. 1985).
57. *See Eales v. Environmental Lifestyles, Inc.*, 958 F.2d 876 (9th Cir. 1992), cert. den. 113 S. Ct. 605.
58. *See Softel, Inc. v. Dragon Medical and Scientific Communications Ltd.*, 891 F. Supp. 935 (S.D. N.Y. 1995). Interestingly, in this case, the court also held that any increase in the infringer's profit may be considered when calculating the profit that must be disgorged to the license holder.
59. *See Harper & Row Publishers, Inc. v. Nation Enterprises*, 471 U.S. 539, 105 S. Ct. 2218 (1985); *Data General Corp. v. Grumman Systems Support Corp.*, 36 F.3d 1147 (1st Cir. 1994).
60. 17 U.S.C. § 504(c)(1).
61. *Id.*
62. 17 U.S.C. § 504(c)(2).
63. The theoretical nature of the relationship between actual and statutory damages is illustrated dramatically when the copyright owner demonstrates that the infringement was willful. *See Peer International Corp. v. Luna Records, Inc.*, 887 F. Supp. 560 (S.D. N.Y. 1995), where the music publisher's president willfully infringed licensed and unlicensed works and was assessed \$10,000 for the licensed works, \$15,000 for the unlicensed works, and \$25,000 that the president used in derivative format without permission even though actual damages were \$4,107. Presumably, this resulted from the court's attempt to find a way to punish the infringer since the statute makes no provision for punitive damages.
64. *See Central Point Software, Inc. v. Nugent*, 903 F. Supp. 1057 (E.D. Tex. 1995).
65. *See Walt Disney Co. v. Powell*, 897 F.2d 565 (D.C.Cir. 1990).
66. 17 U.S.C. § 1202(b).
67. *Universal City Studios, Inc. v. Reimerdes*, 111 F. Supp. 2d 294 (S.D. N.Y. 2000).
68. 17 U.S.C. § 1201(a).
69. 17 U.S.C. § 1201(a)(2).
70. 17 U.S.C. § 1201(a)(3).
71. *Id.*
72. 2000 U.S. Dist. LEXIS 1889 (W.D. Wash. January 18, 2000).
73. *Id.* at 19–21.
74. *Universal City Studios, Inc. v. Reimerdes*, supra note 67.
75. *Universal City Studios v. Corley*, 273 F.3d 429 (2nd Cir. 2002).
76. *Id.* at 446–447.
77. *Id.* at 450–451.
78. 307 F. Supp. 2d 521 (S.D. N.Y. 2004),
79. 320 F.3d 1317 (Fed. Cir. 2003), *writ of certiorari denied*, 539 U.S. 928 (2003).
80. 847 F.2d 255 (5th Cir. 1988).

11 · 48 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

81. *See, e.g., Davidson & Assocs. v. Jung*, 422 F.3d 630, 639 (8th Cir. 2005).
82. 17 U.S.C. § 1201(f)(1).
83. 17 U.S.C. § 1201(f)(2).
84. 17 U.S.C. § 1201(f)(4).
85. 17 U.S.C. § 1201(f)(3).
86. 381 F.3d 1178 (Fed. Cir. 2004), *cert. denied*, 544 U.S. 923 (2005).
87. 387 F.3d 522 (6th Cir. 2004).
88. *Storage Tech. Corp. v. Custom Hardware Eng'g & Consulting, Inc.*, 421 F.3d 1307 (Fed. Cir. 2005).
89. 17 U.S.C. § 1201(i)(1).
90. 17 U.S.C. § 1201(j)(3).
91. 17 U.S.C. § 1201(j)(1).
92. 35 USC § 113 requires the submission of a drawing “where necessary for the understanding of the subject matter to be patented.”
93. 35 U.S.C. §§ 283 and 284.
94. 127 S. Ct. 1746, 1757 (2007).
95. *Webster's Seventh New Collegiate Dictionary* (1967 ed.), p. 644.
96. 2007 CSI/FBI Computer Crime and Security Survey (hereafter the CSI/FBI Survey), pp. 12–13. Although the percentage of organizations reporting Internet abuse is down substantially since this chapter was first published, it nonetheless remains a source of substantial concern. In the same study, 26 percent of respondents reported phishing where the respondent was fraudulently identified as the sender; 25 percent reported misuse of instant messaging and unauthorized access to information; and 17 percent reported theft of customer and/or employee data.
97. SIIA Anti-Piracy 2007 Year in Review (www.siiainc.org/piracy/yir_2007.pdf). According to the SIIA, the source of the financial loss described in the text is the research firm IDC.
98. *See Lamb and Rosen, Global Piracy and Financial Valuation of Intellectual Property*, pp. 11.1–11.3.
99. “The subject matter of copyright … includes compilations.” 17 U.S.C. § 103.
100. *Feist Publications, Inc. v. Rural Telephone Service Co., Inc.* 499 U.S. §§ 340, 361 (1991).
101. *Id.* at 350–351. *See* 17 U.S.C. §§ 101–103.
102. *Id.* at 344, 348–349. *See Ticketmaster Corp. v. Tickets.com, Inc.*, 2000 U.S. Dist. LEXIS 12987 (C.D. Cal. Aug. 10, 2000), *aff'd*, 2001 U.S. App. LEXIS 1454 (9th Cir. Jan. 22, 2001).
103. *Feist Pub., Inc. v. Rural Tel.*, supra note 105, at 352–354, where the court rejected the so-called sweat-of-the-brow doctrine.
104. *Matthew Bender & Co., Inc. v. West Publishing Co.*, 158 F.3d 674, 682 (2d Cir. 1998) (“[t]he creative spark is missing where: (i) industry conventions or other external factors so dictate the selection that any person composing a compilation of the type at issue would necessarily select the same categories of information, or (ii) the author made obvious, garden-variety, or routine selections.”). *See also Silverstein v. Penguin Putnam, Inc.* 368 F.3d 77, 83 (2d Cir. 2004).

NOTES 11 · 49

105. 121 S. Ct. 2381; 150 L. Ed. 2d 500; 2001 U.S. LEXIS 4667; 69 U.S.L.W. 4567 (2001). Note: The party appealing to the Supreme Court is named first.
106. The court found interesting the publishers' decision not to assert a claim of transformative fair use. *Id.* at 2390. See Section 12.1.2.3.3 (transformative use section), supra.
107. 464 U.S. 417 (1984).
108. Transformative fair use was recently applied to the use of Rio devices, which permit individual users to download purchased MP3 music files to a hard drive and then play them either on the PC or a CD. These devices were analogized to the Betamax time shifting discussed in *Sony* and were upheld primarily on that basis. See *Recording Industry Association of America v. Diamond Multimedia Systems, Inc.*, 180 F.3d 1072 (9th Cir. 1999).
109. 17 U.S.C. § 512(c).
110. 17 U.S.C. § 512(k).
111. 17 U.S.C. § 512(i).
112. 17 U.S.C. § 512(c)(1). See *ALS Scan, Inc. v. RemarQ Communities, Inc.*, 239 F.3d 619 (4th Cir. 2001), where the court of appeals determined what notice was sufficient to remove the safe harbor protection. See also *In re Aimster Copyright Litig.*, 252 F. Supp. 2d 634 (N.D. Ill. 2002), aff'd 334 F.3d 643 (7th Cir. 2003), for general discussion of this safe harbor provision, where Aimster had actual knowledge of the infringement by its users and therefore could not avoid liability under the safe harbor.
113. 907 F. Supp. 1361 (N.D. Cal. 1995). The *Netcom* decision predated the DMCA and provided part of the rationale and reasoning used by Congress in drafting and passing Title II of the DMCA. See House Rep. 105-551(I), at 11.
114. The church raised a question of fact about the impact of the ISP's activity on the church's potential market by asserting that the posting of the church's materials on the bulletin board discouraged active participation by existing and potential congregants. Therefore, the court could not find for the ISP as a matter of law.
115. 373 F.3d 544, 555 (4th Cir. 2004). The court went on to state, however, that an ISP "can become liable indirectly upon a showing of additional involvement sufficient to establish a contributory or vicarious violation of the Act. In that case, the ISP could still look to the DMCA for a safe harbor if it fulfilled the conditions therein."
116. See, e.g., *Playboy Enterprises, Inc. v. Frena*, 839 F. Supp. 1552 (M.D. Fla. 1993). The *Frena* decision, insofar as it holds the bulletin board service provider liable for infringement, has been expressly overruled by Title II of the DMCA. See House Rep. 105-551(I), at 11.
117. *Los Angeles Times v. Free Republic*, 2000 U.S. Dist. LEXIS 5669 (C.D. Cal. April 5, 2000). In the *Free Republic* decision, the court recognized the public benefit of posting articles for commentary and criticism but found that the initial postings contained little or no commentary that might transform the article into a new original work. See also *Video Pipeline, Inc., v. Buena Vista Home Entm't, Inc.*, 342 F.3d 191, 199 (3d Cir. 2003), rejecting the fair use defense for an online distributor that made its own movie clip previews and used them as movie trailers by copying short segments of plaintiff's movies in part because the online distributor benefited from the infringement.

11 · 50 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

118. *Storm Impact, Inc. v. Software of the Month Club*, 13 F. Supp. 2d 782 (N.D. Ill. 1998).
119. There are various names for the components of the software programs that actually travel through the Web looking for data, including bots, crawlers, spiders, scrapers, and automated data retrieval systems.
120. *Kelly v. Arriba Soft Corp.*, 336 F.3d 811 (9th Cir. 2003).
121. 416 F. Supp. 2d 828, 838–846 (C.D. Calif. 2006).
122. 508 F.3d 1146 (9th Cir. 2007).
123. 239 F.3d 1004 (9th Cir. 2001), *aff'd*, 284 F.3d 1091 (9th Cir. April 3, 2002).
124. See also *UMG Recordings, Inc. v. MP3.com, Inc.*, 92 F. Supp. 2d 349 (S.D. N.Y. 2000), where the district court held that storing recordings from purchased CDs on MP3.com's servers for retransmission to other users was infringement and not transformative fair use.
125. *Metro-Goldwyn-Mayer Studios, Inc. v. Grokster, Ltd.*, 545 U.S. 913, 125 S. Ct. 125 (2005).
126. *Id.* at 926.
127. *Id.* at 936–937.
128. *Id.* at 937.
129. 100 F. Supp. 2d 1058 (N.D. CA 2000).
130. 30 Cal. 4th 1342; 71 P.3d 296; 1 Cal. Rptr. 3d 32 (2003).
131. 356 F.3d 393 (2d Cir. 2004).
132. 86 F.3d 1447 (7th Cir. 1996).
133. 18 U.S.C. § 1030.
134. Pub. L. 98-474, codified at 18 U.S.C. § 1030.
135. Pub. L. 99-474.
136. National Information Infrastructure Protection Act of 1996, Pub. L. 104–294.
137. 18 U.S.C. § 1030(a)(1).
138. 18 U.S.C. § 1030(a)(2).
139. 18 U.S.C. § 1030(a)(3).
140. 18 U.S.C. § 1030(a)(4).
141. 18 U.S.C. § 1030(a)(5). See *Hotmail Corporation v. Van\$ Money Pie, Inc.*, 1998 WL 388389, 47 U.S.P.Q.2d 1020 (N.D. Cal. 1998).
142. 18 U.S.C. § 1030(a)(6).
143. 18 U.S.C. § 1030(a)(7).
144. Senate Rep. 104-357, pp. 7–8.
145. *Register.com, Inc. v. Verio, Inc.*, 126 F. Supp. 2d 238 (S.D. N.Y. 2000).
146. *Id.*
147. *In Re Intuit Privacy Litigation*, 138 F. Supp. 2d 1272 (2001). But see *U.S. v. Czubinski*, 106 F.3d 1069 (1st Cir. 1997), where the court of appeals found that an IRS employee who accessed private tax information in violation of IRS rules but did not disclose the accessed information could not be prosecuted under 18 U.S.C. §030(a)(4) because he lacked an intent to deprive the affected taxpayers of their right to privacy.
148. *Shurgard Storage Centers, Inc. v. Safeguard Self Storage, Inc.*, 119 F. Supp. 2d 1121 (W.D. Wash 2000).

NOTES 11 · 51

149. *America Online, Inc. v. LCGM, Inc.*, 46 F. Supp. 2d 444 (E.D. Va. 1998).
150. *U.S. v. Middleton*, 231 F.3d 1207 (9th Cir. 2000).
151. Council Directive 95/46, 1995 O.J. (L.281) 31–50 (EC).
152. As a result the United States negotiated with the EU the Safe Harbor Arrangement, administered by the Federal Trade Commission, under which a U.S. company can opt in to compliance with the EU Data Directive.
153. For an updated list, go to www.pirg.org/consumer/credit/statelaws.htm.
154. 149 N.H. 148, 816 A.2d 1001 (2003).
155. 18 U.S.C. §§ 2511(1)(a) and 2502(a).
156. Pub. L. No. 99-508, 100 Stat. 1848 (codified throughout scattered sections of 18 U.S.C.).
157. S. Rep. No. 99-541, at 1 (1986), reprinted in 1986 U.S.C.C.A.N. 3555, 3555.
158. 18 U.S.C. § 2510(12).
159. *Id.* § 2510(4).
160. 18 U.S.C. § 2511(1)(a) (emphasis added); *Konop v. Hawaiian Airlines, Inc.*, 302 F.3d 868, 875 (9th Cir. 2002) (*Konop*) (noting the legislative history of the ECPA indicates that Congress wanted to protect electronic communications that are configured to be private, such as email and private electronic bulletin boards).
161. 18 U.S.C.A. § 511(2)(d). One should note, however, that as a result of the Patriot Act, an order from a U.S. or state attorney general is sufficient to permit the government to install a device to record electronic transmissions for up to 60 days where related to an ongoing criminal investigation. The FBI has in its arsenal a program known as Carnivore that essentially tracks a target's online activity. Recently, Freedom of Information inquiries by the Electronic Privacy Information Center (EPIC, www.epic.org) suggests that the FBI has discontinued use of Carnivore because ISPs, in light of the PATRIOT Act, may be providing information regarding a user's internet traffic directly to the government.
162. *Griggs-Ryan v. Smith*, 904 F.2d 112, 117 (1st Cir. 1990) (holding consent may be implied where the individual is on notice of monitoring of all telephone calls).
163. Federal law allows states to enact their own wiretapping statutes provided that the state statutes are at least as strict as the federal counterpart. Lynn Bernabei, Ethical and Legal Issues of Workplace Monitoring of Employee Communications, 2003 WL 22002093, *2 (April 2003) (hereinafter “Bernabei at __”). Bernabei notes that most states have adopted statutes that mirror the federal statutes and that at least 10 states, including Massachusetts, require the consent of both parties before the employer can record a conversation. *Id.*
164. 18 U.S.C. § 2511(2)(a)(i) (Supp. 2003).
165. *Briggs v. Am. Air Filter Co.*, 630 F.2d 414 (5th Cir. 1980) (holding employer could monitor employee's communication “when [the] employee's supervisor [had] particular suspicions about confidential information being disclosed to a business competitor, [had] warned the employee not to disclose such information, [had] reason to believe that the employee is continuing to disclose the information, and [knew] that a particular phone call is with an agent of the competitor.”).
166. 18 U.S.C. § 2511(1)(a).
167. *Konop*, 302 F.3d at 876 (emphasis added).
168. *Id.*

11 · 52 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

169. E.g., *Wesley Coll. v. Pitts*, 974 F. Supp. 375, 386 (D. Del. 1997) (holding that the act criminalizes only the interception of electronic communications contemporaneously with their transmission, not once they have been stored); *Payne v. Norwest Corp.*, 911 F. Supp. 1299, 1303 (D. Mont. 1995) (holding the appropriation of voicemail or similar stored electronic message does not constitute an “interception” under the act); *Steve Jackson Games, Inc. v. United States Secret Service*, 36 F.3d 457, 461–462 (5th Cir. 1994) (holding that the government’s acquisition of email messages stored on an electronic bulletin board system, but not yet retrieved by the intended recipients, was not an “interception” under the Wiretap Act).
170. See *United States v. Councilman*, 418 F.3d 67, 69–70 (1st Cir. 2005) (*en banc*).
171. *In re Pharmatrak, Inc.*, 329 F.3d 9, 21 (1st Cir. 2003).
172. *Id.*
173. *United States v. Councilman*, 245 F. Supp. 2d 319 (D. Mass. 2003) (Wiretap Act count dismissed against email service provider who was charged with attempting to use electronic communications passing through his service for commercial gain).
174. *Eagle Investment Systems, Corp. v. Tamm*, 146 F. Supp. 2d 105, 112–113 (D. Mass. 2001).
175. USA PATRIOT Act § 209, 115 Stat. at 283; *Konop*, 302 F.3d at 876–878 (“The purpose of the recent amendment was to reduce the protection of voice mail messages to the lower level of protection provided other electronically stored communications.”)
176. 302 F.3d at 876.
177. *Id.* at 879.
178. 18 U.S.C. § 2701 et seq.
179. *Bernabei* at *2.
180. 18 U.S.C. §§ 2701(a)(1), 2707(a) (emphasis added).
181. *Id.* § 2510(17), incorporated by 18 U.S.C. § 2711(1).
182. 18 U.S.C. § 2701(c)(1).
183. 18 U.S.C. § 2701(c)(2).
184. *Theofel v. Fary-Jones*, 359 F.3d 1006 (9th Cir. 2004).
185. *In re DoubleClick, Inc. Privacy Litigation*, 154 F. Supp. 2d 497 (S.D.N.Y. 2001) (*DoubleClick*); *In re Toys R US, Inc. Privacy Litigation*, 2001 U.S. Dist. LEXIS 16947 (N.D. Ca. 2001).
186. 154 F. Supp. 2d at 511–512.
187. *Fraser v. Nationwide Mut. Ins. Co.*, 2003 U.S. App. LEXIS 24856, *19 (3rd Cir. 2003).
188. Jeanie Duncan Fallon, *Open Source Licenses: Understanding the General Public License, Technology Licensing Primer*, p. 248 (2d ed. 2001).
189. Richard Stallman, The GNU Project, available at www.gnu.org/gnu/thegnuproject.html. Stallman also started the Free Software Foundation (FSF) in 1985.
190. John C. Yates and Paul H. Arne, *Open Source Software Licenses: Perspectives of the End User and the Software Developer*, 25th Annual Institute on Computer & Internet Law, vol. 2, p. 104 (2005). It is estimated that thousands of programmers have contributed to Linux.

NOTES 11 · 53

191. It is considered less extreme than the FSF, which basically advocates for an end to proprietary rights as applied to software.
192. See www.opensource.org/licenses/
193. This makes sense especially considering the FSF's vision of free software and its insistence on setting forth those views in the preamble of the GPL.
194. Section 2 of the GPL states: "You must cause any work that you distribute or publish, that in whole or in part contains or is derived from the Program or any part thereof, to be licensed as a whole at no charge to all third parties under the terms of this License."
195. Fallon, at 250.
196. GPL Version 2.
197. *Id.* at 2(c).
198. Lori E. Lesser, *Open Source Software: Risks, Benefits, & Practical Realities in the Corporate Environment*, *Open Source Software: Risks, Benefits, & Practical Realities in the Corporate Environment*, p. 41 (2004).
199. See *id.*
200. *Progress Software Corp. v. MySQL AB*, 195 F. Supp. 328, 329 (D. Mass. 2002).
201. *Id.* The court also noted that MySQL did not demonstrate the likelihood of irreparable harm during the pendency of the case.
202. See Stuart D. Levi and Andrew Woodard, "Open Source Software: How to Use It and Control It in the Corporate Environment," *The Computer Lawyer*, vol. 21 (Aug. 8, 2004). "[A] policy needs to balance the benefits and competitive advantages of open source with the risks of using source code developed by parties with whom the company may not have a formal relationship."
203. See Yates and Arne, *supra* n. 195, p. 107.
204. See Levi and Woodard, *supra* n. 207.
205. Also consider the fact that discovery in the course of litigation would also involve releasing proprietary codes, as IBM was forced to do for some of its products involved in the *SCO* litigation. Although discovery is obviously a different publication from that required under the GPL, it is an important issue to consider.
206. *Id.*
207. The Paris Convention was initially concluded in 1883 and updated in 1967. It is administered by the World Intellectual Property Organization, an agency of the United Nations. The Paris Convention has provisions that apply to patents, trademarks, service marks, industrial designs (similar to design patents), and unfair competition. Approximately 174 nations are now signatories to the Paris Convention.
208. Until the adoption of TRIPS, the Berne Convention was the other major international agreement. Like the Paris Convention, it is administered by the World Intellectual Property Organization. The Berne Convention, first adopted in 1886, has undergone a series of revisions. The Convention includes "every production in the literary, scientific and artistic domain whatever may be the mode or form of its expression." Berne Convention, Art. 2, ¶ 1. Essentially, it assures that a work protected within a member state will also be protected outside of the member state without being subject to discriminating formalities. The number of signatories to the Berne Convention is presently 165 nations.

11 · 54 FUNDAMENTALS OF INTELLECTUAL PROPERTY LAW

209. The WTO effectively began operating on July 1, 1995, as a result of the 1994 Uruguay Round Agreements. The WTO replaces GATT (General Agreement on Tariffs and Trade), which had been in operation since 1950. Congress ratified the Uruguay Round Agreements in December 1994. The WTO has approximately 132 member nations. In 1995, the WTO and the World Intellectual Property Organization (WIPO) signed a joint agreement that provides, among other things, for cooperation in providing legal technical assistance and technical cooperation related to the TRIPS Agreement for developing country members of either of the two organizations. The WIPO has approximately 171 members and is responsible for international cooperation in promoting intellectual property protection around the world. In particular, it looks after various international conventions, such as the Paris Convention and the Berne Convention.
210. Geographical indications are marks or other expressions that state the country, region, or place in which a product or service originates.
211. Industrial designs protect the aesthetic look of the product and are similar but not identical to the United States notion of trade dress. Products may be afforded protection based on novelty or originality of design, depending on national law.
212. Breeder's rights confer protection on new and different plant varieties.
213. Utility models protect the manner in which a product works or functions and as such are different from industrial design, which protects only the aesthetics of the product. Generally, utility models address mechanical functioning, which in the United States is not protectable unless patentable. Thus, the innovation in the United States must be significant to warrant protection.
214. Until 1989, the developing countries largely refused to negotiate standards. Threats by the United States of trade sanctions under the United States Trade Act played a significant role in altering the positions of economically weaker developing countries. In particular, China, India, Taiwan, and Thailand were all investigated.
215. *Feist Publications v. Rural Telephone System*, 499 U.S. 340 (1991).
216. TRIPS, Article 41.
217. TRIPS, Article 61.
218. This section was written by M. E. Kabay.
219. Leahy-Smith American Invents Act of 2011, H.R. 1249, www.govtrack.us/congress/bills/112/hr1249
220. Peter E. Heuser, "Recent Developments in IP Law of Interest to Business Attorneys," Web site of Schwabe, Williamson & Wyatt, February 26, 2013, www.schwabe.com/showarticle.aspx?Show=12770
221. Library of Congress, 2011.
222. Library of Congress, 2012.
223. Preventing Real Online Threats to Economic Creativity and Theft of Intellectual Property Act of 2011, S. 968, www.govtrack.us/congress/bills/112/s968/text
224. Jonathan Weisman, "After an Online Firestorm, Congress Shelves Antipiracy Bills," *New York Times*. January 20, 2012, www.nytimes.com/2012/01/21/technology/senate-postpones-piracy-vote.html?r=0
225. Trevor Timm, "After Historic Protest, Members of Congress Abandon PIPA and SOPA in Doves," *Electronic Frontier Foundation | Deplinks* (blog), January

NOTES 11 · 55

- 19, 2012. <https://www.eff.org/deeplinks/2012/01/after-historic-protest-members-congress-abandon-pipa-and-sopa-droves>
226. Library of Congress, 2011.
227. Electronic Frontier Foundation, Collection of documents about SOPA. *Electronic Frontier Foundation* Web site. December 12, 2011, <https://www.eff.org/search/site/sopa>
228. Mark Hachman, “Inside The Mind of A Patent Troll: If It’s Legal, It Must Be OK,” *readwrite.com*, May 13, 2013, <http://readwrite.com/2013/05/03/inside-the-mind-of-a-patent-troll-if-its-legal-it-must-be-ok>
229. M. E. Kabay, “PanIP has rights: PanIP has patents on e-commerce-related systems,” *NetworkWorld*, April 15, 2003, www.networkworld.com/newsletters/sec/2003/0414sec1.html
230. James E. Bessen and Michael J. Meurer, “The Direct Costs from NPE Disputes: Abstract,” *Social Science Research Network* Web site, June 28, 2012, http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2091210##
231. Loek Essers, “Patent Trolls’ Cost Tech Companies \$29 Billion Last Year, Study Says,” *PCWorld*, June 27, 2012, www.pcworld.com/article/258395/patent_trolls_cost_tech_companies_29_billion_last_year_study_says.html
232. Library of Congress, 2012.
233. Grant Gross, “Senator introduces legislation targeting patent trolls,” *PCWorld*, May 01, 2013, www.pcworld.com/article/2037005/senator-introduces-legislation-targeting-patent-trolls.html
234. Library of Congress, 2013.

Computer Security Handbook, Sixth Edition, Volume I
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

INTRODUCTION TO PART II

THREATS AND VULNERABILITIES

What are the practical, technical problems faced by security practitioners? Readers are introduced to what is known about the psychological profiles of computer criminals and employees who commit insider crime. The focus is then widened to look at national security issues involving information assurance—critical infrastructure protection in particular. After a systematic review of how criminals penetrate security perimeters—essential for developing proper defensive mechanisms—readers can study a variety of programmatic attacks (widely used by criminals) and methods of deception, such as social engineering. The section ends with a review of widespread problems such as spam, phishing, Trojans, Web server security problems, and physical facility vulnerabilities (an important concern for security specialists, but one that is often overlooked by computer-oriented personnel).

The chapter titles and topics in Part II include:

- 12. The Psychology of Computer Criminals.** Psychological insights into motivations and behavioral disorders of criminal hackers and virus writers
- 13. The Insider Threat.** Identifying potential risks among employees and other authorized personnel
- 14. Information Warfare.** Cyberconflict and protection of national infrastructures in the face of a rising tide of state-sponsored and non-state-actor industrial espionage and sabotage
- 15. Penetrating Computer Systems and Networks.** Widely used penetration techniques for breaching security perimeters
- 16. Malicious Code.** Dangerous computer programs, including viruses and worms, increasingly used to create botnets of infected computers for spreading spam and causing denial of service
- 17. Mobile Code.** Analysis of applets, controls, scripts, and other small programs, including those written in ActiveX, Java, and Javascript
- 18. Denial-of-Service Attacks.** Resource saturation and outright sabotage that brings down availability of systems and that can be used as threats in extortion rackets by organized crime

II · 2 THREATS AND VULNERABILITIES

19. **Social-Engineering and Low-Tech Attacks.** Lying, cheating, impersonation, intimidation—and countermeasures to strengthen organizations against such attacks, which have increased drastically in recent years
20. **Spam, Phishing, and Trojans: Attacks Meant to Fool.** Fighting spam, phishing, and Trojans—trickery that puts uninformed victims at serious risk of fraud such as identity theft
21. **Web-Based Vulnerabilities.** Web servers, and how to strengthen their defenses
22. **Physical Threats to the Information Infrastructure.** Attacks against the information infrastructure, including buildings and network media

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 12

THE PSYCHOLOGY OF COMPUTER CRIMINALS

Q. Campbell and David M. Kennedy

12.1 INTRODUCTION	12·1	12.5.3 Asperger Syndrome and Computer Criminals	12·11
12.2 SELF-REPORTED MOTIVATIONS	12·3	12.5.4 Internet Abuse and Computer Crime	12·12
12.3 PSYCHOLOGICAL PERSPECTIVES ON COMPUTER CRIME	12·4	12.6 ETHICS AND COMPUTER CRIME	12·14
12.4 SOCIAL DISTANCE, ANONYMITY, AGGRESSION, AND COMPUTER CRIME	12·4	12.7 CLASSIFICATIONS OF COMPUTER CRIMINALS	12·16
12.4.1 Social Presence and Computer Crime	12·6	12.7.1 Early Classification Theories of Computer Criminals	12·17
12.4.2 Deindividuation and Computer Crime	12·6	12.7.2 Rogers's New Taxonomy of Computer Criminals	12·19
12.4.3 Social Identity Theory and Computer Crime	12·7	12.7.3 Hacktivists and Cyberterrorists: Hacking for a Cause	12·20
12.4.4 Social Learning Theory of Computer Crime	12·8	12.7.4 Dangerous/Malicious Insiders (DI/MI)	12·21
12.5 INDIVIDUAL DIFFERENCES AND COMPUTER CRIMINALS	12·9	12.7.5 Virus Creators	12·22
12.5.1 Antisocial and Narcissistic Personalities	12·9	12.8 RECOMMENDATIONS	12·24
12.5.2 Five-Factor Model of Personality and Computer Criminals	12·10	12.9 FURTHER READING	12·26
		12.10 NOTES	12·26

12.1 INTRODUCTION. Symantec's Internet Security Threat Report for 2011 indicated an 81 percent increase in network attacks compared to 2010; this, coupled with a reported 187 million identities that were exposed due to outsider attacks, suggests that the threat of computer crime is growing to unprecedented levels.¹ For the most part, the industry has relied upon legal and technological solutions to reduce the risks to information security. An alternate approach is to target the human element;

12 · 2 THE PSYCHOLOGY OF COMPUTER CRIMINALS

to try to understand the psychological motivations behind those who would exploit these technologies and to design security system accordingly.² The main drawback to this approach is that the information security field has traditionally relied on outdated stereotypes of computer criminals, which have lead to convoluted, overgeneralized, and inaccurate portrayals of these individuals. Contributing to this is the industry's oversimplification of computer crime and its reliance on a generic, all-encompassing view of the computer criminal. The computer underground is a vast and varied landscape that comprises many different subgroups, some of which are infantile and benign, and others that are criminal and destructive. The purpose of this chapter is to identify the various subgroups of computer criminals and to examine their differing motivations from a psychological perspective. Using theoretical perspectives from social, personality, and clinical psychology, we will review current research on the various subsets of computer criminals, ranging from script kiddies to malicious insiders, and provide recommendations for addressing the problem of computer crime at its source.

The National Institute of Justice defines a *computer criminal* as any individual who uses computer or network technology to plan or perpetrate a violation of the law.³ Although the term *computer hacker* is often used interchangeably with computer criminal, they are not synonymous. The term *hacker* was originally used as an umbrella term to refer to a computer programmer who changes or alters code (i.e., hacks) in a unique or unorthodox fashion to solve a problem or to enhance its use. Such interventions may be legal or illegal depending on the circumstances, intent, outcome, or use of the hacked program.

Although a computer criminal, or *cracker*, sometimes also referred to as a *malicious or criminal hacker*, may fall under this broad definition, these individuals typically alter or exploit technology for destructive purposes or financial gain rather than for benign or creative functions. Common examples of computer crimes include Web page defacements, creation and distribution of viruses, unauthorized access of technology, theft of information, distributed denials of service (DDoS), and so on.

In recent years, computer security analysts have reported that many computer criminals are moving away from the hacking-for-fun-and-notoriety mindset to hacking for profit.⁴ More recently, activist groups (*hacktivists*) have used the Internet as a way of spreading their messages of social and political discord by engaging in digital harassment of their targets. According to the 2012 “Verizon Data Breach Report,” hacktivist groups represent a significant threat to network security, accounting for the majority of data thefts occurring in 2011.^{5,6}

Computer crime is an obvious financial and societal problem that shows no signs of slowing. Researchers suggest that computer attacks will continue to grow in frequency and sophistication as technology continues to evolve. More specialized threats to social networks, peer-to-peer networks (P2P), handheld mobile devices, and nontraditional hardware systems (e.g., networked gaming consoles and point of sale devices), have been identified in the wild with increasing regularity.⁷ Douglas Campbell, president of the Syneca Research Group Inc., states that “the dominant threat to the United States is not thermonuclear war, but the information war.”^{8,9}

One solution that has been offered as an effort to slow this disturbing trend is to examine the motivations of computer criminals from a psychological perspective. Computer-crime researchers suggest that understanding the psychological motivations behind cyber criminals would aid in both cybercrime prevention and protection.^{10,11} Generating a psychological profile of the various subtypes of computer criminals would aid in creating preventive initiatives as well as more effective countermeasures in the fight against computer crime. Since computer crime is not solely a technological

SELF-REPORTED MOTIVATIONS 12 · 3

issue but one involving human agents, psychological theories regarding anonymity, aggression, social learning, and individual difference factors may enable us to better understand the behaviors and motivations of the computer criminal.

Information security consultant Donn Parker asserts that the creation of an effective malicious hacker profile still remains an elusive goal in the information security field.¹² Therefore, the goal of this chapter is to survey past and current literature surrounding the psychological motivations of computer criminals. Theories from criminology, as well as social, personality, and clinical psychology, will be presented in an attempt to explain some of the possible motivations behind computer criminals. Based on these psychological research studies, we will conclude by offering recommendations for attenuating computer crime from the perspective of the perpetrator.

12.2 SELF-REPORTED MOTIVATIONS. Perhaps the simplest approach to understanding the mindset of computer criminals is having the perpetrators describe their motivations in their own words. Using various self-reporting measures, including surveys, open-ended questionnaires, and first-person interviews, researchers have consistently found a number of common accounts used by computer criminals to explain and justify their illicit and sometimes harmful behaviors.^{13,14,15,16}

According to sociologist Paul Taylor, computer criminals report that they are motivated by an interacting mix of six primary categories: addiction, curiosity, boredom, power, recognition, and politics.¹⁷ Using a phenomenological-interpretive interview approach that emphasizes the interviewee's perception of reality, sociologist Orly Turgeman-Goldschmidt similarly found that computer criminals reported curiosity, thrill seeking, the need for power, and the ideological opposition to information restrictions among the motivations for their behaviors.¹⁸ Taylor suggests that the extensive use of computers by these criminals may result from a combination of both compulsive behaviors and intellectual curiosity. From an outsider's perspective, an advanced computer user's need to meet the swiftly changing demands of the computer industry may appear to be an indicator of computer abuse, when in actuality the constant use of technology is a consequence of the field. A relentless curiosity and desire for technological improvement is often used by computer criminals as a motivation for their behaviors.¹⁹ Anecdotal evidence has also suggested that the frustrations that result from restrictive computing environments (e.g., network or Internet filters), coupled with a lack of sufficient intellectual stimulation, contribute to some computer criminals' unauthorized access attempts. Some reformed computer criminals have indicated that once they were provided with more liberal access to technology, they were able to focus their skills on practical and legal endeavors rather than illicit undertakings.²⁰

In one of the most ambitious efforts to understand the mindset of criminal hackers, computer security consultant and reformed computer criminal Raoul Chiesa and colleagues created the Hacker's Profiling Project (HPP).²¹ The aim of the HPP was to utilize criminal profiling techniques to develop a comprehensive profile of criminal hackers. Chiesa developed a questionnaire that was judiciously distributed to known criminal hackers and asked questions regarding personal demographics, technological skill, criminal history, and social relationships. The research revealed again that some of the main motivations of criminal hackers are curiosity, proving their self-worth to themselves and others, and feelings of adventure. These relatively benign motivations are also coupled with feelings of anger, frustration, and rebellion against authority. For many of these individuals, the Internet is viewed as the great equalizer. Because of the reduced social context cues that guide face-to-face interactions, Internet users are judged more on their technological skills rather than their social skills, gender, or

12 · 4 THE PSYCHOLOGY OF COMPUTER CRIMINALS

ethnicity. For a more comprehensive look at the results of the HPP, see the recently published book *Profiling Hackers*.²²

Contrary to their stereotypical portrayals in the news media and in fiction, computer criminals appear to have wide-ranging social networks that exist in both their online and offline environments.^{23,24} Taylor indicates that both the need for power and recognition by their peers may be motivating factors for some cybervandals. Computer criminals report feelings of enjoyment and satisfaction when they prove themselves better than system administrators and their peers. Communications researchers Hyung-jin Woo, Yeora Kim, and Joseph Dominick report in their analysis of Webpage defacements that 37 percent of the prank-related defacements contained messages that bragged or taunted the system administrators. Twenty-four percent of these types of defacements contained statements aimed at obtaining peer recognition and 8 percent contained boastful and self-aggrandizing verbiage.²⁵

12.3 PSYCHOLOGICAL PERSPECTIVES ON COMPUTER CRIME. Although self-reporting analyses can give us some insight into the motivations behind computer criminals, these types of descriptive methodologies typically yield incomplete and sometimes inaccurate results. Unless the causes for our behaviors are obvious, our explicit or consciously held explanations for our actions are often misguided. Research has found that our behaviors are frequently controlled by subtle situational variables and implicit attitudes of which we are not typically aware, and may be distinctly different from the conscious mechanisms we use to explain our actions.²⁶ Therefore, our conscious justifications for our actions may be inaccurate if we are unaware of more subtle cognitive processes. The next section examines more empirically based psychological theories of aggression and deviance to gain a further understanding of the factors that may be influencing the behaviors of computer criminals.

12.4 SOCIAL DISTANCE, ANONYMITY, AGGRESSION, AND COMPUTER CRIME. Many acts of computer crime can be categorized as demonstrations of aggressive behaviors. For example, cracking into a company's Web server and defacing a Web page, or launching a DDoS attack on an organization's computer network, thereby crippling its Internet connection, are common malicious and aggressive acts engaged in by computer criminals. Social psychological theories on hostility and violence suggest that people are more likely to commit acts of aggression when the perpetrator of these acts is anonymous and the threat of retaliation is low.²⁷ Since cybervandals frequently use nicknames (*nicks* or *handles*), stolen accounts, and spoofed Internet Protocol (IP) addresses when they engage in illegal activities, their behaviors may be more aggressive than when they are more easily identifiable. Computer criminals are overly confident that their crimes cannot and will not be traced back to their true identities. Computer criminals who deface Web pages are so confident that they are anonymous that they regularly *tag* the hacked Website by leaving their handles and the handles of their friends, and in some cases, even their Internet email addresses and Web page links.²⁸

Due to the relative anonymity of the Internet and the technical abilities of cybercriminals, which enable them to further obfuscate their identities, the resulting emotional distance may be another factor that contributes to increased aggression online. For example, it is an extremely difficult and tedious task to identify computer criminals who launch DDoS attacks against computer networks. The attacker plants denial-of-service (DoS) programs into compromised shell accounts controlled by a master client. The master client will instruct every slave DoS program to *cooperatively* launch an

SOCIAL DISTANCE, ANONYMITY, AGGRESSION, AND COMPUTER CRIME 12 · 5

attack at the victim's host at a configurable time and date. Thus, the DoS attacks are not launched by the criminal's own computer; rather, the attacks come from innocent networks that have been compromised by the cracker. This additional layer makes it all the more difficult for information-security professionals to locate the attack's perpetrator. Much like the Website vandals, DDoS attackers are also confident that the attacks will not be traced back to their actual identities. Frequently, DDoS attackers will even brag on Internet Relay Chat (IRC) channels and Twitter about how many host nodes they have compromised and against which domain they are planning to launch new attacks.²⁹

Situational influences on behaviors and attitudes work on the Internet much as they do in the real world. However, computer criminals who commit aggressive acts against their innocent victims do not see the immediate consequences of their actions. The computer screen and increased social distance that characterize interactions online can act as an electronic buffer between the attacker and victim. Like the participants in psychologist Stanley Milgram's famous obedience experiment, computer criminals are physically and emotionally removed from their victims while they are committing their harmful actions.³⁰ They do not witness firsthand the consequences of their computerized attacks. Automated cracking and DDoS scripts, coupled with the lack of social presence in computer-mediated interactions, may make it easier to attack an entity that is not only emotionally and physically distant, but also depersonalized (e.g., a system administrator working for a large corporation).

Consistent with social psychologist Albert Bandura's theory of moral disengagement, individuals who engage in unscrupulous behaviors will often alter their thinking in order to justify their negative actions.³¹ According to Bandura, most individuals will not commit cruel or illicit behaviors without first engaging in a series of cognitive justification strategies that allow the person to view those actions as moral and just. Immoral behaviors can be justified by comparing them to more egregious acts, minimizing the consequences of the actions, displacing responsibility, and by blaming the victim themselves. Criminologist Marc Rogers posits that computer criminals may rely on a number of these disengagement strategies in an attempt to reduce the dissonance associated with their malicious activities.³²

Studies conducted by sociologists Paul Taylor and Orly Turgeman-Goldschmidt suggest that many computer criminals are, in fact, engaging in forms of moral disengagement.^{33,34} Their interviewees report that computer crime is driven by a search for answers and spurs the development of new technologies. They further indicate that their electronic intrusions cause no real monetary harm or damage to the victims, and that larger corporations that can afford any financial losses that are incurred from their digital transgressions. Web page crackers will often criticize and publicly taunt the system administrators for not properly securing their computers, suggesting that the victims deserved to be attacked.³⁵ Rogers suggests that this victim-blaming strategy is likely the most common form of moral disengagement that is employed by computer criminals.³⁶

A study conducted by information technology researcher Randall Young and colleagues confirmed that computer criminals have a morally distorted view of their deviant activities, enabling them to socially justify their digital exploits.³⁷ Self-identified computer criminals attending a computer conference reported significantly higher levels of moral disengagement than a control group of university students. The self-identified criminals strongly felt that their digital intrusions were actually helpful to the companies that they invaded and that their friends and families would not think negatively of them if they were caught engaging in illegal computer hacking.

12 · 6 THE PSYCHOLOGY OF COMPUTER CRIMINALS

12.4.1 Social Presence and Computer Crime. Social psychologist Sara Kiesler and colleagues found that during face-to-face (FTF) interactions, conversants implicitly attend to social context cues (e.g., facial expressions and intonations) and use them to guide their social behaviors. Since these social context cues are absent or reduced during computer-mediated interactions, digital communication may be more deregulated than FTF discussions.^{38,39} Kiesler and Lee Sproull suggest that the absence of social-context cues in computer-mediated communication hinders the perception of and adaptation to social roles, structures, and norms.⁴⁰ The reduction of social-context cues in computer-mediated communication can lead to deregulated behavior, decreased social inhibitions, and reduced concern with social evaluation. The most common variables examined in their experiments were hostile language in the form of “flaming” (aggressive, rude, and often ad hominem attacks) and post hoc perceptions of group members (i.e., opinions formed after interacting with members). One empirical study found that group members communicating via computer-mediated communication were more hostile toward one another, took longer to reach decisions, and rated group members less favorably than comparable face-to-face groups.⁴¹ Another experiment reported that there were 102 instances of hostile communication during computer-mediated interactions, compared to only 12 instances of hostile commentary during comparable FTF discussions.⁴²

Based on Kiesler’s findings, computer criminals may be engaging in hostile behaviors partly due to this reduction of available context cues. Crackers who harass and victimize system administrators and Internet users may be engaging in these antisocial activities due to the reduced attention to and concern with social evaluations. There are numerous anecdotal accounts of computer criminals “taking over” IRC channels, harassing people online, deleting entire computer systems, and even taunting system administrators whose networks they have compromised.⁴³ Their criminal and aggressive behaviors may be partially attributed to the reduced social context cues in computer-mediated communication and the resulting changes in their psychological states (i.e., deindividuation) while online.

12.4.2 Deindividuation and Computer Crime. Disinhibited behaviors have also been closely linked to the psychological state of deindividuation. Deindividuation is described as a loss of self-awareness that results in irrational, aggressive, antinormative, and antisocial behavior.^{44,45} The deindividuated state traditionally was used to describe the mentality of individuals who comprised large riotous and hostile crowds (e.g., European soccer riots, mob violence, etc.). Social psychologist Phillip Zimbardo suggested that a number of antecedent variables, often characteristic of large crowds, lead to the deindividuated state. The psychosocial factors associated with anonymity, arousal, sensory overload, loss of responsibility, and mind-altering substances may lead to a loss of self-awareness, lessening of internal restraints, and a lack of concern for social or self-evaluation.⁴⁶

The factors associated with deindividuation also appear to be present during some online activities. For instance, Internet users are relatively anonymous and often use handles to further obscure their true identities. Many of the Websites, software programs, and multimedia files that typify the computing experience are sensory arousing and in some cases can be overstimulating. The Internet can be viewed as a large global crowd that individuals become submersed in when they go online. It is possible that the physical and psychological characteristics associated with the Internet that make it so appealing may also lead individuals to engage in antisocial and antinormative behaviors due to psychological feelings of immersion and deindividuation.⁴⁷

SOCIAL DISTANCE, ANONYMITY, AGGRESSION, AND COMPUTER CRIME 12 · 7

Deindividuation is brought about by an individual's loss of self-awareness, and psychological immersion into a crowd due to the aforementioned antecedents.⁴⁸ The aggressive, hostile, and antinormative actions of computer criminals may be linked to the denindividuated state. Zimbardo found that when participants were deindividuated, operationalized by anonymity, darkness, and loud music, they would administer higher levels of electric shocks to subjects, and for longer lengths of time, than individuated participants. Like Zimbardo's participants, computer criminals may be engaging in hostile and aggressive behavior due to deindividuation—that is, as a direct result of anonymity, subjective feelings of immersion, and the arousing nature of computer and Internet use.^{49,50}

12.4.3 Social Identity Theory and Computer Crime. Social psychologists Martin Lea, Tom Postmes, and Russell Spears have recently developed a social identity model of deindividuation effects (SIDE) to explain the influence of deindividuating variables on behaviors and attitudes during computer-mediated communications.^{51,52,53} According to social identity theory, an individual's self-concept resides on a continuum with a stable personal identity at one end and a social identity at the other. Depending on whether the social self, usually in group situations, or individual self is salient, the beliefs, norms, and actions associated with that particular self-concept will have the greatest influence on the individual's actions and attitudes.⁵⁴ When one of our social identities is salient, the norms associated with that group identity tend to guide and direct our behaviors.

According to the SIDE model, the isolation and visual anonymity that characterizes our online environment serves to enhance our social identities. This increase in social identification with a group may polarize our behaviors and attitudes toward the prevailing norms of that collective.^{55,56,57} Contrary to popular media stereotypes, computer criminals appear to have large social networks and frequently form groups and friendships with other like-minded individuals.⁵⁸ The use of handles and pseudonyms by these individuals combined with their physical isolation from each other may increase their aggressive and criminal tendencies depending on the overall norms associated with their online social groups. If the criminal collective values electronic intrusions and defacements more than programming and coding, then these behaviors will be exhibited to a greater extent by members who strongly identify with that group.

According to Henri Tajfel and John Turner's social identity theory (SIT), we tend to identify with ingroups, or those with whom we share common bonds and feelings of unity.⁵⁹ We have a bias toward our own group members and contrast them with outgroups, whom we perceive as different from those in our ingroup. While this ingroup bias or favoritism may benefit and protect our self-concepts, it may cause us to dislike and unfairly treat outgroup members (e.g., network administrators).

Communications researcher Hyung-jin Woo and colleagues used SIT to explain the motivations behind some Web page defacements.⁶⁰ SIT predicts that when groups are in competition for scarce resources, or feel threatened by outgroup members, there is a tendency for groups to respond aggressively toward each other. Ingroup members see improvements in collective self-esteem and enhanced feelings of group unity when they engage in attacks against outgroup members. Based on these predictions, Woo and colleagues hypothesized that computer criminals who are motivated by outgroup threats will express more aggressive and varied communication in Web page defacements than nonthreatened defacers.⁶¹

A content analysis of 462 defaced Web pages indicated that the majority of the defacements (71 percent) were classified as nonmalicious pranks. The most common

12 · 8 THE PSYCHOLOGY OF COMPUTER CRIMINALS

motivations behind these prankster attacks were to beat the system or its administrator, to gain peer recognition, to brag about accomplishments, and for romantic purposes. Twenty-three percent of the defacements were classified as militant attacks. The motivations behind these attacks were to promote groups associated with nationalism, ethnicity, religion, freedom of information, and anti-pornography.⁶² Consistent with the predictions made by SIT, the militant attacks were characterized by significantly more varied content, obscene language, insults, severe threats, and violent images. SIT may be useful in predicting the frequency and severity of attacks by computer criminals. Although only a minority (23 percent) of Web page defacements in this study resulted from intergroup conflict, those attacks were more severe in nature.⁶³

12.4.4 Social Learning Theory of Computer Crime. Criminologist Marc Rogers suggests that social learning theory (SLT) may offer some insight into the behavior of computer criminals.⁶⁴ According to psychologist Albert Bandura, individuals learn behaviors by observing the actions of others and their associated consequences.⁶⁵ SLT draws from B. F. Skinner's operant-conditioning model of learning where behaviors are learned or extinguished through schedules of reinforcement and punishment. However, Bandura's theory suggests that social learning occurs when an individual simply observes others' behaviors and reinforcements and forms a cognitive association between the two actions. Once the behavior is acquired, the learned actions are subject to external reinforcement, as in operant conditioning or in self-directed reinforcement.⁶⁶ According to the Social Structure and Social Learning model (SSSL), criminals learn deviant behaviors from their associations and subsequent imitations of deviant peers.⁶⁷ Through these peer associations, individuals learn to rationalize criminal behaviors. As Bandura suggests, if these criminal activities are rewarded (e.g., money, increased status, etc.), then the behavior will be strengthened.

Recently, there has been a growing amount of social and media attention focused on information security and computer criminals. Newspapers, magazines, and electronic news sources have reported thousands of incidents, interviews, and commentary related to computer crime. A number of these articles appear to glamorize hacking and the Internet underground.⁶⁸ The articles compare computer criminals to rock-and-roll superstars, James Bond-like spies, and international freedom fighters. Motion pictures and television shows like *The Matrix Trilogy*, *Mission Impossible*, *Hackers*, *Swordfish*, and *The X-Files* have all bestowed mythical qualities on rebellious computer criminals, while media outlets report computer criminals being recruited for high-paying government and industry jobs.⁶⁹ The media's glorification and glamorization of hacking and computer criminals, teaches some individuals that it pays to commit computer crime—at least, from a social learning perspective.

Many crimes involving computers are difficult to investigate and prosecute. The public learns via the media that computer criminals often are afforded fame and notoriety among their peers, and in the information security field, for their illegal activities. There are very few instances of computer criminals' being convicted and serving jail time as a consequence of their actions; usually the criminals are given light sentences. Their notoriety leads to media interviews, book and movie deals, even consulting and public speaking jobs.⁷⁰ Once an action is learned, SLT states that the behavior will be maintained via self-directed and external reinforcement. If computer criminals are rewarded for their illegal activities via the acquisition of knowledge and their elevated status in the hacker community, and the popular media continues to glamorize and focus on the positive consequences associated with computer crime, then the cost and prevalence of these illicit actions will continue to grow. Lending empirical support

INDIVIDUAL DIFFERENCES AND COMPUTER CRIMINALS 12 · 9

to the social learning model, criminologist Thomas Holt found that the four tenets of the SSSL model reliably predicted cyberdeviance in a college population. Students who associated with cybercriminals, imitated cyberdeviants, held morally ambiguous definitions regarding cybercrime, and had their computer-based deviance reinforced, were more likely to engage in online criminal activities.⁷¹

SLT may offer one explanation for the illegal behaviors of computer criminals, especially the marked increase in recent years.⁷² Instead of focusing on the supposedly positive consequences of computer crime, media outlets should stress the negative repercussions of computer crime for both the victims and the perpetrators. Social learning theorists would suggest modeling *appropriate* use of computers and the immediate negative ramifications of cyberdeviance as one element for fighting computer crime.

12.5 INDIVIDUAL DIFFERENCES AND COMPUTER CRIMINALS. Although situational factors can account for some of the behaviors of some computer criminals, one must not discount the impact of personality factors on their illicit activities. Attitudes and behaviors are often the product of both situational influences and individual personality traits.⁷³ It should be noted that there are few empirical studies that engage in a scientific examination of the personality traits of computer criminals, so without concrete evidence, the anecdotal claims regarding pathological traits of cybercriminals should be interpreted with caution. In addition, simply having traits that are *consistent* with a psychological disorder does not mean that one actually *has* the disorder.

12.5.1 Antisocial and Narcissistic Personalities. According to M. E. Kabay, some computer criminals exhibit insincerity and dishonesty in combination with superficial charm and an enhanced intellect, traits that are consistent with the *Diagnostic and Statistical Manual of Mental Disorders IV* (DSM-IV) criteria for *antisocial personality disorder*.⁷⁴ He also notes that some computer criminals commit their illegal behavior for little or no visible rewards despite the threat of severe punishment.

Another central characteristic of antisocial personality disorder is lack of clear insight by perpetrators regarding their behaviors.⁷⁵ Researchers have noted that computer criminals do not view their criminal actions as harmful or illegal.^{76,77} These criminals sometimes rationalize or externalize their behaviors by blaming the network administrators and software designers for not properly securing their computers and programs.

Computer crime researchers Eric Shaw, Keven Ruby, and Jerrold Post also have suggested that some computer criminals demonstrate personality characteristics consistent with some elements of *narcissistic personality disorder*.^{78,79} According to DSM-IV criteria, narcissistic individuals are attention seekers with an exaggerated sense of entitlement.⁸⁰ Entitlement is described as the belief that one is in some way privileged and owed special treatment or recognition.

Shaw and associates suggest that entitlement is characteristic of many “dangerous insiders,” or information technology specialists who commit electronic crimes against their own organizations.⁸¹ When corporate authority does not recognize the work or achievements of an employee to their satisfaction, the criminal insider seeks revenge via electronic criminal aggressions. Anecdotal evidence suggests that outside network intruders also may demonstrate an exaggerated sense of entitlement, as well as a lack of empathy for their victims, also characteristic of narcissistic personality disorder.

One self-identified computer criminal states, “we rise above the rest, and then pull everyone else up to the same new heights ... We seek to innovate, to invent. We, quite seriously, seek to boldly go where no one has gone before.”⁸² Narcissistic individuals

12 · 10 THE PSYCHOLOGY OF COMPUTER CRIMINALS

also frequently engage in rationalization to justify and defend their behaviors.⁸³ Computer criminal Toxic Shock writes, “We are misunderstood by the majority. We are misunderstood, misinterpreted, misrepresented. All because we simply want to learn. We simply want to increase the flow of knowledge, so that everyone can learn and benefit.”⁸⁴ Although it would be a mistake to generalize these hypotheses to the entire population without any empirical support, certain subsets of computer criminals may demonstrate characteristics that are consistent with aspects of both narcissistic and antisocial personality disorders.

12.5.2 Five-Factor Model of Personality and Computer Criminals. In one of the rare empirical studies looking at computer criminals, criminologist Marc Rogers examined the relationship between the five-factor model of personality and self-reported “criminal computer activity.”⁸⁵ The five-factor model formulated by psychologists Robert McCrae and Paul Costa in 1990 suggests that an individual’s personality can be accurately described using five core dimensions: extraversion (e.g., sociable), neuroticism (e.g., anxious), agreeableness (e.g., cooperative), conscientiousness (e.g., ethical), and openness to experience (e.g., nonconforming).⁸⁶

Rogers hypothesized that individuals engaging in computer crime would demonstrate higher levels of:

- exploitation,
- hedonistic morality,
- manipulation,
- antagonism,
- undirected behaviors,
- introversion,
- openness to experiences, and
- neuroticism

than noncriminals would. Three hundred eighty-one psychology students from an introductory psychology class were administered the computer-crime index (CCI), which is a self-report measure of computer-crime activity, along with measures of exploitation, manipulation, moral decision making, and a five-factor personality inventory. Contrary to the researcher’s expectations, individuals who committed computer crimes did not significantly differ from the nonoffenders on any of the five-factor personality measures. However, students who reported engaging in illegal computer activities did demonstrate more exploitative and manipulative tendencies.⁸⁷

In contrast, in a follow-up study Rogers did find that extraversion was a reliable predictor of computer crime behavior. The less extraverted an individual (i.e., more introverted), the more likely they were to engage in illicit computing activities. He suggests that the differences in populations between the two samples, the latter being a Canadian liberal arts college and the former being U.S. students from a technology program, may have contributed to the discrepant findings.⁸⁸ It is clear that further research will need to be conducted in this area in order to clear up the discrepant findings. One possible limitation in both of these studies is that the questionnaires were administered using a pencil-and-paper format, which assesses attitudes and traits when the participants’ *offline* identities are salient.⁸⁹ Internet researchers have long suggested that there is a distinct difference between our online and offline identities.⁹⁰ Therefore,

INDIVIDUAL DIFFERENCES AND COMPUTER CRIMINALS 12 · 11

had the participants been administered measures in an electronic format, when their online identities were more salient, Rogers might have found differing results.

12.5.3 Asperger Syndrome and Computer Criminals. Recently researchers have suggested a possible link between criminal hacking and a relatively new developmental disorder named *Asperger syndrome* (AS).^{91,92,93} AS is a disorder that resides at the mild end of the pervasive developmental disorder (PDD) spectrum, with classic autism at the more severe end. PDDs are characterized by primary developmental impairments in language and communication, social relations and skills, as well as repetitive and intense interests or behaviors.⁹⁴ Unlike autism, individuals who are diagnosed with AS have higher cognitive abilities and IQ scores ranging from normal to superior. Individuals with AS also have normal language and verbal skills, although there are noticeable deficits in social communication. Those diagnosed with AS typically have severe and systematic social skill impairment or underdevelopment, difficulties with interpersonal communication, and repetitive patterns of interests, behaviors, and activities.⁹⁵

According to clinical psychologist Kenneth Gergen, AS individuals must demonstrate social impairment. They may have a lack of desire or inability to interact with peers, and may engage in inappropriate or awkward social responses.⁹⁶ These individuals may have extremely limited or focused interests, and are prone to engage in repetitive routines. Although language development is often normal, these individuals may demonstrate unusual speech patterns (e.g., rate, volume, and intonation). Individuals with AS also may demonstrate clumsy motor behaviors and body language, as well as inappropriate facial expressions and gazing.⁹⁷

One of the noted features of AS that is anecdotally linked to computer criminals is the obsessive or extremely focused area of intellectual interest that the individuals demonstrate. Children with AS often show a preoccupation in areas such as math, science, technology, and machinery. They strive to learn and assimilate as much information as possible about their specialized interest. Researchers have indicated that their preoccupation may last well into adulthood, leading to careers associated with their intellectual interests.⁹⁸ Much of their social communication is egocentric, revolving around their obsessive interests, often leading to strained and difficult social interactions. Although children with AS desire normal peer interaction, their egocentric preoccupations, lack of appropriate social behaviors, and difficulties empathizing with others often leave them frustrated, misunderstood, teased, and sometimes ostracized.⁹⁹

Researchers have noticed similarities in the characteristics associated with AS and traits stereotypically associated with computer hackers.¹⁰⁰ Tony Atwood, an Australian clinical psychologist, suggests that some computer hackers may have a number of characteristics that are associated with AS.¹⁰¹ He notes that many diagnosed AS patients are more proficient at computer programming languages than social language, and that the intellectual challenges that are presented by restricted computers and networks may override the illegal nature of their actions. AS has been used as a successful defense in at least one landmark U.K. court case.¹⁰²

Based on over 200 personal interviews with computer criminals, cybercrime expert Donn Parker reports finding significant similarities between AS sufferers and criminal hackers. Many of the computer criminals that Parker interviewed demonstrated the social awkwardness, atypical prosody, and lack of social empathy during social interactions that are characteristic of AS.¹⁰³ Anecdotal evidence suggests that computer hackers often have an obsessive interest in technology and computers, similar to that seen in individuals with AS, that forms a salient component of both their individual

12 · 12 THE PSYCHOLOGY OF COMPUTER CRIMINALS

and social identities. Due to their egocentric preoccupations, many computer hackers often feel misunderstood and frustrated in face-to-face social situations.

Although the Autism Diagnostic Interview is the most common assessment tool, a number of self-report instruments have been developed to assist in screening disorders on the autism spectrum. To date, sociologist Bernadette Schell has conducted the only empirical study examining Asperger syndrome in a self-identified hacker population.¹⁰⁴ Schell administered the 50-item Autism-Spectrum Quotient (AQ) Inventory to 136 attendees at various well-known hacker conferences (e.g., Defcon¹⁰⁵) between 2005 and 2007. Previous research using the AQ found that diagnosed AS individuals' mean scores were 35.8 compared to a mean score of 16.4 for a control population. Schell found that the mean score for conference attendees in her study was 19.7. However, 11 percent of males in her sample and 1.5 percent of females had mean AQ scores of 32 or higher. Although this study is limited by participants' self-identification as hackers, it does suggest that proposed link between hacker culture and AS may be tenuous at best.¹⁰⁶

To date there has been no clear empirical evidence to suggest a link between AS and computer crime. There is no evidence whatsoever to suggest that Asperger syndrome causes computer hacking. In fact, most sufferers of AS have been characterized as being extremely honest and lawful citizens.¹⁰⁷ It would be a mistake to assume that all computer hackers are suffering from AS or that every AS sufferer is a computer hacker. Characteristics of AS appear more common in computer *hackers* (i.e., those who explore and tinker with computers and technology), rather than in *crackers* who break into computers or use them for illegal activities. At present, there is still no single all-encompassing personality profile that applies to all computer criminals. In fact, many feel that it is inappropriate to try to create a single personality profile that applies to all computer criminals.

12.5.4 Internet Abuse and Computer Crime. Although *technological addiction* is not a disorder recognized by the American Psychological Association (APA), researchers have suggested that some individuals appear to demonstrate disturbed computer use that is similar to other recognized disorders like compulsive gambling or impulse control disorders. Technological addiction is characterized by:

- Excessive use of a particular technology (usually in reference to computers and Internet usage)
- Preoccupation with the technology
- Systematic increase in use
- Failed attempts to curb one's use
- Feelings of malaise and irritation when not using the technology
- Interference with social and professional pursuits due to excessive technology use¹⁰⁸

Information security researchers Kent Anderson and Jerrold Post suggest that some computer criminals appear to have symptoms indicative of potentially pathological computer use.^{109,110} Research by sociologist Bernadette Schell found that hacker conference participants indicated that they spend on average 24 hours a week on hacker-related activities.¹¹¹ Furthermore, Anderson reports that cybercriminals will work for 18 or more hours a day on their computers trying to gain unauthorized access to one

INDIVIDUAL DIFFERENCES AND COMPUTER CRIMINALS 12 · 13

single computer system with little or no external reward for doing so. He also mentions an instance where one U.S. judge even attempted to sentence a computer criminal to psychological treatment for his compulsive computer use.¹¹²

In their interviews with a number of self-identified computer criminals, sociologists Paul Taylor and Tim Jordan indicated that many of their interviewees report experiencing a thrill or rush that isn't comparable to anything that they experience in their real-world interactions when engaging in illegal activities.¹¹³ A number of their respondents reported feelings of depression, anxiety, and impaired social functioning when they are away from their computers. Taylor and Jordan suggest that these abuse-like characteristics may also be combined with feelings of compulsion regarding computers and new technologies. However, the researchers indicated that these compulsive or obsessive-like characteristics may be as much a function of the information technology (IT) field as they are personality traits.¹¹⁴ Unlike most disciplines, the IT field is in a constant state of rapid change. To maintain a level of professional expertise and competence in this area, one must devote a good deal of time and resources to monitoring and adapting to this revolutionary field. This need to keep up with the rapidly changing discipline, combined with the euphoric feelings that some experience when committing illegal activities, may increase the likelihood of technological abuse.

Personality theorists state that for some computer criminals, committing electronic crimes produces an experience similar to that of a chemically induced high. Some computer criminals may commit illegal acts because of the euphoric rush they receive from their actions. Information security researcher August Bequai compares the actions of computer criminals to electronic joyriding.¹¹⁵ Gaining unauthorized access and usage to a computer network is, for these people, similar to that of taking a car on a joyride. One computer cracker interviewed by computer crime researcher Dorothy E. Denning described hacking as "the ultimate cerebral buzz." Other crackers have commented that they received a rush from their illegal activities that felt as if their minds were working at accelerated rates. Some computer criminals have suggested that the euphoric high stems from the dangerous and illegal nature of their activities.¹¹⁶ Lending empirical support to this idea, criminologist Michael Bachmann found that self-identified computer criminals have a heightened propensity to engage in risk behaviors compared to the general population. Computer criminals have compared the feelings they receive from their illegal intrusions and attacks to the rush that is felt by participating in extreme sports like rock climbing and skydiving.¹¹⁷

Researchers have found that some experienced computer users also report sometimes experiencing an altered psychological state known as *flow* while engaging in their technological pursuits. Flow is a psychological state that results in feelings of fulfillment and overall positive affect.¹¹⁸ When individuals become absorbed in a task that matches their skill set, at times they may not be consciously aware of the passage of time or of the differences between the undertaking and their identity. In their study looking at the experience of flow in self-identified computer criminals, psychologists Alexander Voiskounsky and Olga Smyslova found that both inexperienced and highly competent criminals report high levels of flow.¹¹⁹ For the inexperienced criminals, this flow experience may lead them to limit themselves to low-level challenges, and they may remain in this novice stage for a significant amount of time. More experienced crackers who also experience flow may leave the hacking domain once they are no longer presented with suitable challenges for their abilities. Conversely, they may systematically increase their illegal pursuits in attempts to reacquire the flow state, contributing to technological abuse.¹²⁰

12 · 14 THE PSYCHOLOGY OF COMPUTER CRIMINALS

Physical and psychological tolerance can occur when increased amounts of a substance or an activity are needed in order to obtain a *high* or euphoric rush. Tolerance is common in hard drug users who find themselves using increasing amounts of a substance to achieve their original euphoric states. Anecdotal evidence suggests that computer criminals may go through a similar stage of evolution, with each step leading to increased dangerous and riskier behaviors. Many cybercriminals begin by pirating and cracking the copy protection algorithms of software programs. When the *warez* (pirated software) scene loses its thrill, they migrate to chat-room or IRC harassment. The individuals may then begin launching damaging DoS attacks against servers and defacing Websites to obtain that initial rush that originated with simple *warez* trading. As in a substance abuse, the initial euphoric psychological states and resulting tolerance associated with excessive computer use may explain why some computer criminals repeatedly engage in illicit activities, even after they have been caught and punished.

12.6 ETHICS AND COMPUTER CRIME. Researchers have suggested that computer criminals may have an underdeveloped sense of ethics—a moral immaturity—that contributes to their illegal activities.^{121,122,123,124} Because of this ethical immaturity, criminal hackers may think that many of their illegal actions are in fact ethical or beneficial to some degree. Many computer criminals feel that they are ethically entitled to have access to any and all information regardless of legal ownership. Most of these individuals also feel that it is morally right to use inactive computer processing power and time, regardless of who owns the computer system. Computer criminals do not feel that breaking into a computer network should be viewed in the same light as breaking into an individual’s house. Often, computer criminals rationalize their illegal activities and justify their behaviors by blaming the victims for not securing their computer networks properly.¹²⁵

Most computer criminals are adolescents, which may account for the underdeveloped sense of ethics in the community.¹²⁶ Computer scientist Brian Harvey suggests that due to the relative lack of experience and guidance with the computing environment compared to the real world, teenagers and adolescents may be operating at lower levels of moral functioning when online compared to their interactions and decisions in the real world.¹²⁷ Similarly, information security specialist Ira Winkler suggests that computer hackers, because of their generally young age, do not fully understand the repercussions associated with their actions. They also may demonstrate an underdeveloped or complete lack of empathy for their victims.¹²⁸ Computer criminals fail to fully realize the consequences of their electronic intrusions into computer networks. The adolescents do not fully comprehend that their mere presence on a computing network could potentially cost companies thousands of dollars, as well as cost systems administrators their jobs.¹²⁹

Sociologist Orly Turgeman-Goldschmidt further suggests that computer criminals may view their behaviors as nothing more than a new form of social entertainment.¹³⁰ They see their electronic intrusions as a game that provides them with excitement and thrills. The Internet serves an unlimited playground or social center where “netizens” are able to develop new games and forms of social activities. Computer criminals may view their illegal activity as nothing more than fun and thrill seeking, a new form of entertainment that is carried out on an electronic playground.¹³¹

According to Winkler, many criminal hackers learn and develop a sense of computing ethics from their online and offline peers.¹³² In other words, unlike most instruction on morality and ethics that stems from a responsible adult, computing ethics are socially

ETHICS AND COMPUTER CRIME 12 · 15

learned from other adolescents. Computer crackers learn the rules of hacking and computing from elder statesmen in the hacking community who may be no more than a few years older than themselves. Often, in today's society, the younger generation is more knowledgeable about technology than older adults. When adolescents have problems or need guidance in ambiguous situations, many times their parents or other adult role models are unable to offer them the necessary guidance and assistance. Therefore, the youngsters may seek knowledge from their peers, who may or may not offer them the most ethical or wise advice.¹³³ In support of this notion, Vincent Sacco and Elia Zureik found that students who viewed illicit computing behaviors as ethical increased the reported likelihood that they would engage in such actions. Computer crime was least reported when the behavior was seen as being more unethical.¹³⁴

While using the Internet, children are constantly making sophisticated judgments without appropriate adult supervision. Adolescents do not have the same ethical maturity that adults have, yet while using the Internet they are given as much power, authority, and responsibility as ethically mature adults.¹³⁵ Neil Patrick, the leader of one particularly malicious group of phone-system hackers known as the 414s, stated that he did not know that his hacking was illegal or unethical. In fact, asked when, if ever, he began to question the ethics of his actions, Patrick stated that ethics never came into his mind until the Federal Bureau of Investigation agents were knocking on his front door. Patrick and the other young members of the 414s did not see anything wrong with their actions. To them, breaking into proprietary telecommunications networks was more of a game or challenges rather than a criminal act. They saw nothing ethically wrong with their actions.¹³⁶

Psychologist Lawrence Kohlberg developed a three-level theory to explain normal human moral development.

- The first level deals with avoiding punishments and obtaining rewards.
- The second level emphasizes social rules.
- The third level emphasizes moral principles.

Each of his three levels contains two stages that an individual passes through during adolescence on the way to adult moral development.¹³⁷ Computer criminals appear to be operating in the lower three phases of Kohlberg's model: the two stages comprising level 1 and the first stage in level 2. The moral judgments of computer criminals appear to be determined by a need to satisfy their own needs and to avoid disapproval and rejection by others.

In his empirical research on moral development, Rogers found that self-identified computer criminals relied more on hedonistic decision making rather than internal or social morality-based choices, compared to a noncriminal student population. Computer criminals do not appear to be aware of or concerned with the third level of moral development, where moral judgments are motivated by civic respect and one's own moral conscience.¹³⁸ Computer criminals may be functioning at the third level of moral development in the physical world, a level appropriate for teens and adults, and may simultaneously be functioning at lower levels of moral development when their online identities are salient.¹³⁹ Computer criminals acting at these lower levels of morality may be naively engaging in their illegal activities to satisfy their own curiosity and to gain the approval of their peers, without considering larger moral implications of their behaviors.

12 · 16 THE PSYCHOLOGY OF COMPUTER CRIMINALS

According to Shaw and associates, there is a notable lack of ethical regulation and education in organizations, schools, and homes regarding proper computing behavior.¹⁴⁰ Computer criminals who lack ethical maturity fail to realize that their digital actions are sometimes just as damaging as physical aggression. Cybervandals do not see the immediate repercussions of their actions because of the physical distance and lack of social presence in computer-mediated interactions. This ethical immaturity is partially a result of the technology-enhanced knowledge gap between young computer users and their parents. A hacker whose handle was *The Mentor* wrote in the 1986 “Hacker Manifesto,”

This is our world now ... the world of the electron and the switch, the beauty of the baud. We make use of a service already existing without paying for what could be dirt-cheap if it wasn't run by profiteering gluttons, and you call us criminals. We explore ... and you call us criminals. We seek after knowledge ... and you call us criminals. We exist without skin color, without nationality, without religious bias ... and you call us criminals. You build atomic bombs, you wage wars, you murder, cheat, and lie to us and try to make us believe it's for our own good, yet we're the criminals.

Yes, I am a criminal. My crime is that of curiosity. My crime is that of judging people by what they say and think, not what they look like. My crime is that of outsmarting you, something that you will never forgive me for.¹⁴¹

In real-world situations, when parents and teachers strive to instill responsible ethics in adolescents, young adults become capable of making informed decisions regarding ethical dilemmas. However, the same adolescents who demonstrate ethical behavior in the physical world may be ethically bereft in cyberspace, partly due to the lack of adult guidance and instruction. Anecdotal evidence suggests that in today's society, adolescents recreationally use, and are more familiar with, computers and the Internet than their parents. These young adults often learn about Internet-related behaviors and attitudes on their own, or via peer-to-peer interaction. Adolescents are socialized on the Internet by other adolescents, which may lead to a *Lord of the Flies* scenario, where children construct social rules and guidelines to govern their behaviors.¹⁴² These socially constructed norms and guidelines may be both morally and ethically different from real-world norms. Kabay has argued that technological change can take two to three generations for integration of new moral codes into society; by this reasoning, young adults who are growing up with increased awareness of civil behavior in cyberspace will be teaching their own children more appropriate rules of behavior from the earliest ages of the next generation.¹⁴³

12.7 CLASSIFICATIONS OF COMPUTER CRIMINALS. For both ordinary and abnormal behaviors, it is difficult to find one theoretical perspective that can account for every behavior in a given situation. Attitudes and behaviors are the product of the combined influence of an individual's personality and the current social situation. No single theory or theoretical perspective can account for the various types of computerized crimes and the criminals who engage in these activities. There are also many types of computerized crimes, ranging from trading pirated software to cyberterrorism, and many types of computer criminals, ranging from the novice password cracker to the industrial spy. Any theory that would account for the behavior of computer criminals would have to consider, first, the type of illegal activity the person was engaged in and, second, the type of cybercriminal category that the individual falls into.¹⁴⁴ Computer criminals are by nature paranoid and secretive agents who exist in a similar community. They use handles to conceal their true identities and, except for annual hacker

CLASSIFICATIONS OF COMPUTER CRIMINALS 12 · 17

conventions or local meetings, seldom interact with each other in the real world. Therefore, it is difficult for researchers to identify and categorize the various subgroups that exist.

The term *computer hacker* has been both overused and misused as a way of classifying the entire spectrum of computer criminals.¹⁴⁵ The motivations and actions of one subgroup of computer criminals may be entirely different from those of a second group; therefore, it is imperative that in any psychological analysis of computer criminals, the various subcategories be taken into consideration. Many theories have attempted to account for the motivations and behaviors of computer criminals as a whole when the theorists actually were referring to one specific subgroup in the underground culture.¹⁴⁶ Computer criminals are a heterogeneous culture; therefore, one single theory or perspective cannot sufficiently explain all their actions. The fact that researchers traditionally have treated computer criminals as a homogeneous entity has limited the validity and generalizability of their findings. Even researchers who have taken into account the heterogeneous nature of the computing underground have had difficulty with experimental validity. Experimenters have allowed participants to use their own self-classification schemes or attempted to generalize the results of a single subgroup to the entire underground culture.¹⁴⁷

12.7.1 Early Classification Theories of Computer Criminals. Over the past few decades, several researchers have attempted to develop a categorization system for individuals who engage in various forms of computer crime.^{148,149,150,151,152} A comprehensive review of this research is beyond the scope of this chapter. For an extensive review, see Rogers's analysis and development of a new taxonomy for computer criminals.¹⁵³ Bill Landreth, a reformed computer cracker, was one of the earliest theorists to develop a classification scheme for computer criminals.¹⁵⁴ His system divided criminals into five categories based on their experience and illegal activities.

1. The *novice* criminals have the least experience with computers and cause the least amount of electronic disruption from their transgressions. They are considered to be tricksters and mischief-makers; for example, AOL users who annoy chat-room members with text floods and DoS-like *punting* programs that crash AOL sessions using specific font or control code strings.
2. The *students* are electronic voyeurs. They spend their time browsing and exploring unauthorized computer systems.
3. The *tourists*, according to Landreth, commit unauthorized intrusions for the emotional rush that results from their actions. This subgroup of computer criminals is similar to Bequai's electronic joyriders.¹⁵⁵ The tourists are thrill-seekers who receive a cerebral buzz from their illegal behaviors.
4. The *crashers* are malicious computer criminals. This subgroup is composed of the darkside criminals that Kabay refers to.¹⁵⁶ The crashers will crack into networks and intentionally delete and destroy data and cause denials of service.
5. Landreth's final classification of computer criminal is the *thieves*. Criminals who fit into this category commit their illegal actions for monetary gain. These individuals are equivalent to the dangerous insiders that Post, Shaw, and Ruby analyzed.¹⁵⁷ Thieves may work alone, or they may be under contract from both foreign and domestic corporations and governments. Examples include the Russian Business Network or RBN.¹⁵⁸

12 · 18 THE PSYCHOLOGY OF COMPUTER CRIMINALS

Former Australian Army intelligence analyst Nicholas Chantler conducted one of the few empirical examinations of computer criminals and their culture.¹⁵⁹ Published in 1995, this survey-based study attempted to gain a deeper understanding of the underground culture as well as to develop a categorization system for cybercriminals. Chantler posted questionnaires to bulletin board systems (BBSs), Usenet newsgroups, and chat rooms owned or frequented by computer criminals. An analysis of the data yielded five primary attributes—criminal activities, hacking prowess, motivations, overall knowledge, and length of time hacking—that Chantler used to create three categories of computer criminals: *lamers*, *neophytes*, and *elites*.¹⁶⁰

1. *Lamers* have the least technical skill and they have been engaged in their illegal activities for the shortest period of time. This group of criminals is primarily motivated by revenge or theft of services and property.
2. *Neophytes* are more mature than lamers. They are more knowledgeable than the previous category and engage in illegal behaviors in pursuit of increased information.
3. Members of the *elite* group have the highest level of overall knowledge concerning computers and computer crime. They are internally motivated by a desire for knowledge and discovery. They engage in illegal activities for the intellectual challenge and for the thrill they receive from their criminal behaviors.

According to Chantler, the largest proportion of computer criminals at that time, 60 percent, fell into the neophyte category. Thirty percent of computer criminals fell into the elite category, while 10 percent were lamers.¹⁶¹

Information security analyst Donn Parker developed a seven-level categorization scheme for computer criminals.^{162,163} He formalized his scheme, through years of interaction and structured interviews with computer criminals, into the following categories:

1. **Pranksters** are characterized by their mischievous nature.
2. **Hacksters** are motivated by curiosity and a quest for knowledge. Pranksters and hacksters are the least malicious computer criminals.
3. **Malicious hackers** are motivated by a need for disruption and destruction. They receive pleasure from causing harm to computer systems and financial loss to individuals.
4. **Personal problem solvers** commit illegal activities for personal gain. Problem solvers, the most common type of computer criminal according to Parker, resort to crime after failing in legitimate attempts to resolve their difficulties.
5. **Career criminals** engage in their illegal cyberbehaviors purely for financial gain.
6. **Extreme advocates** have strong ties to religious, political, or social movements. Recently, these types of cybercriminals have been dubbed “hacktivists,” a combination of computer hackers and activists.
7. **Malcontents, addicts, and irrational individuals** comprise the final category in Parker’s scheme. Individuals in this category usually are suffering from some form of psychological disorder (e.g., antisocial personality disorder).

CLASSIFICATIONS OF COMPUTER CRIMINALS 12 · 19

12.7.2 Rogers's New Taxonomy of Computer Criminals. After an extensive review of past categorization theories, Rogers has advanced an updated continuum of computer criminals based on his previous work.¹⁶⁴ This continuum comprises eight categories based primarily on the criminals' motivations and technological prowess:

1. **Novice (NV)** criminals have the least amount of technical knowledge and skill. Members of this category are relatively new to the scene and use prewritten and precompiled scripts and tools to commit their computerized crimes. They are primarily motivated by the thrill of lawbreaking and making a name for themselves in the underground.
2. **Cyber Punk (CP)** is the group that most fits the traditional stereotype of hacker. Members of this category are slightly more advanced than the novices. These criminals have the ability to create basic attack scripts and programs. Cyber Punks' typical behaviors include Web page defacements, DDoS attacks, carding, and telecommunication fraud. They are motivated by a need for attention, fame, and monetary gain, usually attained by parlaying their crimes into lucrative jobs and book deals. Winkler suggests that the majority of computer criminals fall into either the cyber punk or newbie categories. He estimates that between 35,000 and 50,000 computer criminals, well over 90 percent of their total estimated number, fall into these categories, whom he dubs clueless.¹⁶⁵
3. **Internals (IN)** consist of disgruntled workers or former workers who hold information technology positions in an organization. Members of this category have an advantage over external attackers due to their job and status within the corporation. Research indicates that internals are responsible for the majority of computer crimes and associated financial loss. Their motivation is typically based on revenge for some perceived wrong (e.g., termination, passed over for a promotion).^{166,167} These internal or malicious insiders are examined more fully later in Chapter 13 of this *Handbook*.
4. **Petty Thieves (PT)** are traditional criminals who have turned to technology as a way of keeping up with the times. These individuals are career criminals whose motivation is primarily financial gain from stealing from banks, corporations, and individuals.
5. **Old Guard (OG)** are computer criminals with advanced technical knowledge and skill. These individuals are responsible for writing many of the exploit programs (e.g., stack overflows, rootkits, etc.) that are used by the less knowledgeable novice and cyberpunk crackers in their cyberattacks; however, they are not criminals in the traditional sense. This group has an underdeveloped sense of ethics regarding privacy and intellectual and personal property and engages in behaviors consistent the traditional hacker ethic and ideology described by Levy.¹⁶⁸ Their illegal behaviors are motivated by a quest for knowledge, curiosity, and intellectual stimulation.
6. **Virus Writers (VW)** do not fit neatly into Roger's Taxonomy primarily due to the lack of research on this group of individuals. The demographics and motivations of virus writers are examined more fully later in this chapter.
7. **Professional Criminals (PC)** are traditionally older and more knowledgeable about technology than the previous categories. Members of these categories may be former government and intelligence operatives who are motivated by

12 · 20 THE PSYCHOLOGY OF COMPUTER CRIMINALS

financial gain. They may often have access to advanced technology and can be adept at industrial espionage. According to Rogers, PCs may be comprised of ex-intelligence agents and are one of the most dangerous types of computer criminals. Their motivation is primarily large-scale financial gain.

8. **Information Warriors** are highly skilled employees who conduct coordinated attacks against information systems in an attempt to cripple or destabilize the infrastructure. This group may be motivated by allegiance and patriotism. We examine this group more fully in the following section.

Based on these eight categories, Rogers has created a baseline circumplex model of computer criminals that aims to further classify computer criminals. Rogers' circumplex groups computer criminals in a circular diagram using two continua: technological skill and motivation. Using this circumplex, the eight categories can be further subdivided, following future empirical work, to more accurately represent individual subgroups of the computer underground.

12.7.3 Hacktivists and Cyberterrorists: Hacking for a Cause. With the emerging protests against the economically advantaged 1 percent and the political upheaval in Middle East, new trends in cybercrime have emerged in the form of a revitalized hacktivist ethos. *Hacktivism* is defined as cause-based hacking for social, political, patriotic, or religious purposes. In general, the perpetrators engage in technological dissent in order to promote freedom of speech and fair distribution of wealth and to combat censorship. Hacktivists will often rally behind a symbol, logo, or flag, whether they are religious, political, or emblematic. The term was first used by groups like Electronic Disturbance Theatre (EDT) and the Cult of the Dead Cow (cDc) in the mid-1990s to establish digital protests and *hack-ins* where the groups would launch coordinated DoS attacks and mass Web page defacements as a means of online protest.^{169,170,171,172} Paralleling the global “Occupy” protest movements, many of these modern hacktivist groups have no central leadership or core philosophy, making their actions and targets difficult to predict. They present an increased danger due to their tendency to single out and target high-profile people and organizations to achieve maximal exposure for their group, cause, and interests. Hacktivists are not motivated by financial gain, but instead prefer to publicize their cause through the harm and embarrassment of their victims.¹⁷³ The methods employed by modern hacktivists have evolved from simple DoS attacks and Web defacements to massive amounts of government, corporate, and personal data theft. The number of hacktivist attacks in 2011 surpassed all of the previous year’s attacks combined that were attributed to social and political motivations.¹⁷⁴

According to Rogers, *cyber-terrorists*, along with *professional criminals*, may present the most danger to individuals and organizations. A cyber-terrorist is defined as, “an individual who uses computer or network technology to control, dominate, or coerce through the use of terror in furtherance of political or social objectives.”¹⁷⁵ Following the September 11, 2001, attacks in the United States, the term cyberterrorism has been frequently overused and misused by the media. An individual who stumbles upon a vulnerable .mil or .gov Website and decides to deface it should not be considered a cyberterrorist unless there is a premeditated intent to cause panic and fear, with the object of obtaining some social end.

Although most of the Web page defacements and network attacks labeled by the media as cyberterrorism would not fit Rogers’s definition, that is not to say that the Internet will not be used as means for terrorist acts in the future.¹⁷⁶ Spurred on by

CLASSIFICATIONS OF COMPUTER CRIMINALS 12 · 21

the specter of cyberterrorism and claims by hackers that they could cripple the Internet in 30 minutes, governments and corporations are investing millions into research aimed at protecting the global computing infrastructure from cyberattacks. Rogers suggests that the relative anonymity of the perpetrator, and the multitude of potential targets that could be simultaneously attacked, makes the Internet a very appealing target for terrorists' actions that will likely be exploited in the near future.¹⁷⁷

Perhaps the most infamous case of cyberterrorism occurred in 2010 when researchers analyzed Stuxnet, one of the most sophisticated electronic worms ever discovered in the wild. Stuxnet was created with the singular goal of creating physical damage to centrifuges in one of Iran's nuclear enrichment facilities. Although no one has taken credit for the Stuxnet worm, computer security researchers suggest that it may have been cyberwarriors working for the United States and/or Israeli governments.¹⁷⁸ Theorists suggest that new breeds of computer criminals, dubbed cyberterrorists and hacktivists, are emerging and are motivated by political or social ideologies related to information freedom, nationalism, ethnic pride, and warfare.^{179,180}

12.7.4 Dangerous/Malicious Insiders (DI/MI). Dr. Eric Shaw and associates classify computer criminals into two categories: outside intruders and dangerous insiders.^{181,182,183} The researchers focus on the critical IT insiders who are typically programmers, technical support staff, networking operators, administrators, consultants, and temporary workers in organizations. Malicious insiders (MI) are a subgroup of such employees who are motivated by greed, revenge, problem resolution, and ego gratification. Shaw and colleagues estimate that the theft of trade secrets costs more than \$250 billion annually for U.S. businesses alone.

Shaw's research, based on corporate surveys and hundreds of investigations by the U.S. Secret Service and Carnegie Mellon University's Computer Emergency Response Team Coordination Center (CERT-CC), has attempted to compile an initial profile of dangerous insiders. Their critical-pathway approach identifies psychological and situational precursors that may predispose an employee to engage in insider espionage¹⁸⁴:

- 1. Personal Predisposition:** medical/psychiatric problems, reduced social skills, previous violations of the law or business ethics, social or professional ties with competitors or adversaries
- 2. Personal Life Stressors:** financial burden, relationship problems, medical issues, legal issues
- 3. Professional Stressors:** demotions, failed promotions, poor performance review, transfer, supervisor disagreements, looming layoffs
- 4. Concerning Behaviors:** workplace violations, conflict, intellectual property (IP) disagreements
- 5. Maladaptive Organizational Responses:** failure to detect, investigate, appreciate, and appropriately respond to concerning behaviors of employees

The generic profile for an MI is a male in his late thirties who has a job in a technology-related area within the company. The majority of thefts occur when an MI has accepted a job elsewhere and within a month of starting their new job. Moore, as cited by Shaw, classifies MIs into two categories, the *Entitled Disgruntled Thief* and the *Machiavellian Leader*.

12 · 22 THE PSYCHOLOGY OF COMPUTER CRIMINALS

- The Entitled Disgruntled Thief typically engages in IP theft due to professional stressors. They steal information that they directly worked on or helped to develop. These individuals feel entitled to the information and typically will use it to help them acquire a new job or to enhance their performance at their already acquired position.
- The actions of the Machiavellian Leader are more premeditated. They are driven more by personal ambition rather than personal or professional stressors. These individuals may also recruit coworkers to assist them in their thefts.

According to Shaw and coauthors, these dangerous insiders typically have introverted personalities.¹⁸⁵ They demonstrate a preference for solitary intellectual activities over interpersonal interaction. Members of this subgroup may have had numerous personal and social frustrations that have hindered their interpersonal interactions and contributed to their antiauthoritarian attitudes.

The malicious subgroup of critical information technologists has been characterized as having a loose or immature sense of ethics regarding computers. The dangerous insiders rationalize their crimes by blaming their company or supervisors for bringing any negative consequences on themselves. They feel that any electronic damage they cause is the fault of the organization for treating them unfairly. The researchers also note that many insiders identify more with their profession than with the company for which they work. This heightened identification with the profession undermines an insider's loyalty to an organization. This reduced loyalty is evidenced by the high turnover rates of jobs in the IT industry. According to Shaw and associates, the unstable bond between insiders and their organizations creates undue tension with regard to security practices and IP rights.¹⁸⁶

Researchers also have suggested that dangerous insiders are characterized by an increased sense of entitlement and hostility toward organizational authority. According to Shaw and coauthors, when an unfulfilled sense of entitlement is combined with previous hostility toward authority, malicious acts or revenge against the organization are typical.^{187,188,189}

12.7.5 Virus Creators. Unlike more traditional forms of computer crime, there has been very little research examining the motivations and behaviors of malware writers, usually referred to as *virus writers*.¹⁹⁰ Despite the name, the group writes other forms of malware such as worms and Trojans. The limited number of research reports on this particular subgroup of computer criminal has relied primarily on one-on-one interviews and surveys, in an effort to understand the actions and motivations of virus creators.^{191,192,193,194}

Using case studies and multiple interviews, researchers Andrew Bissett and Geraldine Shipton examined the factors that influence and motivate virus writers.¹⁹⁵ The researchers suggest that it is difficult to generalize their findings to all virus creators because of the limited published literature and research regarding virus writers. Virus creators appear to demonstrate conscious motivations for their potentially destructive actions that are similar to the motivations of traditional computer criminals. The coders create and distribute their software for reasons of nonspecific malice, employee revenge, ideological motives, commercial sabotage, and information warfare.

Bissett and Shipton's review of anti-virus expert Sarah Gordon's interview with Dark Avenger reveals some of the motivations behind one of the most notorious virus writers.¹⁹⁶ They suggest that Dark Avenger consistently denies responsibility for his

CLASSIFICATIONS OF COMPUTER CRIMINALS 12 · 23

creations and, like traditional computer criminals, engages in victim blaming. Dark Avenger states that it is human stupidity, not the computer, that spreads viruses. The virus writer also appears to self-identify with his malicious code. Dark Avenger seems to project his persona onto viruses in a process called projective identification. During the interview, Dark Avenger stated that the United States could prevent him from entering the country, but it is unable to stop his viruses. Dark Avenger also attempted to justify creating destructive viruses by commenting that most personal computers did not store data of any value, and therefore his malicious programs were not doing any real harm. Similar to the motivations spurring dangerous insiders, the researchers suggest that Dark Avenger creates malicious viruses because he is envious of the work and achievements of other computing professionals. In the interview, Dark Avenger commented that he hates it when people have more powerful computing resources than he does, especially when the individuals do not use the resources for anything that he deems constructive.¹⁹⁷

Sarah Gordon has also examined the ethical development of several virus writers using surveys and structured interviews.^{198,199,200} Her initial four case studies involved an adolescent virus writer, a college-age virus writer, a professionally employed virus writer, and an ex-virus writer. The interviews revealed that all four individuals appeared to demonstrate normal ethical maturity and development consistent with Kohlberg's previously reviewed stage theory.²⁰¹ Gordon suggests that there appear to be many different reasons why individuals create and distribute viruses, including boredom, exploration, recognition, peer pressure, and sheer malice.²⁰²

Gordon suggests that the virus underground in the mid and late 1990s was populated by a second generation or next generation of virus writers whose skill and ability at virus construction is comparable to that of the old school, original virus writers.²⁰³ On the surface, these second-generation creators maintain a public façade that suggests that they are extremely cruel, obnoxious, and more technologically advanced than the previous generation. The next-generation virus writers appear to be more aware of the ethical responsibilities surrounding virus creation and distribution; however, the exact definition of "responsible virus creation and distribution" varies from individual to individual. Many of these next-generation virus writers have considerable technical skill and are motivated by the challenge to defeat the malware countermeasures implemented by antivirus vendors.²⁰⁴

According to Gordon, another group of virus creators populating the virus underground in the 1990s was composed of "new-age" virus writers.²⁰⁵ These individuals are motivated by current trends, such as political activism, virus exchange, freedom of information, and challenges to write the most destructive or sophisticated virus, as opposed to technical exploration. These virus writers are motivated by boredom, intellectual curiosity, mixed messages surrounding the legality of virus creation, and increased access to ever-more-powerful technological resources. Gordon suggests that these new-age virus writers may be older and wiser than the second or next-generation creators. They are very selective as to who has access to their creations, and they do not share their findings or accomplishments with members outside their respective group. Unlike the next-generation creators, new-age virus writers will not stop or grow out of writing viruses, as they are most likely already adults. They will continue to write and distribute more sophisticated viruses in part due to the mixed messages concerning the ethical nature of virus creation propagated by the popular media, academia, and popular freedom of information zeitgeist.²⁰⁶

Researchers have stated that we must be careful not to view virus creators as a homogeneous group.^{207,208} Instead, we must monitor the virus-exchange community

12 · 24 THE PSYCHOLOGY OF COMPUTER CRIMINALS

while pursuing in-depth case histories that may aid in our understanding of virus writers. Education about the ethical nature of virus creation and distribution, and about the repercussions associated with malicious code, may attenuate these potentially destructive activities.

12.8 RECOMMENDATIONS. Psychological theories offer various explanations that may influence criminal activities on the Internet. These situational influences may interact with various individual personality traits to further contribute to illegal behaviors. The current task for researchers is to untangle these personality and situation influences on electronic behavior. They must determine what situations and characteristics are influencing the various subgroups of computer criminals. There is no simple explanation as to why computer criminals engage in hostile and destructive acts. The answer lies in a complex mixture of factors that depends on the social environment and individual personality factors. There are numerous types of computer criminals ranging from script kids to the professional criminal, each with varying personalities and motivators. The interaction of personality variables and environmental factors will determine how a computer criminal reacts in any given situation.

One underlying theme in reducing the overall prevalence of computer crime has been to remove the tangible and psychological reward system that currently surrounds the culture of the Internet underground. Criminal law researchers A. S. Dalal and Raghav Sharma suggest that the information security field has established a pattern of rewarding criminal hackers for their exploits by offering them employment opportunities that further undermines law enforcement efforts to combat computer crime.²⁰⁹ Preventing cybercriminals from profiting from their transgressions via lucrative jobs and book deals may result in a socially learned deterrent to engaging in cybercrime. Thomas Holt and colleagues found strong support for the social learning model of cybercrime in their research: those who received reinforcement in the form of encouragement, praise, and resources from peers and authority figures were more likely to engage in cybercrime activities. This model suggests that removing this reinforcement would significantly reduce the amount and frequency of cyberdeviance, especially in the case of criminals that fall into Rogers' *novice* and *cyberpunk* categories.^{210,211}

Cautioning parents, teachers, and bosses about the dangers of praising adolescents for cyberdeviance and instead instituting consistent punishments should also help to mitigate computer crime for novices. Also, a more proactive strategy that parents and educators could take would be to engage in moral or ethics training on the use of computers and technology before children begin independently using technology on a regular basis. Then, when confronted with ambiguous decisions while using the computers, children would reflexively rely on adults' moral advice and teachings instead of their peers'. Scholastic Incorporated demonstrated in a recent survey that 48 percent of elementary and middle school students did not consider computer hacking to be a crime.²¹² Recently, the U.S. Department of Justice (DoJ) formed a cybecitizen awareness program in an effort to educate parents, teachers, and children about ethical and unethical computing practices. The program seeks to educate individuals through ethics conferences, multimedia presentations, and speaking engagements at schools around the country. This program and others like it may aid to increase moral responsibility and ethical behaviors of adolescents on the Internet.²¹³

Dalal and Sharma also suggest that punishments take into account the motivation and type of cybercrime that is committed, instead of using a rigid approach to sentencing that precludes adapting punishment to the nature of the criminal. For instance, longer jail terms may be a suitable deterrent for criminals who are motivated by financial

RECOMMENDATIONS 12 · 25

gain or malice (e.g., internals, petty thieves, and malicious insiders), but might not be warranted for those who are motivated by intellectual challenges and curiosity (e.g., old guard hackers).

Another problem may be the relatively light sentences that are typically handed down for computer crimes, and the perpetrators' perceptions that they are unlikely to be caught or prosecuted for their transgressions. One study found that the perceived gains achieved by illegal hacking far outweighed the potential costs associated with digital crime. Furthermore, even though computer criminals feel that the punishments for their crimes are severe enough, their feelings that there is a very low probability that they'll be caught minimizes the effects of punishment severity as a deterrent.²¹⁴ Even still, the notoriety that a computer criminal gets from being handed a jail term can be parlayed into a lucrative job once the individual has paid the price for their criminal endeavors.

A more proactive approach to reducing the motivation to engage in illicit activities while using computers could be to establish legal hacking networks for novices to experiment with and explore. For novices and even cyberpunks, safe hacking spaces where they have unlimited access to networks and programs may reduce the psychological reactance that results from strictly enforced computing access and resources. These legal hacking networks would also serve to redirect the users' focus on the technical aspect of computing and traditional hacking endeavors rather than trying to crack open the system. By redirecting adolescents' curiosity toward traditional hacking and technical pursuits into an open network, the users may be less inclined to engage in criminal pursuits.

To combat the threat of malicious insiders, Shaw suggests that businesses also adopt a proactive approach to threat mitigation.²¹⁵ Developing more comprehensive employee-screening methods that include social-media reviews, background checks, substance abuse tests, honesty/ethical tests, and psychological screenings may serve to mitigate the risks of IP theft. Additionally, a more cost-effective and time-efficient method of reducing IP theft would be increasing employee awareness of IP nondisclosure agreements (NDAs) and ensuring that they are aware of the consequences of violating these agreements. Increasing employees' and managers' awareness of intellectual property rights and risk factors, as well as instituting a systematic reporting system for potential violation could mitigate some of the risks of IP theft. Shaw reports that the majority of IP theft was detected by nontechnical means (e.g., employees noticing suspicious activity or similar products being marketed by competitors).²¹⁶

This chapter summarized research on the psychological motivations of computer criminals. The cybercrime landscape is vast in scope and populated by subgroups of computer criminals with their own patterns of motivations, goals, attitudes, and behaviors. Psychological theory and research is one step in trying to generate effective countermeasures to combat the problem of cybercrime.

Although there is no shortage of theories and anecdotal evidence to account for why cybercriminals engage in network intrusions, there is a marked lack of empirical evidence that serves to test and validate these hypotheses. Information security professionals, criminologists, and psychologists must find ways to begin jointly developing, testing, and examining these ideas in order to find concrete solutions for the problem of cybercrime. Such effort will be facilitated by government and private-sector funding to support research into the psychological and social dynamics of the Internet underground. The information security industry as a whole would do well to support these applied scientific efforts.²¹⁷

12 · 26 THE PSYCHOLOGY OF COMPUTER CRIMINALS

12.9 FURTHER READING

- Jaishankar, K., ed. *Cyber Criminology: Exploring Internet Crimes and Criminal Behavior*. Boca Raton, FL: CRC Press, 2011.
- Kirwan, G., and A. Power, eds. *The Psychology of Cyber Crime: Concepts and Principles*. Hershey, PA: IGI Global, 2011.

12.10 NOTES

1. Symantec, *Internet Security Threat Report: 2011 Trends*, Vol. 17, Symantec Website, April 2012, www.symanteccloud.com/en/us/mlireport/b-istr_main-report_2011_21239364.en-us.pdf
2. Marc Rogers, “The Development of a Meaningful Hacker Taxonomy: A Two Dimensional Approach,” *CERIAS Technical Report*, 43 (2005): 1–9.
3. National Institute of Justice, “Electronic Crime Research and Development,” 2006, www.ncj.gov/topics/crime/internet-electronic/
4. Symantec, *Internet Security*, 2012.
5. For more information about computer crimes, see Chapter 2 in this Handbook.
6. Verizon, “Data Breach Investigations Report,” Verizon Website, 2012, www.verizonbusiness.com/resources/reports/rp_data-breach-investigations-report-2012_en_xg.pdf
7. Symantec, *Internet Security*, 2012.
8. For more information about information warfare, see Chapter 14 in this *Handbook*.
9. Douglas Campbell, “A Detailed History of Terrorist and Hostile Intelligence Attacks Against Computer Resources,” 1992. Available at www.syneca.com/publications.htm
10. Donn Parker, “How to Solve the Hacker Problem,” *Journal of the National Computer Security Association*, No. 5 (1994), pp. 4–8.
11. Jerrold Post, “The Dangerous Information Systems Insider: Psychological Perspectives,” (1998). Available email: jmpost@pol-psych.com
12. Parker, “Hacker Problem,” 1994 (pp. 4–8).
13. Paul Taylor, *Hackers: Crime in the Digital Sublime*, (London & New York: Routledge, 1999).
14. Raoul Chiesa, Stefania Ducci, and Silvio Ciappi, *Profiling Hackers: The Science of Criminal Profiling as Applied to the World of Hacking* (Boca Raton, FL: Auerbach Publications, 2008).
15. Orly Turgeman-Goldschmidt, “Hackers’ Accounts: Hacking as a Social Entertainment,” *Social Science Computer Review*, 23, no. 1 (2005): 8–23.
16. Tim Jordan and Paul Taylor, “A Sociology of Hackers,” *Sociological Review*, 46, no. 4 (1998): 757–780.
17. Taylor, *Hackers*, 1999.
18. Turgeman-Goldschmidt, “Hackers’ Accounts,” 2005.
19. Taylor, *Hackers*, 1999.
20. Taylor, *Hackers*, 1999.
21. Chiesa, Ducci, and Ciappi, *Profiling Hackers*, 2008.
22. Chiesa, Ducci, and Ciappi, *Profiling Hackers*, 2008.

NOTES 12 · 27

23. Chiesa, Ducci, and Ciappi, *Profiling Hackers*, 2008.
24. Jordan and Taylor, "A Sociology," 1998.
25. Hyung-jin Woo, Yeora Kim, & Joseph Dominick, "Hackers: Militants or Merry Pranksters? A Content Analysis of Defaced Web Pages," *Media Psychology*, 6, no. 1 (2004): 63–83.
26. David G. Myers, *Social Psychology*, 8th ed. (New York: McGraw-Hill, 2005).
27. Myers, *Social Psychology*, 2005.
28. Woo, Kim, and Dominick, "Hackers: Militants or Merry Pranksters?" 2004.
29. For more information on denial-of-service attacks, see Chapter 18 in this Handbook.
30. Stanley Milgram, *Obedience to Authority* (New York: Harper & Row, 1974).
31. Albert Bandura, "Selective Activation and Disengagement of Moral Control," *Journal of Social Issues*, 46, (1990): 27–46.
32. Marc Rogers, "Modern-Day Robin Hood or Moral Disengagement?" 1999. http://homes.cerias.purdue.edu/~mkr/moral_doc.pdf
33. Turgeman-Goldschmidt, "Hackers' Accounts," 2005.
34. Jordan and Taylor, "A Sociology," 1998.
35. Woo, Kim, & Dominick, "Hackers: Militants or Merry Pranksters?" 2004.
36. Rogers, "Modern-Day Robin Hood," 1999.
37. Randall Young, Lixuan Zhang, and Victor R. Prybutok, "Hacking into the Minds of Hackers," *Information Systems Management*, 24 (2007): 281–187.
38. Sara Kiesler, Jane Siegel, and Timothy McGuire, "Social Psychological Aspects of Computer-mediated Communication," *American Psychologist*, 39 (1984): 1123–1134.
39. Sara Kiesler and Lee Sproull, "Group Decision Making and Communication Technology," *Organizational Behavior and Human Decision Processes*, 52 (1992): 96–123.
40. Kiesler and Sproull, "Group Decision Making," 1992.
41. Kiesler, Siegel, and McGuire, "Social Psychological Aspects," 1984.
42. Kiesler and Sproull, "Group Decision Making," 1992.
43. Turgeman-Goldschmidt, "Hackers' Accounts," 2005.
44. Phillip Zimbardo, "The Human Choice: Individuation, Reason, and Order Versus Deindividuation, Impulse, and Chaos," Nebraska Symposium on Motivation, No. 17 (1969): 237–307.
45. Edward Diener, "Deindividuation, Self-Awareness, and Disinhibition," *Journal of Personality and Social Psychology*, 37 (1979): 1160–1171.
46. Zimbardo, "The Human Choice," 1969.
47. Kiesler, Siegel, and McGuire, "Social Psychological Aspects," 1984.
48. Zimbardo, "The Human Choice," 1969.
49. Zimbardo, "The Human Choice," 1969.
50. M. E. Kabay, "Anonymity and Pseudonymity in Cyberspace: Deindividuation, Incivility and Lawlessness versus Freedom and Privacy," Annual Conference of the European Institute for Computer Anti-virus Research (EICAR), 1998, revised 2001, www.mekabay.com/overviews/anonpseudo.pdf

12 · 28 THE PSYCHOLOGY OF COMPUTER CRIMINALS

51. Tom Postmes, Russell Spears, and Martin Lea, "Social Identity, Normative Content, and 'Deindividuation' in Computer-Mediated Groups," In *Social Identity: Context, Commitment, Content*, ed. N. Ellemers, R. Spears, and B. Doosje (Oxford: Blackwell, 1999), 164–183.
52. Tom Postmes, Russell Spears, and Martin Lea, "Breaching or Building Social Boundaries? SIDE-Effects of Computer Mediated Communication," *Communication Research*, 25 (1998): 689–715.
53. Stephen D. Reicher, Russell Spears, and Tom Postmes, "A Social Identity Model of Deindividuation Phenomena," *European Review of Social Psychology*, 6 (1995): 161–198.
54. Henri Tajfel and John C. Turner, "The Social Identity Theory of Inter-Group Behavior," In *Psychology of Intergroup Relations*, ed. S. Worchel and L. W. Austin (Chicago: Nelson-Hall, 1986), 2–24.
55. Postmes, Spears, and Lea, "Social Identity," 1999.
56. Postmes, Spears, and Lea, "Breaching or Building," 1998.
57. Reicher, Spears, and Postmes, "A Social Identity Model," 1995.
58. Jordan and Taylor, "A Sociology," 1998.
59. Tajfel and Turner, "The Social Identity Theory," 1986.
60. Woo, Kim, & Dominick, "Hackers: Militants or Merry Pranksters?" 2004.
61. Woo, Kim, & Dominick, "Hackers: Militants or Merry Pranksters?" 2004.
62. Woo, Kim, & Dominick, "Hackers: Militants or Merry Pranksters?" 2004.
63. Woo, Kim, & Dominick, "Hackers: Militants or Merry Pranksters?" 2004.
64. Rogers, "Modern-Day Robin Hood," 1999.
65. Albert Bandura, "The Social Learning Perspective: Mechanisms of Aggression." In *Psychology of Crime and Criminal Justice*, ed. H. Toch (New York: Holt, Rinehart & Winston, 1979).
66. Bandura, "Social Learning Perspective," 1979.
67. Thomas J. Holt, George W. Burruss, & Adam M. Bossler, "Social Learning and Cyber-Deviance: Examining the Importance of a Full Social Learning Model in the Virtual World," *Journal of Crime & Justice*, 33, no. 2 (2010): 31–61.
68. M. E. Kabay, "Totem and Taboo in Cyberspace: Integrating Cyberspace into Our Moral Universe," *Journal of the National Computer Security Association* (1996): 4–9. www.mekabay.com/ethics/totem_taboo_cyber.pdf
69. Chiesa, Ducci, and Ciappi, *Profiling Hackers*, 2008.
70. Kabay, "Totem and Taboo," 1996.
71. Holt, Burruss, and Bossler, "Social Learning and Cyber-Deviance," 2010.
72. Bandura, "Social Learning Perspective," 1979.
73. Myers, *Social Psychology*, 2005.
74. American Psychiatric Association, *Diagnostic and Statistical Manual of Mental Disorders*, 4th ed. (DSM-IV) (Washington, DC: American Psychiatric Association, 1994).
75. American Psychiatric Association, *DSM-IV*, 1994.
76. Kabay, "Totem and Taboo," 1996.
77. Eric D. Shaw, Keven Ruby, and Jerrold Post, "The Insider Threat to Information Systems: The Psychology of the Dangerous Insider," *Security Awareness Bulletin*, 2 (1998): 1–10.

NOTES 12 · 29

78. Kabay, "Totem and Taboo," 1996.
79. Shaw, Ruby, and Post, "Insider Threat," 1998.
80. American Psychiatric Association, *DSM-IV*, 1994.
81. Shaw, Ruby, and Post, "Insider Threat," 1998.
82. Toxic Shock, "Another View of Hacking: The Evil That Hackers Do," *Computer Underground Digest*, 2 (1990). Available: <http://cu-digest.org/CUDS2/cud2.html>
83. Shaw, Ruby, and Post, "Insider Threat," 1998.
84. Toxic Shock, "Another View," 1990.
85. Marc Rogers, "Understanding Criminal Computer Behavior: A Personality Trait and Moral Choice Analysis," 2003, Marc Rogers's Website, <http://homes.cerias.purdue.edu/~mkr/CPA.doc>
86. Robert R. McCrae and Paul T. Costa, Jr. "Personality Trait Structure as a Human Universal," *American Psychologist*, 52. (1997): 509–516.
87. Rogers, "Criminal Computer Behavior," 2003.
88. Marcus Rogers, Kathryn Seigfried, and Kirti Tidke, "Self-Reported Computer Criminal Behavior: A Psychological Analysis," *Digital Investigation*, 3S (2006): 116–120.
89. Postmes, Spears, & Lea, "Social Identity," 1999
90. Sherry Turkle, "Identity Crisis," in *Life on the Screen: Identity in the Age of the Internet* (New York: Simon & Schuster, 1995): 255–269.
91. M. J. Zuckerman, "Hacker Reminds Some of Asperger Syndrome," *USA Today*, March 3, 2001, www.usatoday.com/news/health/2001-03-29-asperger.htm
92. Suelette Dreyfus, "Cracking the Hackers' Code," *The Sydney Morning Herald*, August 8, 2002, <http://smh.com.au/articles/2002/08/20/1029114072039.html>
93. Bernadette H. Schell and June Meluychuk, "Female and Male Hacker Conferences Attendees: Their Autism-Spectrum Quotient (AQ) Scores and Self-Reported Adulthood Experiences," in *Corporate Hacking and Technology-Driven Crime: Social Dynamics and Implications*, ed. Thomas J. Holt and Bernadette H. Schell (Hershey, PA: IGI Global, 2011): 144–166.
94. American Psychiatric Association, *DSM-IV*, 1994.
95. American Psychiatric Association, *DSM-IV*, 1994.
96. Stephen Bauer, "Asperger Syndrome," 2001, <http://aspergersyndrome.org/>
97. Bauer, "Asperger Syndrome," 2001.
98. Peter Smith, "The Cybercitizen Partnership: Teaching Children Cyber Ethics," Cybercitizen Partnership, 2000, www.cybercitizenship.org/ethics/whitepaper.html
99. Bauer, "Asperger Syndrome," 2001.
100. Zuckerman, "Hacker Reminds Some," 2001.
101. Dreyfus, "Cracking the Hackers' Code," 2002.
102. Dreyfus, "Cracking the Hackers' Code," 2002.
103. Zuckerman, "Hacker Reminds Some," 2001.
104. Schell and Meluychuk, "Female and Male Hacker AQ Scores," 2011.
105. DEFCON Website, www.defcon.org
106. Schell and Meluychuk, "Female and Male Hacker AQ Scores," 2011.

12 · 30 THE PSYCHOLOGY OF COMPUTER CRIMINALS

107. Bauer, "Asperger Syndrome," 2001.
108. Mark Griffiths, "Internet Addiction: Does It Really Exist?" In *Psychology and the Internet: Intrapersonal, Interpersonal and Transpersonal Applications*, ed. J. Gackenbach, (New York: Academic Press, 1998): 61–75.
109. K. E. Anderson, "International Intrusion: Motives and Patterns," 1994, www.aracnet.com/~kean/Papers/paper.shtml
110. Jerry Post, Eric Shaw, and Keven Ruby, "Information Terrorism and the Dangerous Insider," Paper presented at the InfowarCon, 1998, Washington, D.C.
111. Schell and Meluchuk, "Female and Male Hacker AQ Scores," 2011.
112. Schell and Meluchuk, "Female and Male Hacker AQ Scores," 2011.
113. Jordan and Taylor, "A Sociology," 1998.
114. Jordan and Taylor, "A Sociology," 1998.
115. August Bequai, *Technocrimes* (Lexington, MA: Lexington Books, 1987).
116. Dorothy E. Denning, "Concerning Hackers Who Break into Computer Systems," *CPSR* (1990), <http://cyber.eserver.org/hackers.txt>
117. Michael Bachmann, "The Risk Propensity and Rationality of Computer Hackers," *International Journal of Cyber Criminology*, 4 (2010): 643–656.
118. Csikszentmihalyi, Mihaly, *Flow* (Harper Perennial Modern Classics, 2008).
119. Alexander E. Voiskounsky and Olga V. Smyslova, "Flow-Based Model of Computer Hackers' Motivation," *CyberPsychology & Behavior*, 6, no. 2 (2003): 171–161.
120. Voiskounsky & Smyslova, "Flow-Based Model," 2003.
121. Rogers, "Modern-Day Robin Hood," 1999.
122. Post, Shaw, and Ruby, "Information Terrorism," 1998.
123. Denning, "Concerning Hackers," 1990.
124. Ira Winkler, "Why Hackers Do the Things They Do?" *Journal of the National Computer Security Association*, 7 (1996): 12.
125. Woo, Kim, and Dominick, "Hackers: Militants or Merry Pranksters?" 2004.
126. Denning, "Concerning Hackers," 1990.
127. Brian Harvey, "Computer Hacking and Ethics," Brian Harvey's Website, 1998, www.cs.berkeley.edu/~bh/hackers.html
128. Winkler, "Why Hackers Do," 1996.
129. Harvey, "Computer Hacking and Ethics," 1998.
130. Turgeman-Goldschmidt, "Hackers' Accounts," 2005.
131. Turgeman-Goldschmidt, "Hackers' Accounts," 2005.
132. Winkler, "Why Hackers Do," 1996.
133. Winkler, "Why Hackers Do," 1996.
134. Vincent Sacco and Elia Zureik, "Correlates of Computer Misuse: Data from a Self-Reporting Sample," *Behavior & Information Technology*, 9 (1990): 353–369.
135. Harvey, "Computer Hacking and Ethics," 1998.
136. Harvey, "Computer Hacking and Ethics," 1998.
137. Myers, *Social Psychology*, 2005.
138. Rogers, "Modern-Day Robin Hood," 1999.

NOTES 12 · 31

139. Harvey, “Computer Hacking and Ethics,” 1998.
140. Shaw, Ruby, and Post, “Insider Threat,” 1998.
141. Mentor, “The Hacker Manifesto,” *Phrack Magazine*, January 8, 1986, www.mithral.com/~beberg/manifesto.html
142. William Golding, *Lord of the Flies* (London: Faber and Faber, 1954).
143. Kabay, “Totem and Taboo,” 1996.
144. Denning, “Concerning Hackers,” 1990.
145. Marc Rogers, “A New Hacker Taxonomy,” Marc Rogers’s Website, 2000, http://homes.cerias.purdue.edu/~mkr/hacker_doc.pdf
146. Rogers, “A New Hacker Taxonomy,” 2000.
147. Rogers, “A New Hacker Taxonomy,” 2000.
148. Rogers, “Meaningful Hacker Taxonomy,” 2005.
149. A. N. Chantler, “Risk: The Profile of the Computer Hacker.” Ph.D. dissertation, Curtin University of Technology, 1995.
150. Rogers, “A New Hacker Taxonomy,” 2000
151. Bill Landreth, *Out of the Inner Circle* (Redmond, WA: Microsoft Books, 1985).
152. Donn Parker, *Fighting Computer Crime: A New Framework for Protecting Information* (New York: John Wiley & Sons, 1998).
153. Rogers, “Meaningful Hacker Taxonomy,” 2005.
154. Landreth, *Out of the Inner Circle*, 1985.
155. Bequai, *Technocrimes*, 1987.
156. Kabay, “Totem and Taboo,” 1996.
157. Post, Shaw, and Ruby, “Information Terrorism,” 1998.
158. See B. Guinen and M. E. Kabay, “The Russian Cybermafia: Beginnings,” 2011, [www.mekabay.com/nwss/866_russian_cybercrime_\(guinen\)_part_1.pdf](http://www.mekabay.com/nwss/866_russian_cybercrime_(guinen)_part_1.pdf); M. E. Kabay and B. Guinen, “The Russian Cybermafia: Boa Factory & CarderPlanet,” 2011, [www.mekabay.com/nwss/867_russian_cybercrime_\(guinen\)_part_2.pdf](http://www.mekabay.com/nwss/867_russian_cybercrime_(guinen)_part_2.pdf); and B. Guinen and M. E. Kabay, “The Russian Cybermafia: RBN & the RBS WorldPay Attack,” 2011, [www.mekabay.com/nwss/868_russian_cybercrime_\(guinen\)_part_3.pdf](http://www.mekabay.com/nwss/868_russian_cybercrime_(guinen)_part_3.pdf)
159. Chantler, “Risk,” 1995.
160. Chantler, “Risk,” 1995.
161. Chantler, “Risk,” 1995.
162. Rogers, “A New Hacker Taxonomy,” 2000.
163. Parker, *Fighting Computer Crime*, 1998.
164. Rogers, “Meaningful Hacker Taxonomy,” 2005.
165. Winkler, “Why Hackers Do,” 1996.
166. Eric D. Shaw, “The Role of Behavioral Research and Profiling in Malicious Cyber Insider Investigations,” *Digital Investigation*, 3 (2006): 20–31.
167. Eric D. Shaw and Harley V. Stock, “Behavioral Risk Indicators of Malicious Insider Theft of Intellectual Property: Misreading the Writing on the Wall,” Symantec White Paper, 2011, https://www4.symantec.com/mktginfo/whitepaper/21220067_GA_WP_Malicious_Insider_12_11_dai81510_cta56681.pdf
168. Steven Levy, *Hackers* (New York: Dell Publishing, 1984).

12 · 32 THE PSYCHOLOGY OF COMPUTER CRIMINALS

169. Samantha Murphy, "Culture Lab: Inside the Hacktivists Revolution," *Culture-Lab | New Scientist*, April 12, 2012, www.newscientist.com/blogs/culturelab/2012/04/samantha-murphy-contributorit-is-the.html
170. Thomas J. Holt, "The Attack Dynamics of Political and Religiously Motivated Hackers," In *Cyber Infrastructure Protection*, ed. Tarek Saadawi and Louis Jordan Jr., (Strategic Studies Institute, 2011): 159–180, www.strategicsciencesinstitute.army.mil/pdffiles/PUB1067.pdf
171. Marc Rogers, "The Psychology of Cyber-Terrorism," In *Terrorists, Victims, and Society: Psychological Perspectives on Terrorism and its Consequences*, ed. A. Silke (London: Wiley & Sons, 2003).
172. Alexander Gostev and Costin Raiu, "Kaspersky Security Bulletin. Malware Evolution 2011," Securelist Website, 2012, www.securelist.com/en/analysis/204792217/Kaspersky_Security_Bulletin_Malware_Evolution_2011
173. Holt, "Attack Dynamics," 2011.
174. Verizon, "Data Breach Investigations Report," 2012.
175. Rogers, "Psychology of Cyber-Terrorism," 2003.
176. Rogers, "Psychology of Cyber-Terrorism," 2003.
177. Rogers, "Psychology of Cyber-Terrorism," 2003.
178. Symantec, *Internet Security*, 2012.
179. Holt, "Attack Dynamics," 2011.
180. Rogers, "Psychology of Cyber-Terrorism," 2003.
181. Shaw, Ruby, and Post, "Insider Threat," 1998.
182. Shaw, "Role of Behavioral Research," 2006.
183. Shaw and Stock, "Behavioral Risk Indicators," 2011.
184. Shaw and Stock, "Behavioral Risk Indicators," 2011.
185. Shaw, Ruby, and Post, "Insider Threat," 1998.
186. Shaw, Ruby, and Post, "Insider Threat," 1998.
187. Shaw, Ruby, and Post, "Insider Threat," 1998.
188. Shaw, "Role of Behavioral Research," 2006.
189. Shaw and Stock, "Behavioral Risk Indicators," 2011.
190. Rogers, "Meaningful Hacker Taxonomy," 2005.
191. Andrew Bissett and Geraldine Shipton, "Some Human Dimensions of Computer Virus Creation and Infection," *International Journal of Human Computer Studies*, 52, No. 5 (2000), pp. 1071–5819.
192. Sarah Gordon, "The Generic Virus Writer," 4th International Virus Bulletin Conference, Jersey, U.K. (September 1994). <http://vxheaven.org/lib/static/vdat/epgenvr2.htm>
193. Sarah Gordon, "The Generic Virus Writer II," 6th International Virus Bulletin Conference, Brighton, U.K. (September 1996). www.research.ibm.com/antivirus/SciPapers/Gordon/GVWII.html
194. Sarah Gordon, "Virus Writers: The End of Innocence?" Presented at the 10th International Virus Bulletin Conference (September 2000). <http://vxheaven.org/lib/asg12.htm>
195. Bissett and Shipton, "Some Human Dimensions," 2000.
196. Bissett and Shipton, "Some Human Dimensions," 2000.

NOTES 12 · 33

197. Bissett and Shipton, "Some Human Dimensions," 2000.
198. Gordon, "Generic Virus Writer," 1994.
199. Gordon, "Generic Virus Writer II," 1996.
200. Gordon, "Virus Writers," 2000.
201. Myers, *Social Psychology*, 2005.
202. Gordon, "Generic Virus Writer II," 1996.
203. Gordon, "Generic Virus Writer II," 1996.
204. Gordon, "Generic Virus Writer II," 1996.
205. Gordon, "Generic Virus Writer II," 1996.
206. Gordon, "Generic Virus Writer II," 1996.
207. Rogers, "Meaningful Hacker Taxonomy," 2005.
208. Gordon, "Virus Writers," 2000.
209. A. S. Dalal, & Raghav Sharma, "Peeping into a Hacker's Mind: Can Criminological Theories Explain Hacking?" *ICFAI Journal of Cyber Law* 6, no. 4 (2007): 34–47.
210. Rogers, "Meaningful Hacker Taxonomy," 2005.
211. Holt, Burruss and Bossler, "Social Learning and Cyber-Deviance," 2010.
212. Gordon, "Virus Writers," 2000.
213. Gordon, "Virus Writers," 2000.
214. Young, Zhang, and Prybutok, "Hacking into the Minds," 2007.
215. Shaw and Stock, "Behavioral Risk Indicators," 2011.
216. Shaw and Stock, "Behavioral Risk Indicators," 2011.
217. M. E. Kabay, "Time for Industry to Support Academic INFOSEC," M. E. Kabay's Website, 2004. www.mekabay.com/opinion/endowed_chairs.pdf

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 13

THE INSIDER THREAT

Gary L. Tagg, CISSP

13.1 INTRODUCTION	13·1	13.3 MITIGATING THE INSIDER THREAT	13·7
13.2 THREATS FROM INSIDERS	13·2	13.3.1 System and Asset Inventories	13·8
13.2.1 How Common Are Insider Attacks?	13·2	13.3.2 Data Loss Prevention (DLP)	13·8
13.2.2 Examples of Insider Attacks	13·3	13.3.3 Internal Honeypots	13·10
13.2.3 Types of Insider Threats	13·5		
13.2.4 Internet-Based Systems	13·6		
13.2.5 Service Provider Threats	13·7	13.4 CONCLUDING REMARKS	13·10
13.2.6 System Administration Threats	13·7	13.5 FURTHER READING	13·11
		13.6 NOTES	13·11

13.1 INTRODUCTION. An insider is someone who has been given a role within an organization and has access to premises and/or internal systems and information. There are many types of roles; some are core to the business and performed by employees, while others are non-core and contracted out to service providers such as cleaning, maintenance, or information technology (IT). The categories of insiders may be classified as:

Current staff work directly for and under the control of the organization's management. This category includes employees as well as temporary staff, contractors, and consultants. Most of these people are located on the organization's premises and are connected to the internal network with access to internal information.

Departing staff represent one of the highest risks to an organization. These people consist of employees who have resigned or are planning to do so, temporary staff, contractors or consultants coming to the end of their contract, as well as all the people whose employment or services are being ended by the organization. These people still have access to internal information, and may be motivated to take that information with them when they leave, or to commit sabotage in revenge for perceived wrongs.

Former staff are those who are no longer employed by or providing services to the organization. This group still has insider knowledge, and without mitigating controls, they can do substantial damage long after they have left the organization. Former staff can be highly motivated to attack their former employers.

13 · 2 THE INSIDER THREAT

Service-providers: Organizations have many roles that are necessary for the smooth running of the business but that are not core to the company's mission. Examples of these roles are cleaning services, maintenance, and IT. Over the last 20 years, service providers have steadily moved up the service stack to perform core functions as well. It is now common for service providers to perform business operations and even to be the first point of contact with customers. Although not all service providers will need access to premises or internal systems, there are many roles that do.

The key distinction between service providers and staff is that staff (who can be provided by service providers) are working under the control of the organization's management and subject to its policies and procedures; in contrast, service providers are providing a specified service, and all personnel are the responsibility of the service provider. This situation increases risk to the organization, as it has little control over service-provider staff, but the threats are the same as from internal staff.

Partners: Organizations are partners when they work together on a business venture. Access to information goes up to senior business leaders within the partners. Partnerships may be short or long term, and partners may be competitors at the same time as being partners. This situation creates risk for the organization, which has to share information relevant to the partnership but not other internal information.

The term *partner* is sometimes used where service providers are performing key functions for an organization and success of the service is essential. For these key contracts, it is good practice to treat the service provider more like a partner than a vendor, but for insider threat purposes they are service providers.

13.2 THREATS FROM INSIDERS. Successful organizations need to have people in many different roles, with the primary distinction between business roles (business management, sales, customer services, and product design) and support/infrastructure roles (IT, finance, logistics, human resources, and the like).

The people in each role need access to information and systems to perform their role, and the key point to make at this stage is that the impact an insider attack makes on the organization is related to the insider role. A salesman leaving an organization to join a competitor will have had access to customer- and product-related information, whereas someone in product design is likely to have access to valuable intellectual property on current and future products. Customer service center staff are likely to have access to customers' personally identifiable information (PII) that can be used to commit identity fraud.

IT is a high-risk area, and is essential to the efficient running of the business. Most of an organization's information is stored in its IT systems, and without proper controls, IT administrators could have access to everything.

13.2.1 How Common Are Insider Attacks? According to the 2011 CyberSecurity Watch Survey,¹ 21 percent of electronic crime events were conducted by insiders, 58 percent by outsiders, and 21 percent unknown. Of the insider events, 76 percent were dealt with internally without legal action or law enforcement, which is likely a reason that external attacks are more often in the news.

The 2012 U.K. Information Security Breaches Survey² provides some interesting statistics.

- 6 percent of the responding 447 organizations reported staff sabotage of systems, with some organizations experiencing incidents on a weekly basis

THREATS FROM INSIDERS 13 · 3

- 19 percent of large organizations reported that staff used systems to commit theft or fraud, and this figure had doubled since the survey in 2010
- With regard to other incidents caused by insiders, large organizations were more likely to experience incidents, with 82 percent reporting incidents versus 45 percent for small organizations
- 61 percent of large organizations reported unauthorized access to systems or data
- 45 percent reported a breach of data protection laws or regulations
- 36 percent reported misuse of confidential information
- 47 percent reported loss or leakage of confidential information

The 2013 U.K. Information Security Breaches Survey³ includes the following relevant information for this chapter (quoting directly from the Executive Summary—note U.K. spelling):

- 36% of the worst security breaches in the year were caused by inadvertent human error (and a further 10% by deliberate misuse of systems by staff)
- 57% of small businesses suffered staff-related security breaches in the last year (up from 45% a year ago)
- 17% of small businesses know their staff broke data protection regulations in the last year (up from 11% [in the 2012 report])
- 14% of large organisations had a security or data breach in the last year relating to social networking sites
- 9% of large organisations had a security or data breach in the last year involving smartphones or tablets
- 4% of respondents had a security or data breach in the last year relating to one of their cloud computing services
- 4% of the worst security breaches were due to portable media bypassing defences

13.2.2 Examples of Insider Attacks. To get an understanding of the types of insider events, a number of organizations maintain databases and publish results. The FBI maintains an Insider Threat page on their Website⁴ containing a list of prosecuted insider theft cases. A short summary of these cases follows:

- Retired research scientist conspired with current and former employees to steal trade secrets from his former employer, and sell them to companies in China.
- Employee working for two different U.S. companies stole trade secrets from both companies which were used to benefit Chinese universities.
- A research scientist stole trade secrets from her employer and made them available for sale through her own company.
- An employee stole customer and employee lists, contract information, and other trade secrets to provide to a foreign government, but instead gave them to an undercover FBI agent.
- A computer programmer working for a financial firm copied proprietary software during his last few days at the company.
- An employee who was fired had kept copies of trade secrets. These trade secrets were then sold to a rival company.

13 · 4 THE INSIDER THREAT

- An employee stole and attempted to sell trade secrets that provided everything needed to start a competing business.
- Over a two-day period, an employee copied hundreds of technical documents from her employer, which were found in her luggage during a check at the airport.
- Spies working for U.S. defense companies stole internal information about the space shuttle, Delta IV rocket, the C-17 military plane, and submarine propulsion systems, along with other information. This information was provided to the Chinese government.

In a summary of the “Top Five Insider Attacks of the Decade,” the Linux.com editorial board listed the following cases⁵:

- Roger Duronio was convicted in 2006 to eight and a half years in federal prison⁶ for his actions in 2002 when, apparently in a fit of pique at not receiving what he considered an adequate annual bonus, he sabotaged the computer systems of his employer, UBS PaineWebber. “UBS was hit on March 4, 2002, at 9:30 in the morning, just as the stock market opened for the day. Files were deleted from up to 2,000 servers in both the central data center in Weehawken, N.J., and in branch offices around the country. Company representatives never reported the cost of lost business but did say it cost the company more than \$3.1 million to get the system back up and running. Duronio worked at UBS as a systems administrator until he quit a few weeks before the attack. Witnesses testified that he quit because he was angry that he didn’t receive as large an annual bonus as he expected. Investigators found copies of the malicious code on two of his home computers and on a printout sitting on his bedroom dresser.”⁷
- In 2005, a sting operation by a *The Sun* reporter from Britain netted him confidential details of more than a thousand “... accounts, passports, and credit cards ...” from NatWest and Barclays banks. This led to investigations of call centers in India, when criminals “... boasted of being able to provide details of as many as 200,000 bank accounts in a month, which, he further said, came from more than one call center.”⁸
- The “Athens Affair” was discovered when investigators looked into the apparent suicide of an electrical engineer, Costas Tsalikidis, in his Athens apartment. The inside job by Tsalikidis, the head of network planning, and unknown others compromised the Vodafone-Panafon company (“Vodafone Greece”) using malware and may have resulted in monitoring and recording of conversations involving “... the prime minister, his defense and foreign affairs ministers, top military and law-enforcement officials, the Greek EU commissioner, activists, and journalists.”⁹
- San Francisco network administrator Terry Childs locked other employees out of the city’s network in July 2008 because he claimed that his supervisor was unqualified to have administrative control. He was sentenced to four years in California state prisons. “Prosecutors characterized the former network administrator as a power hungry control freak who couldn’t be managed.”¹⁰
- Bradley Manning, a 22-year-old U.S. Army intelligence analyst, was arrested in May 2010 and charged in July 2010 with leaking nearly half a million classified U.S. videos and cables to the WikiLeaks project.¹¹ He was charged with “aiding the enemy” in February 2012¹² and accused of aiding the terrorist group al

THREATS FROM INSIDERS 13 · 5

Qaida through his actions.¹³ In January 2013, the trial was rescheduled until June 2013¹⁴ and Manning was denied the opportunity to justify his actions using a whistleblower defense.¹⁵

- *The Sunday Times* reported in 2012 that “Confidential personal data on hundreds of thousands of Britons is being touted by corrupt Indian call centre workers, an undercover investigation has discovered. Credit card information, medical and financial records are being offered for sale to criminals and marketing firms for as little as 2p.” Records of 500,000 Britons were apparently for sale, many of which were supposedly less than 72 hours old and included “... sensitive material about mortgages, loans, insurance, mobile phone contracts, and Sky Television subscriptions ...”¹⁶

13.2.3 Types of Insider Threats. There are three main classifications of insider threats: accidental, malicious, and nonmalicious.

13.2.3.1 Accidental Threats. Accidental threats are generally caused by mistakes; for example, staff may not follow operating procedures due to carelessness, disregard for policies, or a lack of training and awareness of the right thing to do. An example is a customer service representative who accidentally breaches client confidentiality by emailing client information to the wrong email address. Such errors may be caused by the use of email clients that have an auto-complete feature on the email address; staff under pressure to keep up with the volume of work may not notice the error before they send the data. This error is a particularly high risk for financial services organizations where in some jurisdictions a client confidentiality breach is a criminal offense.

Other typical examples include a database administrator who accidentally deletes a database table during maintenance, a systems administrator allowing a programmer to modify a production system without proper approvals, or an operator reformatting a disk drive without having two full, verified backups.

13.2.3.2 Malicious Threats. Malicious threats deliberately try to damage the organization or to benefit the attacker. Disgruntled IT administrators can sabotage IT systems, bringing an organization to a halt. There have been many incidents where both current and former administrators have deliberately caused system issues for various motives: enjoying the lifestyle of traveling around the world in luxury to fix the problems they created, extorting money from the organization, or simply causing as much damage as possible.

Some company information is highly valuable and specifically targeted by attackers. PII is one category that is sometimes illegally copied by staff to conduct identity theft and fraud, or to sell it to criminals. Another category is intellectual property (IP) such as trade secrets. Staff may take this information to help them with their next job or to sell to competing companies. This crime is thought to be common with IT developers who often seek to take their source code with them. Industrial espionage sponsored by rival companies or foreign governments is another common threat to IP (see Chapter 11 in this *Handbook* for examples).

Information can leave an organization by being copied to removable storage such as USB flash or hard drives and CD/DVD writers. Portable 3TB drives are now readily available, which means entire databases can be copied to a drive measuring less than 7 × 5 inches in size. With gigabit ethernet becoming standard, the time required to copy

13 · 6 THE INSIDER THREAT

this data is rapidly decreasing. Other common channels include emailing attachments to external email addresses, uploading files to external email services and to Internet Websites, and using cloud backup and cloud storage tools. To address these data-leakage channels, products known as data loss prevention (DLP) systems are increasingly being installed in organizations.¹⁷

There are also the common physical threats, such as taking printed information from people's desks or from the office printers. Where logical access controls are strong, staff intent on stealing information can take photographs of documents or information on screen with the high-resolution cameras in today's mobile phones.

There are some incidents where the motive may be conscience based as well as having a desire to damage an organization. Since the financial crisis in 2008, governments have increased their efforts to reduce tax avoidance, and are aggressively pursuing the use of foreign tax havens. There have been a number of publicized incidents where insiders have provided lists of offshore clients and accounts to country tax authorities.^{18,19}

13.2.3.3 Nonmalicious Threats. Nonmalicious threats are actions taken deliberately by people without intent to damage the organization. Often, the motive is to increase productivity, and the mistakes occur due to a lack of training or awareness of policies, procedures, and the risk. There have been many incidents in which staff loaded internal information onto Internet-based systems, some of which have no access controls. This error makes the information available to anyone who uses the sites and is often found and indexed by search engines.

One incident occurred in the early days of cloud computing. A drug researcher was given a lengthy lead time by his internal IT department for the delivery of infrastructure to conduct simulations. What he did instead was use his credit card to buy time on cloud-based systems and ran the simulations there instead. This would have put valuable intellectual property at risk had these Internet systems been compromised. This information was also sitting on the cloud provider's storage systems; what if the cloud vendor were to fail to reinitialize the storage when reallocating released storage to another customer, and the next user of the storage were to understand the value of what came pre-loaded on their system—and be dishonest?

Another example that is often reported in the media is the loss of PII when staff copy information to laptop computers or to removable storage devices such as USB drives or CDs/DVDs, which are then lost or stolen.²⁰

One common insider threat that can happen for both malicious and nonmalicious reasons is to email internal information to their home email address. Once on the staff member's own computer, the information is vulnerable to theft, successful attack on the computer or email account, or recycling of the machine by donating it to charity, or giving it away to family or friends. There have been many media reports of people finding sensitive information on secondhand computer hard disks.²¹

The nonmalicious motive for this practice is often to enable the employee to work from home, or the information is needed for a business trip. The malicious motive is to take the information with them, for reasons covered earlier.

13.2.4 Internet-Based Systems. The growing trend to use outsourced applications available over the Internet rather than internal systems can contribute to insider crime. With internal systems, when someone leaves an organization, they no longer have physical access to premises, and their network and application user-ids are supposed to be disabled immediately, removing access to corporate networks, computers,

MITIGATING THE INSIDER THREAT 13 · 7

and applications. Appropriate action makes it difficult for even maliciously motivated former staff to access confidential systems and their information.

But with Internet-based systems, the former staff member may still have access to confidential data via any Internet connection. It is very difficult to ensure all application accounts are promptly closed when someone leaves, creating a high risk of continued access to these systems. The people concerned may be deterred by the risk of being prosecuted for misusing their rights after leaving the organization, but their risk can be reduced if they use a former colleague's log-on account.

Organizations don't have to avoid using these external services; indeed, these external services are often among the best available, and can be provided at a much lower cost than hosting internally. However, security officers must coordinate closely with their human resources departments and the IT people responsible for maintaining the lists of external providers to ensure that access by former employees to external services is shut down as quickly as access to internal systems.

13.2.5 Service Provider Threats. Even organizations with strong personnel controls that reduce the risk from internal staff may have no protection from external service-provider staff. Although organizations commonly include their policies, standards, and procedures into contracts and treat this precaution as sufficient to address risk, the service provider may not consistently perform all the required controls, either through attempts to maximize profit or through poor management.

For example, a service provider faced with high staff turnover might bypass preemployment checks to get replacement staff quickly onto a service, particularly if there are service performance issues. Similarly, poor standards for handling termination of employment at the service provider could lead to compromise of client information.

In cases where an entire IT service is being provided by a service provider, the service provider's staff and IT administrators may have access to all the client organization's information. Email is a particular risk that may contain a great deal of an organization's intellectual property and that has powerful search facilities to easily target individuals and specific information.

13.2.6 System Administration Threats. System administrators have privileged access to an organization's IT infrastructure and, in poorly managed systems, may have access to critical and sensitive information on the systems. System administrators may deliberately attack the availability of an organization's systems even if they cannot access the data themselves. Administrators may destroy individual data sets, applications, or entire systems and networks—and may be able to destroy system backups to make the organization's recovery more difficult. This risk is most acute in smaller organizations, in which one person may manage the servers, applications, networks, email, backups, and perhaps even user administration. In larger organizations, it is much easier to segregate these roles as they should be, thus limiting access and therefore the damage one person can cause.

Even in large organizations with thousands of servers with effective segregation implemented, an administrator may be able to submit a job that can be run simultaneously on every server to wipe every hard disk. This type of threat highlights the importance of controlling the access and scope of administrators and preventing unauthorized changes from being implemented.

13.3 MITIGATING THE INSIDER THREAT. This section describes at a high level a few of some specific mitigating controls and how they address the insider threat.

13 · 8 THE INSIDER THREAT

For more details on the wider range of controls, refer to the detailed information in other chapters in this *Handbook* such as Chapter 15 (Penetrating Computer Systems and Networks), Chapter 28 (Identification and Authentication), Chapter 31 (Content Filtering and Web Monitoring), Chapter 45 (Employment Practices and Policies), Chapter 52 (Application Controls), Chapter 53 (Monitoring and Control Systems), Chapter 54 (Security Audits and Inspections), Chapter 55 (Cyber Investigation), Chapter 56 (Computer Incident Response Teams), and Chapter 68 (Outsourcing and Security).

13.3.1 System and Asset Inventories. If you don't know what you have, you can't manage it. Without an inventory of servers, you can't ensure they are patched, enabling an insider to compromise them and use them for their own purposes. Without a list of applications running on each server, you may have unnecessary servers sitting on your network being used for unauthorized purposes. The active systems on the network also need to correlate to the inventory. An administrator who doesn't remove a redundant system in your Internet DMZ could use it to bypass all of your network perimeter controls after he has left the organization.

In particular, one of the most dangerous tools for insider crime is the unauthorized and undetected wireless access point—easily purchased from any electronic store for at low cost and capable of transferring data from an internal network to unauthorized devices within the corporate facilities or even to external agents within a modest radius outside the building. See Chapter 33 in this *Handbook* for discussion of wireless network security.

13.3.2 Data Loss Prevention (DLP). There have been technology adoptions over the last 10 years which have made it much easier for staff to either accidentally or deliberately breach the confidentiality of organizational and client information. Data loss prevention (DLP) is a class of IT security system increasingly being implemented by organizations to help address these risks. A typical DLP system needs to address the following vulnerabilities at a minimum:

- Mobile storage devices/removable media, such as USB memory sticks and hard drives, mobile phones, memory cards, CD/DVD writers, along with infrared, Bluetooth, FireWire, and SCSI-connected storage.
- File uploads—including encrypted data—to external Websites via the standard protocols within Web browsers, such as ftp, http, and https. These controls may be enforced at both the Internet gateway as well as on the desktop.
- Detection of when laptops are not on the corporate network and preventing files being copied to noncorporate file shares.

For the majority of staff, the DLP system can be configured to block attempts to copy and upload data to these data-leakage channels. However, for most of these channels, there are going to be some people who have a genuine business need to copy and upload information, resulting in exceptions to the policy. The DLP system can help manage this risk by creating a log of what has been copied to support incident investigation or to enable a review of what a staff member has copied out of the organization. These facilities can be particularly useful when an employee hands in his resignation.

MITIGATING THE INSIDER THREAT 13 · 9

13.3.2.1 Email Data Leakage. Once the low-hanging fruit of removable storage and file uploads are blocked, staff begin to export information via email. The motive may not be malicious, but it still results in the organization's losing control over the information. To address this risk, the DLP system can be configured to report on people sending attachments to home email addresses, which is the usual destination for information, or to any unauthorized email address. These reports can then be used to increase staff awareness about policy or to support disciplinary processes for deliberate data leakage.

However, the true benefit from DLP comes when the business areas are engaged, because DLP can drive a process of identifying and defining the critical information that has to be protected. For example, the IT development group could have a DLP rule that blocks emails containing source code, to prevent developers taking their work with them when they leave. Even if source code is placed in an encrypted zip file to try and avoid detection, the metadata (e.g., filenames and identifiers of file creators) may still be readable and can trigger the rule.

For other areas of the business, client information, business strategy, business results, intellectual property, and PII such as credit card numbers and Social Security numbers can be configured into the DLP system.

13.3.2.2 Data Leakage Using Cloud Storage. Widespread availability of external data storage facilities (e.g., Dropbox and Google Drive) adds to the complexity of DLP. Careful application of Web monitoring and blacklisting specific URLs may be helpful, but determined opponents of the regulations may circumvent such methods using a variety of proxy avoidance Websites which mask the destination of the HTTP request. Security administrators should be on the lookout for new sites so they can add them to the corporate blacklists for outbound communications through their firewalls.

13.3.2.3 Difficulties with DLP. DLP is not a panacea. It is difficult if not impossible to prevent someone determined to leak data, but with a DLP system, one can make it difficult for them to do so without detection, and this barrier usually deters the majority of people.

With email DLP, it is very difficult to identify safe blocking rules that prevent data leakage. For example, if employees routinely send emails to customers, some of these customers are going to be using their home email addresses, which means that a hard block (blacklisting) of emails sent to home email addresses won't be feasible.

13.3.2.4 Legal Issues with DLP. With the implementation of DLP, an organization progresses from investigating reported incidents to actively monitoring for breaches of company policy. One of the major issues with implementing a global DLP system is that active monitoring may be subject to privacy and workplace laws, and failure to comply with these laws in some countries is a criminal offence.

Additionally, with customer and organization information being captured in DLP logs, consideration needs to be given as to where the log files are stored and who will be reviewing them. As an example from the financial services industry, some information is price sensitive and there are strict requirements on who can access it, to prevent insider dealing. Client information within the log files may be covered under banking secrecy laws and regulations. If the log files for one country are stored on a server in another country, then outsourcing regulations need to be complied with as well, and planners must identify the most stringent regulations to ensure efforts for

13 · 10 THE INSIDER THREAT

compliance. Data protection laws have been enacted in many countries that require data subjects to give their agreement to their information being processed and for defined purposes; therefore, monitoring of communications may need to be included in customer contracts and other data protection declarations.

Despite all of these hurdles, for most countries, it is possible to roll out a DLP system. As part of the project planning stage, the organization needs to commission a legal investigation for each country to understand whether it is lawful to actively monitor corporate email, but more importantly, the preconditions for any monitoring to lawfully take place. For countries that forbid email monitoring or any form of workplace surveillance, it is usually possible to implement blocking controls on writing to removable storage along with uploads to external Websites, provided no log files are kept.

13.3.2.5 Remote Access Solution. A commonly implemented remote access solution is the provision of Webmail from home computers. With Webmail you need to configure the system to prevent staff opening attachments within local applications on the home computer. If this is not blocked, staff can copy corporate information by saving the open attachment to their local hard drive. In addition, any use of the Remote Desktop Protocol (RDP) also needs to be properly configured to prevent local USB resources being mapped to the remote computer.

13.3.3 Internal Honeypots. Honeypots are attractive systems or nodes that appear to have valuable confidential data. Roger Grimes, author of a textbook about honeypots,²² writes:

One of the best things you can do to get early warning of internal attackers is to implement honeypots. ... because they're low cost and low noise—and they work!

Because the internal attackers you seek could be trusted IT employees, the entire project must be kept secret. It should be known only to the sponsor, management, and the implementers. Often, we give the project a boring code name like Marketing Business Development, which is used in all documents and e-mails, avoiding terms having anything to do with honeypots. Don't even tell the network security people about it, in so far as you can and still have an operational project. Then take a few computers that are destined for de-provisioning or the scrap heap and turn them into your honeypots.²³

The point of such a honeypot is that there should be zero access to it: It serves no function and there is never a reference to it in any internal or external documentation. Thus, anyone who accesses it is either doing so by pure accident or is violating policies governing unauthorized access to internal data (there is no one authorized to access the honeypot). Detailed logging should be in place at all times to provide forensic information immediately; however, administrators should ensure that they can actually visualize exactly what is being done in real time, not simply rely on the detailed log files for after-the-fact analysis. An early-warning system should immediately alert system administrators to the problem (preferably as the access is in progress) via screen messages, voice messages, and text messages.

13.4 CONCLUDING REMARKS. Fighting the insider threat need not be solely reactive. In addition to the entire spectrum of information assurance measures covered throughout this *Handbook*, maintaining an environment of security awareness and of encouragement, trust, fair-dealing, and long-term commitment to the welfare of employees must remain among the very best approaches to reducing the risk of insider crime.

NOTES 13 · 11**13.5 FURTHER READING**

- Department of Defense. *DoD Insider Threat Mitigation: Final Report of the Insider Threat Integrated Process Team*. DoD. 2000. www.dtic.mil/cgi-bin/GetTRDoc?AD=ADA39138
- Noonan, T. & E. Archuleta. *The National Infrastructure Advisory Council's Final Report and Recommendations on the Insider Threat to Critical Infrastructures*. NIAC. 2008. www.dhs.gov/xlibrary/assets/niac/niac_insider_threat_to_critical_infrastructures_study.pdf
- Silowash, G., D. Gappelli, A. Moore, R. Trzeciak, T. J. Shimeall, L. Flynn. *Common Sense Guide to Mitigating Insider Threats*, 4th ed. CMU/SEI. 2012. Technical Report CMU/SEI-20120TR-012. www.sei.cmu.edu/reports/12tr012.pdf

13.6 NOTES

1. CMU/SEI 2011
2. PwC 2012
3. PwC 2013
4. See www.fbi.gov/about-us/investigate/counterintelligence/the-insider-threat
5. Linux.com editorial staff, "Top Five Insider Attacks of the Decade," *LINUX.COM*, January 13, 2011, www.linux.com/news/technology-feature/security/397143-top-five-insider-attacks-of-the-decade
6. Sharon Gaudin, "Ex-UBS Systems Admin Sentenced To 97 Months In Jail," *InformationWeek*, December 13, 2006, www.informationweek.com/ex-ubs-systems-admin-sentenced-to-97-mon/196603888
7. Sharon Gaudin, "Ex-UBS Sys Admin Found Guilty, Prosecutors To Seek Maximum Sentence," *InformationWeek*, July 19, 2006, www.informationweek.com/ex-ubs-sys-admin-found-guilty-prosecutor/190700064
8. SiliconIndia News, "Indian call centers selling U.K.'s secrets," *SiliconIndia News*, June 23, 2005, www.siliconindia.com/shownews/Indian_call_centers-selling_UKs_secrets-nid-28560-cid-2.html
9. Vassilis Prevelakis and Diomidis Spinellis, "The Athens Affair: How some extremely smart hackers pulled off the most audacious cell-network break-in ever," *IEEE SPECTRUM*, June 29, 2007, <http://spectrum.ieee.org/telecom/security/the-athens-affair>
10. Robert McMillan, "Network admin Terry Childs gets 4-year sentence," *NetworkWorld*, August 17, 2012, www.networkworld.com/news/2010/080710-network-admin-terry-childs-gets.html
11. Chris McGreal, "US private Bradley Manning charged with leaking Iraq killings video," *The Guardian*, July 6, 2010, www.guardian.co.uk/world/2010/jul/06;bradley-manning-charged-iraq-killings-video
12. Karen McVeigh, "Bradley Manning defers plea after being formally charged with aiding the enemy: No date set for WikiLeaks suspect's trial but Manning's lawyer says he would object to any delay in the trial beyond June," *The Guardian*, February 23, 2012, www.guardian.co.uk/world/2012/feb/23;bradley-manning-defer-plea-charges
13. Associated Press, "Bradley Manning aided al-Qaida with WikiLeaks documents, military says: Manning, charged with aiding the enemy, accused of indirectly

13 · 12 THE INSIDER THREAT

- aiding terrorist group by leaking thousands of documents,” *The Guardian*, March 3 2012, www.guardian.co.uk/world/2012/mar/15/bradley-manning-wikileaks
14. Adam Gabbatt and Ed Pilkington, “Bradley Manning trial delayed until June after sentence reduction granted: Judge reschedules US soldier’s trial to give more time for review of classified information related to WikiLeaks case,” *The Guardian*, January 9, 2013, www.guardian.co.uk/world/2013/jan/09/bradley-manning-trial-delayed
 15. Ed Pilkington, “Bradley Manning denied chance to make whistleblower defence: Judge rules that Manning will not be allowed to present evidence about his motives for the leak—a key plank of his defence,” *The Guardian*, January 17, 2013. www.guardian.co.uk/world/2013/jan/17/bradley-manning-denied-chance-whistleblower-defence
 16. Tom Gardner, “Indian call centres selling YOUR credit card details and medical records for just 2p,” *MailOnline*, March 18, 2012, www.dailymail.co.uk/news/article-2116649/Indian-centres-selling-YOUR-credit-card-details-medical-records-just-2p.html
 17. Eric Ouellet, “Magic Quadrant for Content-Aware Data Loss Prevention,” Gartner, Inc., January 3, 2013, <https://www.ca.com/us/register/forms/collateral/Magic-Quadrant-for-Content-Aware-Data-Loss-Prevention-2013.aspx>
 18. BBC 2012
 19. BBC 2011
 20. BBC 2008
 21. BBC 2009
 22. Roger A. Grimes, *Honeypots for Windows*, Apress, 2005.
 23. R. A. Grimes, “Honeypots: A sweet solution to the insider threat: A honeypot can be a cheap, easy, and effective warning system against the trusted insider gone bad,” *InfoWorld*, May 1, 2009, www.infoworld.com/d/security-central/honeypots-sweet-solution-insider-threat-922?page=0,0

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 14

INFORMATION WARFARE

Seymour Bosworth

14.1 INTRODUCTION	14·2	14.5 WEAPONS OF CYBERWAR	14·21
14.2 VULNERABILITIES	14·2	14.5.1 Denial of Service and Distributed Denial of Service	14·21
14.2.1 Critical Infrastructure	14·2	14.5.2 Malicious Code	14·22
14.2.2 Off-the-Shelf Software	14·3	14.5.3 Cryptography	14·23
14.2.3 Dissenting Views	14·4	14.5.4 Psychological Operations	14·23
14.2.4 Rebuttal	14·4	14.5.5 Physical Attacks	14·24
14.3 GOALS AND OBJECTIVES	14·5	14.5.6 Biological and Chemical Weapons and Weapons of Mass Destruction	14·25
14.3.1 Infrastructure	14·5	14.5.7 Weapons Inadvertently Provided	14·25
14.3.2 Military	14·5		
14.3.3 Military Offensives	14·8		
14.3.4 Government	14·8		
14.3.5 Energy Systems	14·9		
14.3.6 Transportation	14·11		
14.3.7 Commerce	14·12	14.6 DEFENSES	14·25
14.3.8 Financial Disruptions	14·13	14.6.1 Legal Defenses	14·25
14.3.9 Medical Security	14·14	14.6.2 Forceful Defenses	14·26
14.3.10 Law Enforcement	14·15	14.6.3 Technical Defenses	14·27
14.3.11 International and Corporate Espionage	14·15	14.6.4 In-Kind Counterattacks	14·27
14.3.12 Communications	14·16	14.6.5 Integration of Cyberwarfare into Military Planning	14·27
14.3.13 Destabilization of Economic Infrastructure	14·16	14.6.6 Cooperative Efforts	14·28
14.4 SOURCES OF THREATS AND ATTACKS	14·17	14.7 SUMMARY	14·29
14.4.1 Nation-States	14·17	14.8 FURTHER READING	14·29
14.4.2 Cyberterrorists	14·19		
14.4.3 Corporations	14·21	14.9 NOTES	14·30

14 · 2 INFORMATION WARFARE

Information warfare is the offensive and defensive use of information and information systems to deny, exploit, corrupt, or destroy, an adversary's information, information-based processes, information systems, and computer-based networks while protecting one's own. Such actions are designed to achieve advantages over military or business adversaries.

—Dr. Ivan Goldberg, Institute for Advanced Study of Information Warfare

14.1 INTRODUCTION. Until recently, warfare was conducted by armed forces representing adversarial nations, or by revolutionary elements opposing their own governments. Today, although such conflicts still exist around the world, the ubiquitous nature of computers and associated technology has created new forces, new threats, new targets, and an accompanying need for new offensive and defensive weapons. Information warfare (IW), also known as *e-warfare* or *cyberwar*, is actually, or potentially, waged by all U.S. armed forces and by those of other nations, as well as by commercial enterprises, by activist groups, and even by individuals acting alone.

Conventional wars, whether large or small, are regularly reported by the news media. Information war, however, are largely ignored except by those with a professional interest in the field. One reason for this is that conventional warfare is a matter of life or death; photos and eyewitness accounts are dramatic reminders of human cruelty and mortality. In contrast, IW has so far been conducted bloodlessly, with only economic and political consequences. However, it is becoming increasingly evident that IW may soon be conducted in ways that could equal or exceed the death and destruction associated with conventional weapons.

Conventional wars are fought by known combatants with clearly defined allies and antagonists, but IW often is waged by unknown entities with uncertain allegiances and goals. IW may be conducted on many fronts simultaneously, with wars fought within wars, and with both civilian and military targets devastated.

The motives for conventional warfare were almost always territorial, religious, political, or economic. These are still important, but to them must be added the psychological motivations of groups and individuals—groups far more widely distributed and less easily overcome.

This chapter discusses information warfare in terms of the vulnerabilities of targets, participants' objectives, sources of threats and attacks, weapons used, and defenses against those weapons.

14.2 VULNERABILITIES. Until recently, concerns over the security of the technological infrastructure in technologically advanced nations have been viewed with skepticism. However, by the mid-1990s, opinion leaders in government, industry, and the security field were coming to grips with widespread vulnerabilities in the critical infrastructure.

14.2.1 Critical Infrastructure. In 1998, President Bill Clinton circulated Presidential Decision Directive 63, which outlined his administration's policy on critical infrastructure protection:

Critical infrastructures are those physical and cyber-based systems essential to the minimum operations of the economy and the government. ... They include, but are not limited to, telecommunications, energy, banking and finance, transportation, water systems and emergency services, both government and private.¹

VULNERABILITIES 14 · 3

Having defined the very broad, vital areas that require protection, the paper went on to describe succinctly their vulnerability:

The United States possesses both the world's strongest military and its largest national economy. Those two aspects of our power are mutually reinforcing and dependent. They are also increasingly reliant upon certain critical infrastructures and upon cyber-based information systems....

Because of our military strength, future enemies, whether nations, groups or individuals, may seek to harm us in non-traditional ways including attacks within the United States. Our economy is increasingly reliant upon interdependent and cyber-supported infrastructures and non-traditional attacks on our infrastructure and information systems may be capable of significantly harming both our military power and our economy.

A few examples of specific weaknesses were given by Jack L. Brock, Jr., director, Government-wide and Defense Information Systems, United States General Accounting Office:

In May 1999 we reported that, as part of our tests of the National Aeronautics and Space Administration's (NASA) computer-based controls, we successfully penetrated several mission-critical systems. Having obtained access, we could have disrupted NASA's ongoing command and control operations and stolen, modified, or destroyed systems software and data.

In August 1999, we reported that serious weaknesses in Department of Defense (DOD) information security continue to provide both hackers and hundreds of thousands of authorized users the opportunity to modify, steal, inappropriately disclose, and destroy sensitive DOD data.²

Although these "attacks" were carried out one at a time, and without malicious intent, it is apparent that they, and many others, could have been launched simultaneously and with intent to inflict the maximum possible damage to the most sensitive elements of the national infrastructure.

In a memorandum to its chairman, describing a report of the Defense Science Board Task Force on Defensive Information Operations, Larry Wright stated in 1999 that:

The threats to the DoD infrastructure are very real, non-traditional and highly diversified.... The vulnerabilities of these United States are greater than ever before, and we know that over twenty countries already have or are developing computer attack capabilities. Moreover, the Department of defense should consider existing viruses and "hacker" attacks to be real "Information Operations or Warfare," what early aviation was to Air Power. In other words, we have not seen anything yet!³

The report concluded that "[i]t is the view of this task force that DoD cannot today defend itself from an Information Operations attack by a sophisticated nation state adversary."

14.2.2 Off-the-Shelf Software. One characteristic of almost all military and civilian infrastructures is that they share, with more than 100 million computers, a single ubiquitous operating system, and many of the same applications programs, such as word processors, spreadsheets, and database software. These commercial off-the-shelf (COTS) products are available around the world, to friend and foe alike, and they appear to be more intensively studied by malefactors than by their security-inadequate producers. Each of these products presents entry points at which one common vulnerability may be exploited to damage or destroy huge portions of the national infrastructure.

14 · 4 INFORMATION WARFARE

Until, and unless, this software is rendered significantly more resistant to attack, all of its users remain at risk.

14.2.3 Dissenting Views. Not every influential observer concurs in these possible scenarios. Dr. Thomas P. M. Barnett, a professor and senior decision researcher at the Decision Support Department, Center for Naval Warfare Studies, U.S. Naval War College, voiced a fairly typical disagreement in 1999:

If absence makes the heart grow fonder, network-centric warfare is in for a lot of heartbreak, because I doubt we will ever encounter an enemy to match its grand assumptions regarding a revolution in military affairs. The United States currently spends more on its information technology than all but a couple of great powers spend on their entire militaries. In a world where rogue nations typically spend around \$5 billion a year on defense, NCW is a path down which only the U.S. military can tread.⁴

14.2.4 Rebuttal. It may be of some benefit to have spokespersons for this unworried viewpoint, but their opinions must be weighed against those, for example, of Scott Henderson, of the Navy-Marine Corps intranet, who said: “One of our critical capabilities will be how we are to defend our information and our information systems from an adversary’s attack.”⁵ He stated that successful intrusions, or attacks, on Navy computer systems increased from 89 in 2000 to 125 by mid-2001, an annualized increase of 80 percent. Those figures did not include successful attacks that went undetected or unsuccessful attempts that may have identified a weak point from which to launch future and probably more successful, attacks.

A highly significant factor in IW is its asymmetric nature. The barriers to entry for attackers are low; their weapons can be inexpensive, easily obtained, highly effective, easily automated, and used with negligible risk of personal harm. In contrast, defensive measures are extremely costly in time, money, and personnel and they may be ineffective against even unsophisticated attackers using obsolete computers.

The 2013 *Annual Report to Congress: Military and Security Developments Involving the People’s Republic of China* from the U.S. Office of the Secretary of Defense includes the following evaluation of China’s “Anti-Access/Area Denial (A2/AD)” capabilities:

An essential element, if not a fundamental prerequisite, of China’s emerging A2/AD regime is the ability to control and dominate the information spectrum in all dimensions of the modern battlespace. PLA authors often cite the need in modern warfare to control information, sometimes termed “information blockade” or “information dominance,” and to seize the initiative and gain an information advantage in the early phases of a campaign to achieve air and sea superiority. China is improving information and operational security to protect its own information structures, and is also developing electronic and information warfare capabilities, including denial and deception, to defeat those of its adversaries. China’s “information blockade” likely envisions employment of military and non-military instruments of state power across the battlespace, including in cyberspace and outer space. China’s investments in advanced electronic warfare systems, counter-space weapons, and computer network operations (CNO)—combined with more traditional forms of control historically associated with the PLA and CCP systems, such as propaganda and denial through opacity, reflect the emphasis and priority China’s leaders place on building capability for information advantage.⁶

Considering the nature and extent of already successful attacks against major elements of U.S. military and civilian infrastructures, there appears to be no justification for discounting the views of those who believe that IW, in both its offensive and defensive roles, must be accorded the attention that surrounds any potentially cataclysmic force. This *Handbook*, especially Chapters 16, 17, 18, 20, and 21, contains many

GOALS AND OBJECTIVES 14 · 5

examples of viruses, worms, and other malware that have created massive disruptions in very large networks. The worst-case scenarios presented here should serve to awaken a measured response in those who may have been unaware or unconcerned.

14.3 GOALS AND OBJECTIVES. Attacking forces, in information warfare, will always have a variety of strategic and tactical motives behind their actions; defensive forces generally have only one tactical goal—to blunt the enemy’s attack and, if possible, to counterattack. Only after this is accomplished, and the nature of the attackers has been studied, can strategies for long-range operations be determined and effected.

14.3.1 Infrastructure. Depending on the target, an attacker’s goals may vary widely, but attackers generally want to damage, subvert, or destroy the infrastructure. In doing so, an attacker would hope to bring government, the economy, and military operations to a standstill or at least to reduce their efficiency and effectiveness to instill fear, uncertainty, and doubt (FUD), and ultimately to induce widespread chaos that could cost many lives.

Although this view is entirely appropriate to wars between nations or to campaigns by terrorists, it must be tempered when considering commercial warfare, whose main goal is competitive financial advantage.

14.3.2 Military. Today, information warfare is a vital concern of area commanders under battlefield conditions. They must obtain complete, accurate, and timely information about their opponents’ actions, intentions, weaknesses, and resources while denying the same to their adversaries. The ultimate objective for all of these activities is to support the military tactics that will maximize the enemy’s body count, or at least to render its defenses ineffective, so that surrender becomes the only viable option. The other side of the coin, defensive tactics, are aimed at preventing enemies from accomplishing their objectives.

In the United States, the Joint Chiefs of Staff (for Army, Navy, Marine Corps, Coast Guard, and Air Force) have formulated the *Joint Doctrine for Operations Security* to be followed by all commanders of combatant commands in planning, preparation, and execution of joint operations. The publication states:

Operations Security (OPSEC) is a process of identifying critical information and subsequently analyzing friendly actions attendant to military operations and other activities, to: (a) identify those operations that can be observed by adversary intelligence systems; (b) determine what indicators adversary intelligence systems might obtain that could be interpreted or pieced together to derive critical information in time to be useful to adversaries; and (c) select and execute measures that eliminate or reduce to an acceptable level the vulnerabilities of friendly actions to adversary exploitation.⁷

OPSEC is a process that could be applied to every element of civilian infrastructure, as well as to the military, although all sources of information commonly used by the military are not available to the civilian sector. Other military code words for intelligence activities are:

- HUMINT (human intelligence) is the most widely used source of information, as it has always been for both the civilian and military sectors. HUMINT is often the only source capable of direct access to an opponent’s plans and intentions. Some

14 · 6 INFORMATION WARFARE

intelligence gathering is quite open, but covert or clandestine operations must be conducted in secrecy, so as to protect the sources of confidential information.

- SIGINT (signals intelligence) is obtained from communications (COMINT), electronics (ELINT), and foreign instrumentation signals (FISINT).
- COMINT (communications intelligence) is information intended for others and intercepted without leaving a trace.
- ELINT (electronic intelligence) derives technical or geographic location data from an opponent's electromagnetic radiations, other than those that arise from communications or from nuclear detonations or radioactive sources. The primary ELINT sources are radars (radio detection and ranging).
- FISINT (foreign instrumentation signals intelligence) is obtained from intercepting and analyzing metered performance parameters electronically transmitted from sources such as a ballistic missile.
- MASINT (measurement and signatures intelligence) is scientific and technical in nature. Its purpose is to identify distinctive features associated with a source, emitter, or sender so as to facilitate subsequent identification or measurement. These features include wavelength, modulation, time dependencies, and other unique characteristics derived from technical sensors.
- IMINT (imagery intelligence) is produced by photography, infrared sensors, lasers, radars, and electro-optical equipment. This equipment, operated from land, sea, air, or space platforms, provides strategic, tactical, and operational information.
- TECHINT (technical intelligence) is derived from the exploitation and analysis of captured or otherwise acquired foreign equipment.
- OSINT (open source intelligence) is available to the general public from news media, unclassified government publications, public hearings, contracts, journals, seminars, and conferences. The World Wide Web has become an important tool of OSINT.

The *Joint Doctrine for Operations Security* lists several generic military activities with some of their associated critical information. It must be the objective of all information warfare to acquire this critical information about their opponents while denying such information to them:

- *Diplomatic negotiations* include military capabilities, intelligence verification, and minimum negotiating positions.
- *Political-military crisis management* includes target selection, timing considerations, and logistic capabilities and limitations.
- *Military intervention* requires information about intentions, military capabilities, forces assigned and in reserve, targets, and logistic capabilities and constraints.
- *Counterterrorism* involves forces, targets, timing, strategic locations, tactics, and ingress and egress methods.
- *Open hostilities* information involves force composition and disposition, attrition and reinforcement, targets, timing, logistic constraints, and location of command and control (C²) nodes.

GOALS AND OBJECTIVES 14 · 7

- *Mobilization* requires information about an intent to mobilize before public announcement, impact on military industrial base, impact on civilian economy, and transportation capabilities and limitations.
- *Intelligence, reconnaissance, and surveillance* information includes purpose and targets of collection, timing, capabilities of collection assets, and processing capabilities.

In addition to the Joint Chiefs' doctrines, the Department of Defense and each individual branch of service have been charged with the responsibility for establishing task forces, advisory groups, training and awareness programs, and virtual information networks to mobilize IW forces and to bring into being a strong defense against enemy attack.

Further evidence of the importance of military information and the vulnerabilities that exist at this time is contained in the 2001 report of the Secretary of Defense to the President and the Congress:

Information superiority is all about getting the right information to the right people at the right time in the right format while denying adversaries the same advantages. The United States enjoys a competitive advantage in many of the technical components of information superiority, but the U.S. also has vulnerabilities stemming from its increasing dependence on high technology. Experiences from Somalia to the Balkans have shown that low technology adversaries also can wage effective information campaigns, especially in urban environments.

In the Information Age, the opportunities and obstacles to achieving national security objectives often are informational in nature. Information superiority is a principal component of the transformation of the Department. The results of research, analyses, and experiments, reinforced by experiences in Kosovo, demonstrate that the availability of information and the ability to share it significantly enhances mission effectiveness and improves efficiencies. Benefits include: increased speed of command, a higher tempo of operations, greater lethality, less fratricide and collateral damage, increased survivability, streamlined combat support, and more effective force synchronization. Kosovo also highlighted the shortage of assets for intelligence, surveillance, and reconnaissance, as well as the need for more secure interoperability and information protection, especially within coalitions.

To ensure that the above prerequisites are in place, DoD is developing appropriate policy and oversight initiatives, actively pursuing opportunities to improve international cooperation in the areas of Command, Control, Communication, Computers, Intelligence, Surveillance, and Reconnaissance (C4ISR) and space-related activities, partnering with industry, and working to anticipate and understand the implications of emerging information technologies.

The quality of DoD's infostructure will be a pacing item on the journey to the future. The ability to conceive of, experiment with, and implement new ways of doing business to harness the power of Information Age concepts and technologies depends upon what information can be collected, how it can be processed, and the extent to which it can be distributed. The ability to bring this capability to war will depend upon how well it can be secured and its reliability. DoD envisions an infostructure that is seamless with security built-in, one that can support the need for increased combined, joint, and coalition interoperability, leverages commercial technology, and accommodates evolution.⁸⁷

Although not as well publicized as are the U.S. defensive efforts, equal attention, time, and resources are being expended on actual and possible offensive operations. Every objective, every tactic, and every recommendation just mentioned, and some too sensitive to discuss here, are subjects for study and implementation of offensive strategies and tactics aimed at enemies, present and future.

14 · 8 INFORMATION WARFARE

14.3.3 Military Offensives. The nature of U.S. offensive planning has been described in detail in a top-secret presidential memo published by *The Guardian* (U.K.) on June 7, 2013.

The memo states that:

The United States has an abiding interest in developing and maintaining use of cyberspace as an integral part of U.S. national capabilities to collect intelligence, and to deter, deny, or defeat any adversary that seeks to harm U.S. national interests in peace, crisis, or war. ... The United States Government shall conduct DCEO [Defensive Cyber Effects Operations] and OCEO [Offensive Cyber Effects Operations] under this directive consistent with its obligations under international law, including with regard to matters of sovereignty and neutrality, and as applicable, the laws of armed conflict. This directive pertains to cyber operations, including those that support or enable kinetic, information, or other types of operation.⁹

In this context, “kinetic operations” are a recognized euphemism for warfare, and this memo provides the U.S. Government with a broad mandate to protect U.S. interests, as it sees them, with any means available to it without a declaration of war.

At the time of this writing (July 2013), the full implications of this serious security breach cannot be evaluated. At the very least, it will provide the critics of U.S. policy with reinforcement for their view of the United States as a militaristic, even terrorist force in international relations. For defenders of U.S. policies, it will represent an ordered, rational, response to threats that must be contemplated and guarded against.

14.3.4 Government. The objectives of government, at every level, must be to protect the lives and welfare of its constituencies. Any breakdown in an essential government function may produce marked unrest, rioting, vandalism, civil disobedience, and possibly much bloodshed.

Just as in the military, government must be able to defend itself against an information attack waged by any enemy of the established order. Although not every element of government is perceived by all to perform a useful function, there are agencies without which it would be virtually impossible to sustain a developed nation’s day-to-day activities.

At the federal level, civil servants’ salaries, Social Security payments, tax collections and disbursements, military expenditures, lawmaking, and a myriad of other functions and activities can be carried out only with the active and pervasive use of computers and computer networks. In the past, some of these computer operations have been penetrated by hackers, crackers, and political dissidents, but only one at a time. It does not require a science fiction writer to imagine what the effect would be if simultaneous attacks were successfully launched against major federal government agencies.

At state levels, although the effects would be more constrained geographically, a great deal of damage could be done to emergency response units, to police and judiciary functions, and to health and welfare services. All of these depend on computerized functions that are protected even less than those of federal agencies.

For municipalities and even smaller governments, zoning enforcements and other local functions can be suspended without serious consequences, but police radio and computer networks are easily penetrated, and their ability to maintain law and order compromised.

GOALS AND OBJECTIVES 14 · 9

As demonstrated by many previous incidents, government functions at any level are susceptible to information warfare. Natural events, Murphy's law (what *can* go wrong *will* go wrong), poorly configured systems, flawed operating systems and application programs, together with inadequate security measures underlie the vulnerability of government systems.

14.3.5 Energy Systems.¹⁰ Oil, coal, and hydroelectric systems are critical elements of the national infrastructure. Supervisory control and data acquisition (SCADA) security in the electric power industry suffers from widespread misconceptions and a breakdown in communications between administrators and security experts. In brief,

- Attacks on electric power plants and the distribution grid may not result in the catastrophic scenarios painted by the promoters of panic, but any interruption in electric power delivery can cause widespread infrastructure disruption.
- SCADA systems controlling electric generators and distribution systems are not, in fact, isolated by air gaps from the Internet.
- On the contrary, vulnerability analysis teams have systematically and repeatedly demonstrated that power companies are unaware of the reality of their interconnectedness and vulnerabilities.
- There are documented cases of industrial espionage, sabotage, denial of service, and malware attacks on electric power grid SCADA systems.
- SCADA systems have been considered too stable to bother updating with current patches; as a result, they are consistently vulnerable to exploits of current (and even ancient) vulnerabilities.
- Many SCADA systems were developed without consideration of security, secure coding, or integration of security dimensions of software quality assurance.
- Government and academia have significant projects in place to advance SCADA security, but acceptance by industry is modest at best. Academics engaged in SCADA security research are doing a good job of reaching other academics through peer-reviewed presentations at academic conferences; they are less successful in reaching managers at power companies.
- Pressure is rising in the public sphere, in government circles, among security practitioners, and within the electric power industry to come to grips with the need for improved cybersecurity.
- The electric power industry must coordinate its efforts to implement well-established standards for protecting computer systems and networks in all its SCADA systems and related networks. In addition, the industry should implement cyber situational awareness solutions to integrate multiple inputs from SCADA and network sensors that will permit intelligent, agile response to attacks and effective forensic analysis of those attacks.

The electric power industry has become a fundamental underpinning of twenty-first century life. In a landmark report on "The Electricity Economy," author Jesse Berst and colleagues describe the convergence of growing demand, an increasing dependence on computerized SCADA systems, and the inevitable complexity of interactions among elements controlled by diverse entities with limited coordination.¹¹ To illustrate the

14 · 10 INFORMATION WARFARE

growth in electricity demands, the report's Table 1 shows global electricity demands of 2.06 terawatts (TW) in 1950 versus 3.8 TW in 2000 and a predicted 6.99 TW in 2050. The proportion of electricity as a percentage of global energy utilization was 10.4 percent in 1950 and 25.3 percent in 2000; by 2050 it may reach 33.7 percent. The authors add,

Today we depend on electricity for basic needs such as food, water, shelter, communication, employment, and health care. Those needs are served by infrastructures for food preservation, water treatment, heat and light, phone service, Internet, offices, factories, hospitals and emergency response, to name a few. Yet all of those essentials degrade or disappear without electricity.¹²

In October 1997, the President's Commission on Critical Infrastructure Protection (the "Marsh Report" named after Commission Chairman Robert T. Marsh) included the following warning:

Prolonged disruption in the flow of energy would seriously affect every infrastructure.

The significant physical vulnerabilities for electric power are related to substations, generation facilities, and transmission lines. Large oil refineries are also attractive targets. The increase in transportation of oil via pipelines over the last decade provides a huge, attractive, and largely unprotected target array. Oil and gas vulnerabilities include lines at river crossings; interconnects; valves, pumps, and compressors; and natural gas city gates. Large metropolitan areas could be deprived of critical fuel for an extended period by a properly executed attack.

The widespread and increasing use of Supervisory Control and Data Acquisition (SCADA) systems for control of energy systems provides increasing ability to cause serious damage and disruption by cyber means. The exponential growth of information system networks that interconnect the business, administrative, and operational systems contributes to system vulnerability.¹³

Electrical power systems have been harmed through data leakage, industrial espionage, insider threats and sabotage.¹⁴ Incidents described in the reference include

- 2006 Japan's Power Plant Security Info Leaked Onto Internet
- 2007 Egypt Accuses Nuclear Employee of Spying
- 2007 Former Nuclear Plant Engineer Allegedly Took Data to Iran
- 2007 Saboteur of California Power Grid Gained Access Despite Warning
- 2009 Fired Nuclear-Power-Plant Employee Arrested for Hacking Systems
- 2009 (Former) IT Consultant Confesses to SCADA Tampering

Hackers and malware writers and distributors have also attacked power systems.¹⁵ Cases summarized in the reference include

- 2000 Hacker Shocks Electric Company
- 2003 Slammer Worm Crashes Ohio Nuclear Plant Network
- 2006 National Nuclear Security Administration Computers Hacked; Info on 1500 Taken
- 2010 Stuxnet Worm Attacks SCADA Vulnerabilities

GOALS AND OBJECTIVES 14 · 11

The Stuxnet incident is a significant development in information warfare. The European Network and Information Security Agency (ENISA) published a detailed analysis of the case in 2010:

Stuxnet is a specialised malware targeting SCADA systems running Siemens SIMATIC® WinCC or SIMATIC® Siemens STEP 7 software for process visualisation and system control. SCADA in general refers to computer systems that monitor and control industrial processes, such as, e.g., those in nuclear power plants, or in facilities for water treatment.

This highly sophisticated malware uses several vulnerabilities in the underlying Windows® operating system for infection and propagation. Infection works via USB-drives or open network shares. A root kit component hides the content of the malware on infected WinCC systems. An infected system can usually be controlled remotely by the attacker. In the end this means that the attacker has full control of the respective facility.¹⁶

ENISA also published a report on securing computer-controlled energy-distribution systems (*smart grids*) in December 2012.¹⁷

14.3.6 Transportation. Airplanes, trains, trucks, and ships are all likely targets for physical and information warfare. Because all of them are necessary to support the infrastructure by transporting personnel and materials, any disruption can cause severe problems. Because all of these transportation systems increasingly rely on sophisticated telecommunications and computing resources, they are subject to information warfare.

14.3.6.1 Aviation. The most visible, and potentially the most vulnerable, component of the transportation infrastructure is the aviation industry. Unlike the fly-by-the-seat-of-your-pants technology of aviation's early days, today's airplanes and the systems that dispatch and control them in flight are almost totally dependent on electronic communications and instruments, both analog and digital.

To a great extent, almost every airplane depends on its global positioning system (GPS) to determine its position in space, its course, speed, bearing to an airfield, and other important functions. Airplanes generally are required to fly at certain altitudes, in specific corridors, avoiding restricted areas, bad weather, and other aircraft. These requirements are met by a combination of GPS, ground and airborne radar, internal instruments, and communications from ground controllers. In the original design of these types of equipment, little or no consideration was given to security; as a result, all of them are susceptible to information warfare attacks.

The accuracy and reliability of GPS and airborne radar, however, has led federal aviation authorities to consider implementing a system wherein ground controllers and published restrictions would no longer determine altitude, speed, clearance distances, and other flight parameters. Instead, pilots would have the option to choose any flight parameter that they believed to be safe. This new system is intended to increase the number of flights that can safely traverse the limited airspace. It is undoubtedly capable of doing so, but at the same time, it will greatly increase the dangers of flight should information warfare be waged against airplanes and the aviation infrastructure.

14.3.6.2 Railroads. Less so than airplanes, but not to a negligible degree, trains are possible targets of IW. Train movements; switch settings, communications between engineers, trainmen, and control centers are all carried on by insecure radio communications and wired lines. Attacks against any or all of these can prevent the railroads from carrying out their important functions, possibly by causing disastrous wrecks.

14 · 12 INFORMATION WARFARE

14.3.6.3 Trucking. The great majority of domestic goods shipments are carried by tractor-trailer trucks. Foodstuffs, especially, depend on this relatively fast, reliable means of transportation. If even a short disruption were to be caused by IW, untold quantities of foodstuffs would rot in the fields, as would additional stockpiles awaiting distribution from central warehouses. Data for scheduling, routing, locating trucks, setting times and locations of pickup and delivery, and performing maintenance could be prevented from reaching their destinations.

14.3.6.4 Shipping. Ships are indispensable means for transporting vast quantities of materials over long distances. Navigational data, such as position, speed, course to steer, and estimated time of arrival, are a few of the parameters determined by computers and GPS on virtually every ship afloat. Conventional radar and communications by VHF and high-frequency radio are in common use, with satellite communications becoming more prevalent, despite an early start that met with technical and economic difficulties.

Radar and communications jamming are old established weapons of IW, as is interception of critical information. Little attention has been paid to security in designing or operating this equipment, and that places ships at great risk, as does the threat of physical attacks.

14.3.6.5 Other Transportation Vulnerabilities. Recognizing the importance of transportation to a nation's infrastructure, IW attackers could create wide-ranging disruptions if they were to intercept and successfully prevent receipt of critical information within the transportation industry. Recently, as a leader in new technology, the Port Authority of New York and New Jersey has begun converting to a wireless infrastructure at its many airports, train stations, bus terminals, tunnels, bridges, and shipping facilities. It requires no stretch of the imagination to predict what a determined attacker might accomplish in damaging or destroying such an infrastructure. The danger is especially great in light of the general lack of security from which wireless transmissions suffer.

When the World Trade Center (WTC) was destroyed by terrorist action, the Port Authority's offices in the WTC were completely destroyed, and more than 70 of its employees were officially listed as deceased or missing. Although that catastrophe pointed up the need for greater physical security, it also demonstrated how the Internet can be used in emergency situations. The Port Authority site, www.panynj.gov, was used to convey operational messages to the public as well as information for tenants, employees and prospective employees, vendors, suppliers, contractors, and the media.

14.3.7 Commerce. In 1924, in an address to the American Society of Newspaper Editors, President Calvin Coolidge said: "After all, the chief business of the American people is business. They are profoundly concerned with producing, buying, selling, investing, and prospering in the world. I am strongly of the opinion that the great majority of people will always find these are moving impulses of our life."¹⁸

Now, 90 years later at the time of writing, these statements are no less true. Producing, buying, selling, and investing are the commercial means by which U.S. citizens and guest workers can hope to achieve prosperity. Although not recognized earlier, infrastructure is the glue that ties these functions together and permits them to operate efficiently and economically.

GOALS AND OBJECTIVES 14 · 13

If these bonds were to be broken, American business would come to a virtual standstill; it is that reality which makes the commercial infrastructure so inviting a target. Without complete, accurate, and current information, no investors would put their money at risk, and no transactions would take place among producers, buyers, and sellers.

In a populace lacking food, utilities, prescription drugs, money, and other necessities, civil disorder would be widespread. With the breakdown of commerce and the citizenry's unwillingness or inability to perform their customary functions, government at every level might cease to operate. This, in turn, would make military defensive actions highly problematic, and an enemy that combined IW with conventional force attacks would be difficult to resist.

On a less catastrophic level, there have been several cases of deliberate stock manipulation by means of insertion of false information into various news channels; an enemy could cause significant disruption in the stock market by forcing a few key stocks into unwarranted declines. In addition, the widespread use of automated trading tools that respond to significant drops in specific shares or in particular aggregate stock indexes could precipitate major economic problems in the developed world.

14.3.8 Financial Disruptions. Money is the lifeblood of every developed nation. For an economy to be healthy, its money supply, like the body's blood supply, must be strong, healthy, and free flowing. For an IW attacker, disruptions in the enemy's money supply and in its free flow are important objectives. Likely targets in the financial infrastructure include payment systems, investment mechanisms, and banking facilities.

14.3.8.1 Payment Systems. Every government employee, every member of the armed forces, every office worker, factory hand, service worker, engineer, and retired person—in fact, almost every individual in the United States—depends on regular receipt of funds necessary for survival. Paychecks, dividends, welfare and unemployment benefits, commissions, payments for products, and fees for services comprise most of the hundreds of millions of daily checks, direct deposits, and wire transfers without which most people would be unable to purchase their essential needs—assuming that products and services were available to meet those needs.

The great majority of payroll systems are computerized. Many of them, including those of the federal and state governments, depend on a few centralized computer payroll services. Even if those services were not damaged by infrastructure attacks, the banks on which payroll funds are drawn might be. This would halt, or at least impede, the cutting of physical checks, the direct deposits, cash withdrawals, wire transfers, and any other means by which payments are made. Such a situation has never occurred within the United States except in small local areas and for only brief periods of time. No one can predict what the consequences would be for a widespread attack, and surely no one would want to find out.

For more on banking payment systems, see Section 14.3.8.3.

14.3.8.2 Investment Mechanisms. Various stock, bond, and commodity exchanges provide the principal means by which individual, institutional, and corporate entities easily and expeditiously can invest in financial instruments and commodity goods.

With few exceptions, each exchange has all of its computers and communications located within a single facility, with connections to tens of thousands of terminals

14 · 14 INFORMATION WARFARE

worldwide. Disruption in these systems would not have as disastrous an effect as would a payment system disruption, but it would not be long before a breakdown in investment mechanisms would produce a commercial meltdown.

Because of the vast sums of money involved, exchange systems largely have been hardened against intrusion, and some have remote, redundant facilities, but there have been instances where hardware and software problems as well as physical exploits have brought down an exchange infrastructure.

14.3.8.3 Banking. The banking industry is the foundation of the modern financial system and, by extension, both American and foreign capitalist economies. At some point, every important financial transaction is conducted through the banking system. As such it is vital to economic health. With the advent of information warfare, the electronic, interdependent nature of banking—and finance in general—combined with its critical nature, makes the banking system a likely target for a strategic attack against a country. This is a new viewpoint for an industry focused on crime, traditional financial crises, and the more recent phenomenon of low-level hacking. It is critical, however, that we master this viewpoint and adapt our banking industry to it, for the threats information warfare poses are different from traditional bank security threats and will increase as the age of information warfare develops. Focused correctly, a well-prepared attack could cause chaos throughout the international system.¹⁹

The ubiquitous banking system is as highly automated and as security conscious as any element of the world's infrastructure. With ATMs, online banking, funds-transfer networks, and check clearing, banks are integral to virtually every commercial transaction.

As an example of the scope of banking operations involving money transfers, FED-WIRE, operated by the Federal Reserve Board, serves approximately 9,000 depository institutions as of January 2012, providing transfers that are immediate, final, and irrevocable.²⁰ In 2011, it processed 127 million transactions with a total value in excess of \$663 trillion. The average daily volume in 2011 was over 506,000 transactions valued at more than \$2.6 trillion.²¹

The Clearing House Interbank Payment System (CHIPS) “is responsible for over 95% of USD cross-border transactions, and nearly half of all domestic wire transactions, totaling \$1.5 trillion daily [in 2012].”²²

The Society for Worldwide Interbank Funds Transfer (SWIFT) “is a member-owned cooperative through which the financial world conducts its business operations with speed, certainty, and confidence. More than 10,000 financial institutions and corporations in 212 countries trust us every day to exchange millions of standardized financial messages. This activity involves the secure exchange of proprietary data while ensuring its confidentiality and integrity.”²³ Information about dollar value is not made public, but the amounts are known to be huge and the traffic enormous. For example, by the end of September 2012, the annualized number of messages transferred for the year was over 339 million and the annualized volume of files transferred among institutions was over 22 million.²⁴

If any of these systems were to be attacked successfully, the consequences for the financial well-being of many nations would be disastrous. Despite intensive efforts to safeguard the networks, attacks could be launched against the central computers, the computers of each user, and the networks that connect them.

14.3.9 Medical Security. In hospitals, as in group and private medical practice, the primary functions are carried out in a decentralized mode, making large-scale

GOALS AND OBJECTIVES 14 · 15

attacks impracticable. However, ancillary functions, such as sending invoices to the government, to health maintenance organizations, and to individuals, for services provided, and placing orders for drugs and supplies, all require interconnections with centralized computers.

Although the medical profession is often slow to adopt new infrastructure elements, network-connected computers have been mandated at least for payments, and they are becoming increasingly popular for maintaining patient data, for research, and for other functions. There have been reports that hospital systems have been penetrated, with prescriptions switched and HIV-negative patients advised that their test results were positive.

See Chapter 71 of this *Handbook* for more detail about medical information security.

14.3.10 Law Enforcement. The objectives of law enforcement are to facilitate the apprehension of criminals and wrongdoers. To accomplish this, facilities in common use include computers in every squad car connected to precinct headquarters and networks that interconnect local, state, federal, and international databases. With local law enforcement, it is clear that jamming, or noise interference on emergency channels, or denial of computer services would greatly exacerbate the effects of physical attacks. At worst, a state of panic and chaos might ensue.

Another attack on law enforcement could be flash crowds—groups of people gathered into a single physical location through instructions sent electronically.²⁵ One commentator wrote in 2004,

... [T]raining people to assemble on command in large numbers at, say, shoe stores, piano showrooms or restaurants for no good reason other than the fun of being part of a huge crowd is a perfect setup for creating an army of willing, mindless drones who will congregate on command at the site of a terrorist attack or at places where their presence will interfere with response to criminal or terrorist activities. Want to rob a bank in peace and quiet? Set up a conflict between two instant crowds to draw the police to an instant riot.²⁶

14.3.11 International and Corporate Espionage. Espionage has been a recognized military activity since at least the biblical story of Joshua, one of 12 spies sent to explore the land of Canaan.²⁷ However, its application to civilian commerce dates only from the Industrial Revolution. Since then, industries and indeed nations have prospered to the extent that they could devise and retain trade secrets. In the United States, the unauthorized appropriation of military secrets has been legally proscribed since the country's inception, with penalties as severe as death, during wartime.

Only recently have economic espionage and the theft of trade secrets become the subjects of law, with severe penalties whether the law is broken within or outside of the United States or even via the Internet.

The Economic Espionage Act of 1996 was signed into law by President Clinton on October 11, 1996. Section 1832 provides that:

- (A) Whoever, with intent to convert a trade secret, that is related to or included in a product that is produced for or placed in interstate or foreign commerce, to the economic benefit of anyone other than the owner thereof, and intending or knowing that the offense will injure any owner of that trade secret, knowingly—
 - (1) Steals, or without authorization appropriates, takes, carries away, or conceals, or by fraud, artifice, or deception obtains such information;
 - (2) Without authorization copies, duplicates, sketches, draws, photographs, downloads, uploads, alters, destroys, photocopies, replicates, transmits, delivers, sends, mails, communicates, or conveys such information;

14 · 16 INFORMATION WARFARE

- (3) Receives, buys, or possesses such information, knowing the same to have been stolen or appropriated, obtained, or converted without authorization;
- (4) Attempts to commit any offense described in any of paragraphs (1) through (3); or
- (5) Conspires with one or more other persons to commit any offense described in any of paragraphs (1) through (3), and one or more of such persons do any act to effect the object of the conspiracy,

Shall, except as provided in subsection (b), be fined under this title or imprisoned not more than 10 years, or both. (b) Any organization that commits any offense described in subsection (a) shall be fined not more than \$5,000,000.²⁸

Although the foregoing lists all of the actions that are proscribed, it is not specific as to which assets are to be protected as trade secrets. For this, see the Defense Security Service paper, "What Are We Protecting?"²⁹ There, the five basic categories of People, Activities/Operations, Information, Facilities, and Equipment/Materials are expanded into 42 specific assets, with the admonition that every company official must clearly identify to employees what classified or proprietary information requires protection. Only if the owner has taken reasonable measures to keep such information secret, and the information derives actual or potential economic value from not being generally known to or readily obtainable through proper means, will the courts view it as a trade secret.

For further information on intellectual property, including trade secrets, see Chapters 11 and 42 in this *Handbook*.

14.3.12 Communications. Communications are the means by which all elements of a civilization are tied together. Any significant destruction of communications media would disrupt the most important segments of society. Without adequate communications, transactions and services would come to a complete halt. In the United States, communications have been disrupted frequently, but fortunately, the infrastructure has been so vast and so diverse that the consequences have rarely been more than temporary. Even after the World Trade Center disaster of September 11, 2001, when Verizon's downtown telephone facilities centers were heavily damaged, service was restored within four days to the New York Stock Exchange and to other important users in the area.

Contrary to popular belief, the Internet is so widely used and concentrated in so few backbone points that a coordinated attack actually could destroy its functioning. For many years, backup facilities have included redundant computers and all of their associated peripherals, often in remote locations. Too often, however, alternate communications facilities are not provided. Unless this is rectified, the same disaster that brings down one installation could disable all.

14.3.13 Destabilization of Economic Infrastructure. A major difference between wealthy, developed nations and poor, undeveloped countries lies in the strength of their economic infrastructures. The existence of strong capital markets, stable banking and lending facilities, and efficient payment processes, all tied together by fast, technically advanced communications capabilities, is essential to healthy, growing economies.

At opposite ends of this spectrum lie Afghanistan and the United States. The perpetrators of the attacks on the World Trade Center and the Pentagon, identified as Osama bin Laden and his Al-Qaeda organization, operating out of Afghanistan, chose as their targets the symbols and the operating centers of America's military operations and of its economic infrastructure.

SOURCES OF THREATS AND ATTACKS 14 · 17

The recession of the late 2000s and early 2010s has illustrated the vulnerability of global economic systems to disruption. With hundreds of thousands thrown out of work and with investment capital drying up, the entire economic infrastructure of the world has already suffered a great blow without direct cyberattacks. Every effort must be bent toward preventing attacks that imperil the economic infrastructure of the world. Security can no longer be the duty of a few technical people; it has become everyone's responsibility.

14.4 SOURCES OF THREATS AND ATTACKS. The actual and potential originators of information warfare are numerous and powerful. One need not be paranoid to feel that an attack may come from any direction. This section lists sources that have already proven their capabilities for conducting cyberwar.

14.4.1 Nation-States. U.S. military preparations for cyberwar have been described in Section 14.3.2. This section details some of the measures that another great power—China—is effecting toward the same ends. Most of the material is from a paper entitled “Like Adding Wings to the Tiger: Chinese Information War Theory and Practice.”³⁰

14.4.1.1 China and Information Warfare. Although China is a nuclear power, it does not yet have the arsenal necessary to threaten a superpower like the United States. However, it can do so with its IW forces; adding wings to the tiger makes it more combat worthy. Nor is Chinese IW entirely theoretical. On August 3, 2000, the *Washington Times* reported that hackers suspected of working for a Chinese government institute took large amounts of unclassified but sensitive information from a Los Alamos computer system. A spokesman stated that “an enormous amount of Chinese activity hitting our green, open sites” occurs continuously.³¹

According to an article in the Chinese Armed Forces newspaper, the *Liberation Army Daily*, their first attack objectives will be the computer networking systems that link a country’s political, economic, and military installations, as well as their general society.³² A further objective will be to control the enemy’s decision-making capability in order to hinder coordinated actions.

Expanding on Mao Zedung’s theory of a People’s War, IW can be “carried out by hundreds of millions of people using open-type modern information system.”³³ In this war, combatants can be soldiers or teenagers, or anyone who has a computer as a weapon.³⁴ Ironically, China, with its long-standing fear of outside information as a possible spur to counterrevolutionary action, now views arming large numbers of intelligent people with computers and access to the Internet as a necessary survival measure. It remains to be seen just how many personal computers will be made available and how China will ensure that they will be used only as the government intends.

The “Annual Report to Congress on the Military Power of the People’s Republic of China” from the U.S. Department of Defense has been issued every year since 2002. Reading through all the reports provides valuable perspective on the DoD view of information warfare capabilities of the People’s Republic of China (PRC) and the People’s Liberation Army (PLA). The 2011 edition included the following analysis:

- Cyberwarfare Capabilities. In 2010, numerous computer systems around the world, including those owned by the U.S. Government, were the target of intrusions, some of which appear to have originated within the PRC. These intrusions were focused on exfiltrating information. Although this alone is a serious concern, the accesses and skills required for these intrusions are similar to those necessary to conduct computer network attacks. China’s 2010 *Defense*

14 · 18 INFORMATION WARFARE

White Paper notes China's own concern over foreign cyberwarfare efforts and highlighted the importance of cybersecurity in China's national defense.

- Cyberwarfare capabilities could serve PRC military operations in three key areas. First and foremost, they allow data collection through exfiltration. Second, they can be employed to constrain an adversary's actions or slow response time by targeting network-based logistics, communications, and commercial activities. Third, they can serve as a force multiplier when coupled with kinetic attacks during times of crisis or conflict.
- Developing capabilities for cyberwarfare is consistent with authoritative PLA military writings. Two military doctrinal writings, *Science of Strategy* and *Science of Campaigns*, identify information warfare (IW) as integral to achieving information superiority and an effective means for countering a stronger foe. Although neither document identifies the specific criteria for employing computer network attack against an adversary, both advocate developing capabilities to compete in this medium.
- The *Science of Strategy* and *Science of Campaigns* detail the effectiveness of IW and computer network operations in conflicts and advocate targeting adversary command and control and logistics networks to impact their ability to operate during the early stages of conflict. As the *Science of Strategy* explains,
—In the information war, the command and control system is the heart of information collection, control, and application on the battlefield. It is also the nerve center of the entire battlefield. [Emphasis ours.]³⁵

In parallel with its military preparations, China has increased diplomatic engagement and advocacy in multilateral and international forums where cyberissues are discussed and debated. Beijing's agenda is frequently in line with the Russian Federation's efforts to promote more international control over cyberactivities. China has not yet agreed with the U.S. position that existing mechanisms, such as International Humanitarian Law and the Law of Armed Conflict, apply in cyberspace. China's thinking in this area is evolving as it becomes more engaged.”

14.4.1.2 Strategies. The People's Liberation Army (PLA) with 1.5 million reserve troops has been carrying out IW exercises on a wide scale. One such exercise, in Xian Province, concentrated on conducting information reconnaissance, changing network data, releasing information bombs, dumping information garbage, disseminating propaganda, applying information deception, releasing clone information, organizing information defense, and establishing spy stations.³⁶ The antecedents of these tactics can be found in a book of unknown authorship, first mentioned about 1,500 years ago, entitled *The Secret Art of War: The 36 Stratagems*. Strategy 25 advises:

Replace the Beams with Rotten Timbers. Disrupt the enemy's formations, interfere with their methods of operations, change the rules which they are used to following, and go contrary to their standard training. In this way you remove the supporting pillar, the common link that makes a group of men an effective fighting force.³⁷

The 36 stratagems deserve close study; many of them are obviously in use even today by China and others. For example, strategy 3 says:

Kill with a Borrowed Sword. When you do not have the means to attack your enemy directly, then attack using the strength of another.

Lacking the weapons to attack the United States directly, the perpetrators of the WTC attack used the airliners belonging to their targets.

Strategy 5 says:

Loot a Burning House. When a country is beset by internal conflicts, when disease and famine ravage the population, when corruption and crime are rampant, then it will be unable to deal with an outside threat. This is the time to attack.

SOURCES OF THREATS AND ATTACKS 14 · 19

Some of the strategies might well be employed by the United States. For example, strategy 33 advises:

The Strategy of Sowing Discord. Undermine your enemy's ability to fight by secretly causing discord between him and his friends, allies, advisors, family, commanders, soldiers, and population. While he is preoccupied settling internal disputes his ability to attack or defend, is compromised.

To accomplish this, IW may prove to be an effective weapon.

14.4.1.3 Training. Several high-level academies and universities have been established to conduct IW instruction for the PLA. In addition, training is planned for large numbers of individuals to include:

- Basic theory, including computer basics and application, communications network technology, the information highway, and digitized units
- Electronic countermeasures, radar technology
- IW rules and regulations
- IW strategy and tactics
- Theater and strategic IW
- Information systems, including gathering, handling, disseminating, and using information
- Combat command, monitoring, decision making, and control systems
- Information weapons, including concepts, principles of soft and hard destruction, and how to apply these weapons
- Simulated IW, protection of information systems, computer virus attacks and counterattacks, and jamming and counterjamming of communications networks³⁸

It is doubtful that all of these training objectives have been accomplished, but there seems to be a major commitment to do so, and sooner rather than later.

China and the United States are only two of the nations that are openly preparing for, and actually engaged in, information warfare. It is obvious that many others are similarly involved and that these measures, combined with conventional weapons or weapons of mass destruction, have the potential to elevate warfare to a destructive level never before possible and hardly conceivable.

14.4.2 Cyberterrorists

“Cyberterrorism” means intentional use or threat of use, without legally recognized authority, of violence, disruption, or interference against cybersystems, when it is likely that such use would result in death or injury of a person or persons, substantial damage to physical property, civil disorder, or significant economic harm.³⁹

Cyberterrorists, those who engage in cyberterrorism, generally are able to carry out the same sort of cyberwar as nation-states; in fact, they may be state-sponsored. The major difference is that terrorist attacks are usually hit-and-run, where nations are capable of sustained and continuous operations. Although conventional warfare always was carried out in an overt fashion, it is the nature of IW that it can be engaged in without a declaration of war and without any clear indication of who the attacker actually is. In fact, it may not be recognized that a war is being conducted; it may

14 · 20 INFORMATION WARFARE

seem only that a series of unfortunate, unconnected natural failures of computers and communications are disrupting an economy.

Terrorists, especially when state-sponsored, would be very likely to conceal their IW activities in this manner, so as to avoid the retribution that would inevitably follow. However, some terrorists would publicly take credit for their actions, in order to bolster their apparent strength and to gather added support from like-minded individuals and organizations.

The seriousness of terrorist threats after 9/11 resulted in Executive Order 13228 of October 8, 2001, establishing the Office of Homeland Security and the Homeland Security Council.⁴⁰ The mission of the Office was to “develop and coordinate the implementation of a comprehensive national strategy to secure the United States from terrorist threats or attacks.” Its function was “to coordinate the executive branch’s efforts to detect, prepare for, prevent, protect against, respond to, and recover from terrorist attacks within the United States.”

The Department of Homeland Security was mandated by Congress on January 24, 2003, and was fully formed on March 1, 2003. Celebrating its tenth anniversary in 2013, the department employs more than 200,000 people dedicated to fulfilling its mission.

On February 15, 2005, Michael Chertoff was sworn in as the second secretary. His five goals:

1. Protect our Nation from Dangerous People
2. Protect our Nation from Dangerous Goods
3. Protect Critical Infrastructure
4. Strengthen our Nation’s Preparedness and Emergency Response Capabilities
5. Strengthen and Unify Operations and Management^{41²⁶}

On April 30, 2008, Secretary Chertoff, recognizing new realities, said:

[T]he technology of the 21st Century is changing so rapidly that many of our rules and procedures, which were built at a time that we had a certain kind of communication system and a certain kind of analog set of processes, that legal structure seems woefully inadequate to a digital age when the movement of communications is not rooted in any one place and when it's very difficult to take the concepts which made a lot of sense in the days of the rotary telephone and apply them in the world of voice over internet protocols.⁴²

In March 2010, a report on comments by FBI Director Robert Mueller and former White House “terrorism czar” Richard Clarke included this summary:

“As you well know, a cyber-attack could have the same impact as a well-placed bomb,” Mueller said. “In the past 10 years, Al-Qaeda’s online presence has become as potent as its in-world presence.”

Al-Qaeda uses for the Internet range from recruiting members and inciting violence to posting ways to make bio-weapons and forming social-networks for aspiring terrorists, according to Mueller. “The cyber-terrorism threat is real and rapidly expanding,” Mueller said. “Terrorists have shown a clear interest in hacking skills and combining real attacks with cyber attacks.”

Threats are also rising from online espionage, with hackers out for source code, money, trade, and government secrets, according to the FBI. “Every major company in the U.S. and Europe has been penetrated—it’s industrial warfare,” said Richard Clarke, who was a White House adviser under three prior U.S. presidents. “All the little cyber-devices that the companies here sell have been unable to stop them. China and Russia are stealing petabytes of information.”

WEAPONS OF CYBERWAR 14 · 21

Clarke, now a partner at Good Harbor Consulting firm, was among the RSA panelists discussing cyber-warfare. “Nation states have created cyber-warfare units. They are preparing the battlefield,” Clarke said. “We have the governments of China and Russia engaging in daily activities successfully that the U.S. government and private industry are not stopping and they are stealing anything worth stealing....”

Mueller urged computer security professionals to join in a united, international alliance with law enforcement agencies to battle enemies in cyberspace. He credited such teamwork with resulting in the recent arrest of three men in Spain suspected of running a network of nearly 13 million computers secretly infected with malicious software and used for nefarious deeds.

Mueller called on victims of cyber-attacks to break the pattern of remaining silent out of fear that reporting crimes would hurt their positions in the marketplace. “Maintaining the code of silence will not benefit you or your clients in the long run,” Mueller said. “We must continue to do everything we can together to minimize and stop these attacks.”⁴³

The challenges faced by the Department of Homeland Security are multitudinous and complex. Whether it proves effective in reducing or eliminating terrorism within the United States will depend on solving the problems of overlapping authorities, inertia, incompatible databases, turf wars, funding, management, the predictability of terrorist actions, and a host of political and technological issues.

14.4.3 Corporations. The threats aimed at or directed by corporations are far less deadly than those of the military or of terrorists, but they are no less pervasive. Thefts of data, denial of service, viruses, and natural disasters traditionally have been at the heart of individual corporate security concerns. These concerns have not abated, but to them have been added fears that attacks on large segments of the information infrastructure are more likely to create damage than is an attack against any single enterprise. To guard against this, every installation should operate behind strong firewalls and effective access controls.

In the wake of the September 11 attacks, Richard Clarke, who had been National Coordinator for Security, Infrastructure Protection, and Counterterrorism since May 1998, was appointed to a new post. As special advisor to the president for cyberspace security, Mr. Clarke warned that terrorists are out to hurt our economy and that they can use viruses in massive, coordinated attacks against corporate IT systems. He recommended, at a minimum, that disaster recovery plans include near-online, offsite backup facilities and redundant communications paths.

14.5 WEAPONS OF CYBERWAR. The weapons used in information warfare have existed for many years, but newer and more malevolent versions are produced with increasing frequency. For this reason, system security cannot be considered as static, but rather as part of an ongoing process that must be continuously monitored and strengthened. This section briefly describes the most common and most dangerous IW weapons, with references to other chapters where more detailed information is available.

14.5.1 Denial of Service and Distributed Denial of Service. Denial of service (DoS) and distributed denial of service (DDoS) are means by which computers, network servers, and telecommunications circuits can be partially or completely prevented from performing their designated functions. Any computer element that has been designed for a specific maximum capacity, if flooded by messages or data inputs that greatly exceed that number, can be slowed or even brought to a complete halt.

14 · 22 INFORMATION WARFARE

A DoS attack is carried out by a single computer that has been programmed to overwhelm the target system's capacity, usually by generating, automatically, a very large number of messages. A DDoS attack is implemented by planting a small program on hundreds or thousands of unaware computers. At a signal from the attacker, all of the agents (sometimes called zombies or daemons) are caused to send many messages simultaneously, thus flooding the victim's system or preempting all of its bandwidth capacity.

On April 26, 2007, a page-one article in the *New York Times* reported on what some Estonian authorities described as the first war in cyberspace. It was precipitated by the removal from a park in Tallinn of a bronze memorial to the Soviet soldiers of World War II. It was believed, but not proven, that the Russian government, or individual activists, had used DDoS attacks to bring down computers propagating the Websites of the Estonian president, prime minister, and Parliament as well as of banks and newspapers. The attacks were finally brought under control with the help of experts from NATO, the European Union, the United States, Finland, Germany, Slovenia, and Israel. Details of many DoS and DDos attacks and the recommended defenses are contained in Chapter 18 of this *Handbook*.

14.5.2 Malicious Code. Malicious code includes viruses, worms, and Trojan horses, as described in Chapter 16. Mobile code, such as Java, ActiveX, and VBScript, was developed to increase the functionality of Websites, but all three, as described in Chapter 17, also can be used maliciously.

There have been innumerable instances where malicious code has been used to damage or deface Websites, both civilian and military. Apparently, all of these exploits have been perpetrated by single individuals or by very small groups of unaffiliated crackers. However, in the event of actual cyberwar, it seems certain that large groups of coordinated, technically knowledgeable attackers will attempt to wreak havoc on their opponents' infrastructures through the use of malicious code.

In 2012, news reports indicated that Flame, malware with a relationship to Stuxnet, may, similarly, have been developed by U.S. and Israeli cyberwarfare specialists. Flame was not used for sabotage but rather for data theft.⁴⁴ It was held to be responsible for stealing Iranian passwords, network descriptors, and even data files. Armed with this information, Stuxnet was allegedly used to destroy Iranian centrifuges enriching uranium as a component of atomic weapons. It accomplished this by controlling the controllers that set the speed of the centrifuges. When the speeds were set excessively high, they damaged or destroyed the centrifuges.

Because most of the hardware elements of these systems are commercially available, and because the software is readily duplicated, the threats to military and commercial applications are imminently capable of extensive and costly attacks. Just as U.S. military and governmental agencies, and most of their allies, are engaged in large-scale operations to develop defensive capabilities, it is essential that all commercial enterprises exert major efforts to do the same. Initiatives have begun to form close working relationships between government and the private sector. Also, industry groups have begun advocating relaxation of those laws that prohibit close cooperation between competitors. This will be necessary before information can be shared as required to strengthen the infrastructure. Similarly, groups are requesting that shared information be protected from those who would use the Freedom of Information Act to force disclosure.

Every prudent organization will support these initiatives and will work with appropriate government agencies and industry groups to ensure its own survival and the welfare of the country itself.

WEAPONS OF CYBERWAR 14 · 23

14.5.3 Cryptography. Military operations, since the earliest recorded times, have utilized cryptography to prevent critical information from falling into enemy hands. Today, information is a vastly more important resource than ever before, and the need for cryptography has increased almost beyond measure. Not only the military, but indeed every financial institution, every competitive commercial enterprise, and even many individuals feel impelled to safeguard their own vital information. At the same time, access to the secret information of enemies and opponents would provide inestimable advantages.

Recognizing this, powerful supercomputers, directed by mathematicians, theoretical scientists, and cryptographers, are being applied to improving the processes of encryption and decryption. The most notable achievement in the recent past was the British construction of a computerized device to break the German Enigma code. The information thus obtained has been widely credited with a significant role in the outcome of World War II.

The development of effective mechanisms for spreading computations over millions of personal computers has greatly reduced the time required for brute force cracking of specific encrypted messages; for example, messages encrypted using the 56-bit Digital Encryption Standard (DES) were decrypted in four months using 10,000 computers in 1997, 56 hours using 1,500 special-purpose processors in 1998, and 22 hours using 100,000 processors in 1999.⁴⁵

A major issue, yet to be resolved, is the strength of cryptographic tools that may be sold domestically or exported overseas. The contending forces include producers of cryptographic tools who believe that if the strength of their product is in any way restricted, they will lose their markets to producers in other countries with more liberal policies. Similarly, proponents of privacy rights believe that unbreakable cryptographic tools should be freely available.

The countervailing view is that virtually unbreakable cryptographic tools shipped overseas will inevitably find their way into the hands of unfriendly governments, which may use them in conducting cyberwars against us. Domestically, law enforcement agencies believe that they should have “back-door” entry into all cryptographic algorithms, so that they may prevent crimes as wide-ranging as embezzlement, drug trafficking, and terrorism.

As domestic crimes and terrorist attacks grow in number and intensity, it seems certain that at least a few civil liberties, including privacy rights, may be infringed. The hope is that an optimum balance will be struck between the need for security and the core values of our democracy.

For more on privacy in cyberspace, see Chapter 69 in this *Handbook*.

14.5.4 Psychological Operations. Psychological operations (PSYOP) may be defined as planned psychological activities directed to enemy, friendly, and neutral audiences in order to influence their emotions, motives, attitudes, objective reasoning, and behaviors in ways favorable to the originator. The target audiences include governments, organizations, groups, and individuals, both military and civilian.

One of the most potent weapons in information warfare, PSYOP attempts to:

- Reduce morale and combat efficiency within the enemy’s ranks
- Promote mass dissension within, and defections from, enemy combat units and/or revolutionary cadres

14 · 24 INFORMATION WARFARE

- Support our own and allied forces cover and deception operations
- Promote cooperation, unity, and morale within one's own and allied units, as well as within friendly resistance forces behind enemy lines⁴⁶

The information that accomplishes these ends is conveyed via any media: by printed material such as pamphlets, posters, newspapers, books, and magazines, and by radio, television, personal contact, public address systems, and of increasing importance, through the Internet.

A classic example of successful PSYOP application was the deception practiced prior to the Allied invasion of the European mainland. Through clever “leaks,” false information reached Germany that General Patton, America’s most celebrated combat commander, was to lead an army group across the English Channel at Pas de Calais. As a consequence, German defensive forces were concentrated in that area. For weeks after the Normandy invasion was mounted, Hitler was convinced that it was just a feint, and he refused to permit the forces at Calais to be redeployed. Had this PSYOP failed, and had more of Germany’s defensive forces been concentrated in Normandy, the Allied landing forces might well have been thrown back into the sea.

Although generally considered not to involve a PSYOP action, the September 11 attacks and the subsequent spread of anthrax spores made clear that a physical action can have the greatest and most far-reaching psychological effects. Beyond mourning the death of almost 3,000 innocent civilians, the new sense of vulnerability and powerlessness caused great psychological trauma throughout the nation and much of the Western world. The full consequences to the travel, entertainment, and hospitality industries, as well as to every segment of the world economy, are likely to be both disastrous and long-lasting.

A major, integrated, expert PSYOP mission to restore morale and encourage behavior can halt or reverse a downward spiral, but worldwide recessions and acts of nature, such as cyclones, hurricanes, and earthquakes, can do more than PSYOP actions to demoralize a nation.

14.5.5 Physical Attacks. Prior to September 11, 2001, physical attacks, as a part of cyberwar, were generally considered in the same light as attacks against any military objective, and defensive measures were instituted accordingly. In the civilian sector, starting with student attacks against academic computers in the 1960s and 1970s, there have been occasional reported physical attacks against information processing resources. Although access controls have been almost universally in place, their enforcement often has been less than strict.

Another indication of the susceptibility of the information infrastructure to physical attack is the prevalence of “backhoe attacks” in which construction crews accidentally slice through high-capacity optic cables used for telecommunications and as part of the Internet backbones.⁴⁷ The signs indicating where not to dig can serve as markers for those targeting single points of failure.

A related vulnerability is undersea telecommunications cables, which are unprotected against accidental—or deliberate—damage from ship anchors and from other objects or tools. Breaks in these cables can interrupt the Internet and telephone networks on a global scale.⁴⁸

The destruction of the WTC and a portion of the Pentagon have brought the possibility of additional physical attacks very much into the forefront of cyberwar thinking, for both the military and the civilian infrastructures. Car bombings and packaged bombs

DEFENSES 14 · 25

had become almost commonplace, especially in the Mideast. Successful attacks had been launched against U.S. embassies and troop barracks, as well as against Israel, England, Spain, and France. To guard against such actions, perimeter defenses were widened, and in some areas personal searches at strategic points were instituted.

These defenses have proven to be of limited value, and suicide bombers seem to be increasing in numbers and in the effectiveness of their weapons. The use of commercial aircraft, fully loaded with fuel, as manned, guided missiles was apparently never considered prior to 11 September. After that date, there has been widespread recognition that protective measures must be taken that will prevent a recurrence of those tragic events. Airport security has become a direct federal responsibility, under a new Transportation Security Administration in the Department of Transportation. On November 11, 2001, President Bush signed a bill that requires all airport baggage screeners to be U.S. citizens and to undergo criminal background checks before becoming federal employees. At many airports, security is provided by private contractors. The protective measures in common use are considered to be pointless, inconvenient, and ineffective by many travelers. Although even minimal safeguards against known weapons are being debated, there appears to be little thinking directed toward other types of attacks that might even now be in the planning stage.

14.5.6 Biological and Chemical Weapons and Weapons of Mass Destruction. Although the use of these weapons can affect every element of society, they have a particular potency in destroying the infrastructure of a targeted nation. The WTC attacks have had long-lasting psychological effects, but the results of the anthrax dissemination may be even more deeply traumatic. Already, the presence of anthrax spores has interfered with the functioning of the Congress, the Supreme Court, the U.S. Postal Service, hospitals, and other institutions. Although the furor over these attacks, as well as their incidence, has dissipated, there may be even more such attacks in the future. Unless any future culprit is apprehended quickly, and countermeasures taken immediately, damage to the infrastructure could be extensive.

14.5.7 Weapons Inadvertently Provided. There are many widespread vulnerabilities to computer systems that are not created as weapons, but whose presence makes the targets of cyberwar highly vulnerable. Poor software designs and inadequate quality control create opportunities for attackers to damage or destroy information, and the information systems themselves. Chapters 38 to 40 of this *Handbook* are especially useful in identifying and eliminating these sources of security vulnerabilities.

14.6 DEFENSES. A variety of defenses may be employed both to prevent attacks and to mitigate their effects. Because each of these defenses may have only limited utility, it is evident that new and more effective defenses must be developed.

14.6.1 Legal Defenses. As a defense against IW attacks or as a framework for apprehending and prosecuting attackers, the international legal system has been generally ineffective. The reasons for this include:

- Information warfare is not prohibited under the United Nations (UN) Charter, unless it directly results in death or property damage.
- Laws that are not recognized and enforced lose their power to compel actions.
- There is little or no police power to enforce those few laws that do exist.

14 · 26 INFORMATION WARFARE

- The issue of sovereignty as it relates to transborder communications is unresolved.
- Neither the United States nor any other major power has pressed for international laws to govern information warfare. This may be attributed to the fact that such laws, while desirable for defense, would impair the nation's own offensive operations.
- Many nations do not recognize cyberwar attacks as criminal actions.
- In many lands, political considerations determine judicial outcomes.
- Few countries support extradition of their citizens even when indicted for terrorist or criminal activities.
- Terrorists, drug cartels, the international mafia, and even individual hackers have every reason to circumvent the law, and usually possess the resources that enable them to do so.
- Identifying attackers may be difficult or even impossible.
- New technologies arrive at a rate much faster than appropriate legislation.

Further acting to constrain law as a deterrent is the fact that there has been no universal acceptance of definitions for IW-relevant terminology: Attacks, acts of war, aggression, hostilities, combatants, crimes, criminals—all remain vague concepts. Until such terms, as applied to IW, are clearly defined, there can be no legal strictures against them.

The difference between acceptable and unacceptable targets is obscured by the dual-use, civilian and military, characteristics of infosystems and infrastructures. Similarly, it is difficult to condemn denial of service, when peacetime boycotts and economic sanctions are widely applied to further economic or political ends.

Clearly, legal defenses against cyberwar are inadequate at this time. Whether the United States will pursue effective international legislation remains doubtful, until the question of building adequate defenses, without hobbling offensive operations, is resolved.

14.6.2 Forceful Defenses. If IW attacks are accepted as acts of war, the use of retaliatory military force would be highly likely. The strategic and tactical decisions that would follow are well beyond the scope of this chapter, but six considerations are relevant.

1. The United States is growing reluctant to engage in combat without the sanction of the United Nations and without the concurrence of major allies. If the provocation is limited to an IW attack, it may be difficult to build a coalition or even to avoid UN condemnation.
2. The identity of the attacker may be unclear. Even after the September 11 attacks, the United States had no enemy that admitted culpability. As a consequence, the United States could not declare war on any nation or state but could only declare a war on “terrorism.”
3. The attacker may be misidentified. Through the use of “spoofing” and routing an attack through unaware nations, the anonymous culprit may escape detection, while blame falls on an innocent victim.
4. There may be difficulty in determining whether a particular event is an act of information warfare or simply the result of errors, accidents, or malfunctions.
5. The attackers may not be a foreign government, against whom war can be declared, but a criminal organization, a disaffected group, activists, commercial competitors, or even individuals bent on mischief.

DEFENSES 14 · 27

6. The United Nations, and international sentiment in general, requires that military force only be used in response to armed attack and, further, that the response be proportional to the attack that provoked it.

In light of these considerations, it seems unlikely that information warfare, unless it results in catastrophic injuries and deaths, will be met by a forceful reaction. The top secret Presidential memo leaked on June 7, 2013, and discussed in Section 14.3.3 clearly indicates that armed actions may now be initiated without the impetus of injuries and deaths.

14.6.3 Technical Defenses. The technical defenses against IW are many and varied. Almost the entire contents of this volume are applicable to safeguarding against cyberwar attacks. These same measures can prove equally effective in defending against IW, criminals, activists, competitors, and hackers.

14.6.4 In-Kind Counterattacks. A cyberwar defense that has been used often is an in-kind counterattack, where flaming is met by flaming, DDoS by DDoS, site defacement by site defacement, and propaganda by propaganda. Recent examples include exchanges between Israelis and Arabs, Kashmiris and Indians, Serbs and Albanians, Indians and Pakistanis, Taiwanese and Chinese, and Chinese and Americans.

Although there may be personal satisfaction in originating or responding to such attacks, the net effect is usually a draw, and, therefore, in-kind attacks generally have been short-lived. In the future, such attacks may no longer be the output of only a few individuals, but may be mounted by large numbers of similarly minded cyberwarriors, organized into coordinated groups, with sophisticated tools and with covert or overt state sponsorship.

In that event, the asymmetric nature of the adversaries' infrastructures would be telling. Clearly, if the Taliban, for example, were to mount another full-scale cyber-terrorist attack against the United States with the help of their supporters throughout the world, the effects could be devastating. Although the United States might mount a highly sophisticated in-kind response, it probably would have no effect on the Taliban's organization, its economy, its military effectiveness, or its ability to carry out suicide missions, biological warfare, or other physical attacks. A great and powerful nation may lack the ability to destroy a small, primitive, almost nonexistent infrastructure.

A serious problem with any kind of counterattack is that the origins of cyberattacks through the Internet using Internet Protocol version 4 (IPv4) are easily masked, or *spoofed*, because of the lack of mandatory, verifiable source authentication. It is possible that poorly analyzed data about the supposed attackers could lead to a counterattack against innocent victims.

14.6.5 Integration of Cyberwarfare into Military Planning. The United States Cyber Command (USCYBERCOM) was initiated in 2009 as a subset of the United States Strategic Command and became fully operational in its Fort Meade headquarters in 2010.⁴⁹ Its mission and focus are defined as follows:

- Mission: USCYBERCOM is responsible for planning, coordinating, integrating, synchronizing, and directing activities to operate and defend the Department of Defense information networks and when directed, conducts full-spectrum military cyberspace operations (in accordance with all applicable laws and regulations) in order to ensure U.S. and allied freedom of action in cyberspace, while denying the same to our adversaries.

14 · 28 INFORMATION WARFARE

- Focus: The command is charged with pulling together existing cyberspace resources, creating synergy that did not previously exist and synchronizing war-fighting effects to defend the information security environment.
- The Command centralizes direction of cyberspace operations, strengthens DoD cyberspace capabilities, and integrates and bolsters DoD's cyberexpertise. USCYBERCOM improves DoD's capabilities to ensure resilient, reliable information and communication networks, counter cyberspace threats, and assure access to cyberspace. The command works closely with interagency and international partners in executing the cybermission.

USCYBERCOM has initiated training exercises for U.S. military forces. The first *Cyber Flag* exercise was held in 2011 at Nellis Air Force Base; the second was in 2012.⁵⁰ The 2012 "... exercise saw approximately 700 participants, up from last year's 300, and doubled the network size. All participants had a specific role to play, playing the part of a U.S. team or role-playing an adversary."

In January 2013, plans surfaced for a major investment by the Pentagon in cybersecurity:

The expansion would increase the Defense Department's Cyber Command by more than 4,000 people, up from the current 900, an American official said. Defense officials acknowledged that a formidable challenge in the growth of the command would be finding, training, and holding onto such a large number of qualified people.

The Pentagon "is constantly looking to recruit, train and retain world class cyberpersonnel," a defense official said Sunday.

"The threat is real and we need to react to it," said William J. Lynn III, a former deputy defense secretary who worked on the Pentagon's cybersecurity strategy.

As part of the expansion, officials said the Pentagon was planning three different forces under Cyber Command: "national mission forces" to protect computer systems that support the nation's power grid and critical infrastructure; "combat mission forces" to plan and execute attacks on adversaries; and "cyberprotection forces" to secure the Pentagon's computer systems.⁵¹

14.6.6 Cooperative Efforts. Although the United States has been moderately successful in building coalitions in support of military operations, it has shown little inclination to build an international consensus dealing with information warfare. This may be so because of the legal difficulties outlined in Section 14.6.1 or because any prohibitions against offensive cyberwar will limit United States options. Nevertheless, whether by treaty, convention, agreement, or UN directive, technical people, diplomats, and statesmen of all well-intentioned countries should work together to define unacceptable and harmful actions and to devise means for detecting, identifying, and punishing those who transgress.

In June 2013, the North Atlantic Treaty Organization (NATO) representatives participated in a discussion of "how best to defend against cyber threats."⁵² NATO posted this summary of its stance:

Against the background of rapidly developing technology, NATO is advancing its efforts to confront the wide range of cyber threats targeting the Alliance's networks on a daily basis. NATO's Strategic Concept and the 2012 Chicago Summit Declaration recognise that the growing sophistication of cyber attacks makes the protection of the Alliance's information and communications systems an urgent task for NATO, and one on which its security now depends.

In June 2011, NATO adopted a new cyber defence policy and the associated Action Plan, which sets out a clear vision of how the Alliance plans to bolster its cyber efforts. This policy

FURTHER READING 14 · 29

reiterates that any collective defence response is subject to decisions of the North Atlantic Council, NATO's principal political decision-making body.

The revised policy offers a coordinated approach to cyber defence across the Alliance. It focuses on the prevention of cyber attacks and building resilience. All NATO structures will be brought under centralised protection and NATO will enhance its capabilities to deal with the vast array of cyber threats it currently faces, including through integrating them into the NATO Defence Planning Process. This way Allies will ensure that appropriate cyber defence capabilities are included as part of their planning to protect information infrastructures that are critical for core Alliance tasks. The revised cyber defence policy also stipulates NATO's principles on cyber defence cooperation with partner countries, international organisations, the private sector and academia.⁵³

14.7 SUMMARY. The potential for information warfare to damage or destroy the infrastructure of any nation, any corporation, or, in fact, any civilian, governmental, or military entity is unquestionable. Until now, the only incidents have been isolated and sporadic, but the possibility of sustained, coordinated, simultaneous attacks is strong. If these attacks are combined with physical, chemical, or biological warfare, the effects are certain to be devastating.

Although the types of potential attackers, and the probable weapons they will use, are well known, the available defenses do not at this time offer any great assurance that they will be effective. The United States and many of its allies are engaged in great efforts to remedy this situation, but formidable obstacles are yet to be overcome. The military is generally better prepared than the civilian sector, but much of the military's infrastructure is woven into and dependent on transportation, communications, utilities, food production and distribution, and other vital necessities that are owned by private enterprises.

Recent terrorist attacks and the probability of future offensives should serve as an immediate impetus to devote whatever resources are needed to combat the threats to our way of life and, in fact, to our very existence.

14.8 FURTHER READING

- Armistead, E. L. *Information Operations: Warfare and the Hard Reality of Soft Power*. Potomac Books, 2004.
- Armistead, E. L. *Information Warfare: Separating Hype from Reality*. Potomac Books, 2007.
- Armistead, E. L. *Information Operations Matters: Best Practices*. Potomac Books, 2010.
- Arquilla, J., and D. Ronfeldt, eds. *In Athena's Camp: Preparing for Conflict in the Information Age*. RAND Corporation, 1997. Available free in parts as PDF files from http://rand.org/pubs/monograph_reports/MR880/
- Campen, A. D., and D. H. Dearth, eds. *Cyberwar 3.0: Human Factors in Information Operations and Future Conflict*. Fairfax, VA: AFCEA International Press, 2000.
- Clarke, R. A. *Cyber War: The Next Threat to National Security and What to Do About It*. Ecco, 2010.
- Cohen, F. *World War 3: We Are Losing It and Most of Us Didn't Even Know We Were Fighting in It—Information Warfare Basics*. Fred Cohen & Associates, 2006.
- Denning, D. E. *Information Warfare and Security*. Addison-Wesley, 1998.
- Erbschloe, M., and J. Vacca. *Information Warfare*. McGraw-Hill, 2001.
- Greenberg, L., S. E. Goodman, and K. J. Soo Hoo. *Information Warfare and International Law*. National Defense University Press, 1998.

14 · 30 INFORMATION WARFARE

- Hall, W. M. *Stray Voltage: War in the Information Age*. U.S. Naval Institute Press, 2003.
- Henry, R., and C. E. Peartree, eds. *The Information Revolution and International Security*. Center for Strategic and International Studies, 1998.
- Kahn, D. *The Codebreakers*. Scribner, 1996.
- Kramer, F. D., S. H. Starr, and L. Wentz, eds. *Cyberpower and National Security*. Potomac Books, 2009.
- Lesser, I. O., B. Hoffman, J. Arquilla, D. Ronfeldt, and M. Zanini. *Countering the New Terrorism*. RAND Project Air Force, 1999. Available free in parts as PDF files from http://rand.org/pubs/monograph_reports/MR989/
- Macdonald, S. *Propaganda and Information Warfare in the Twenty-First Century: Altered Images and Deception Operations*. Routledge, 2007.
- Marsh, R. T., chair. *Critical Foundations: Protecting America's Infrastructures. The Report of the President's Commission on Critical Infrastructure Protection*. 1997; www.fas.org/sgp/library/pccip.pdf
- Price, A., and C. A. Horner. *War in the Fourth Dimension: U.S. Electronic Warfare, from the Vietnam War to the Present*. London, UK: Greenhill Books/Lionel Leventhal, 2001.
- Rattray, G. J. *Strategic Warfare in Cyberspace*. MIT Press, 2001.
- Reveron, D. S. *Cyberspace and National Security: Threats, Opportunities, and Power in a Virtual World*. Georgetown University Press, 2012.
- Rid, T. *Cyber War Will Not Take Place*. Oxford University Press, 2013.
- Rosenzweig, P. *Cyber Warfare: How Conflicts in Cyberspace Are Challenging America and Changing the World*. Praeger, 2013.
- Schwartzau, W. *Information Warfare: Chaos on the Electronic Superhighway*, 2nd ed. Thunder's Mouth Press/Perseus Publishing Group, 1996.
- Zalmay, K., and J. P. White, eds. *Strategic Appraisal: The Changing Role of Information in Warfare*. McGraw-Hill, 1999.

14.9 NOTES

1. W. J. Clinton, "Critical Infrastructure Protection," Presidential Decision Directive 63, May 22, 1998, www.fas.org/irp/offdocs/pdd/pdd-63.htm
2. J. L. Brock, "Critical Infrastructure Protection: Fundamental Improvements Needed to Assure Security of Federal Operations," GAO/T-AIMD-00-7, Testimony before the Subcommittee on Technology, Terrorism and Government Information, Committee on the Judiciary, U.S. Senate, October 6, 1999, www.gao.gov/archive/2000/ai00007t.pdf
3. L. Wright, "Protecting the Homeland: Report of the Defense Science Board Task Force on Defensive Information Operations 2000 Summer Study, Vol. II." Office of the Undersecretary of Defense for Acquisition, Technology, and Logistics (March 2001), www.acq.osd.mil/dsb/reports/dio.pdf (URL inactive)
4. T. P. M. Barnett, "The Seven Deadly Sins of Network-Centric Warfare," United States Naval Institute Proceedings 125, No. 1 (January 1999): 36–39, www.usni.org/magazines/proceedings/1999-01/seven-deadly-sins-network-centric-warfare
5. G. G. Gilmore, "Navy-Marine Corps Intranet Girds for Cyber-Attacks," Armed Forces Press Service, July 6, 2001, www.defenselink.mil/news/newsarticle.aspx?id=44745

NOTES 14 · 31

6. U.S. Office of the Secretary of Defense, "Annual Report to Congress: Military and Security Developments Involving the People's Republic of China 2013," www.defense.gov/pubs/2013_China_Report_FINAL.pdf
7. Joint Chiefs of Staff, "Joint Doctrine for Information Operations," Joint Publication 3-13, 2006, www.dtic.mil/doctrine/jel/new_pubs/jp3_13.pdf (URL inactive)
8. W. S. Cohen, Annual Report to the President and the Congress: Secretary of Defense, 2001, www.dod.mil/execsec/adr2001/index.html, Chapter 8: "Information Superiority and Space," www.dod.mil/execsec/adr2001/Chapter08.pdf (URL inactive)
9. "Obama tells intelligence chiefs to draw up cyber target list—full document text: Eighteen-page presidential memo reveals how Barack Obama has ordered intelligence officials to draw up a list of potential overseas targets for US cyber attacks." *The Guardian*, June 7, 2013, www.guardian.co.uk/world/interactive/2013/jun/07/obama-cyber-directive-full-text
10. This section uses verbatim extracts of articles published by M. E. Kabay originally published in *Network World Security Strategies*. Used with permission.
11. Jesse Berst, "The Electricity Economy: New Opportunities from the Transformation of the Electric Power Sector," White Paper, Global Environment Fund & Global SmartEnergy, August 2008, 55, www.terrawatts.com/electricity-economy.pdf p. 12.
12. Berst, "The Electricity Economy," p. 19.
13. Robert T. Marsh, *Critical Foundations: Protecting America's Infrastructures*, U.S. Government, Washington, DC: President's Commission on Critical Infrastructure Protection, 1997, 192. www.fas.org/sgp/library/pccip.pdf
14. M. E. Kabay, "Attacks on Power Systems: Data Leakage, Espionage, Insider Threats, Sabotage," *Network World Security Strategies*, September 13, 2010, [www.mekabay.com/nwss/828c_attacks_on_power_systems_\(2\).pdf](http://www.mekabay.com/nwss/828c_attacks_on_power_systems_(2).pdf)
15. M. E. Kabay, "Attacks on Power Systems: Hackers, Malware," *Network World Security Strategies*, September 8, 2010, [www.mekabay.com/nwss/828c_attacks_on_power_systems_\(1\).pdf](http://www.mekabay.com/nwss/828c_attacks_on_power_systems_(1).pdf)
16. ENISA, "Stuxnet Analysis," Press Release, July 10, 2010, www.enisa.europa.eu/media/press-releases/stuxnet-analysis
17. ENISA, "New Report on Smart Grids Cyber Security Measures; A Risk-Based Approach Is Key To Secure Implementation, According to EU Agency ENISA," Press Release, December 19, 2012, www.enisa.europa.eu/media/press-releases/smart-grids-cyber-security-measures-a-risk-based-approach-is-key-to-secure-implementation
18. See www.presidency.ucsb.edu/ws/?pid=24180
19. S. M. Parker, "Information and Finance: A Strategic Target," CommSec, 1997, <http://all.net/books/iw/iwarstuff/www.commsec.com/security/infowarfare.htm>
20. Federal Reserve Financial Services, "All Fedwire Participants," 2012, <https://www.federalreserve.gov/fpddir.txt>
21. Federal Reserve Board, "Fedwire Funds Service 2011 Annual Summary," 2011, www.federalreserve.gov/paymentsystems/fedfunds_ann.htm
22. Clearing House Interbank Payments System, "About CHIPS," www.chips.org/about/pages/033738.php
23. SWIFT "Company Information," www.swift.com/about_swift/company_information/index.page

14 · 32 INFORMATION WARFARE

24. SWIFT, "SWIFT in Figures—FIN Traffic," 2012, www.swift.com/about_swift/company_information/swift_in_figures/archive/2012/SIF_2012_09.pdf (URL inactive)
25. Ann Zimmerman and Miguel Bustillo, "'Flash Robs' Vex Retailers." *Wall Street Journal*, October 21, 2011, <http://online.wsj.com/article/SB10001424052970203752604576643422390552158.html>
26. M. E. Kabay, "Critical Thinking and Disintermediation," 2007, www.mekabay.com/opinion/critical_thinking.pdf
27. Numbers 13:16, 17
28. Public Law 104-294, "Economic Espionage Act of 1996"; www.law.cornell.edu/usc-cgi/get_external.cgi?type=pubL&target=104-294
29. "Counterintelligence: What Are We Protecting?" Defense Security Service, 1998, www.dss.mil/portal>ShowBinary/BEA%20Repository/new_dss_internet/isp/count_intell/what_protecting.html or <http://tinyurl.com/5fwafb>
30. T. L. Thomas, "Like Adding Wings to the Tiger: Chinese Information War Theory and Practice," Foreign Military Studies Office, Fort Leavenworth, KS, 2000; www.iwar.org.uk/iwar/resources/china/iw/chinaiw.htm
31. B. Gertz, "Hackers Linked to China Stole Documents from Los Alamos," *Washington Times*, August 3, 2000, p. 1
32. Shen Weiguang, "Checking Information Warfare Epoch Mission of Intellectual Military," *Jiefangjun Bao*, February 2, 1999, p. 6, as translated and downloaded from the Foreign Broadcast Information System (FBIS) Web site on February 17, 1999; www.opensource.gov (registration restricted to U.S. federal, state and local government employees and contractors).
33. Wei Jencheng, "New Form of People's Warfare," *Jiefangjun Bao*, June 11, 1996, p. 6, as translated and reported in FBIS-CHI-96-159, August 16, 1996.
34. Shen Weiguang (1995). "Focus of Contemporary World Military Revolution—Introduction to Research in IW," *Jiefangjun Bao*, November 7, 1995, p. 6, as translated and reported in FBIS-CHI-95-239, December 13, 1995, pp. 22–27.
35. M. E. Kabay, "US DoD Annual Estimates of Information Warfare Capabilities and Commitment of the PRC 2002-2011," 2013, www.mekabay.com/overviews/dod_prc_iw.pdf
36. Qianjin Bao, December 10, 1999, provided by William Belk via email to Timothy L. Thomas. According to Mr. Thomas, Mr. Belk is the head of a skilled U.S. reservist group that studies China.
37. Quotation from S. H. Verstappen, *The Thirty-Six Strategies of Ancient China* (Books and Periodicals, 2000). As described at www.chinastrategies.com/home36.htm
38. Zhang Zhenzhong and Chang Jianguo, "Train Talented People at Different Levels for Information Warfare," *Jiefangjun Bao*, February 2, 1999, as translated and downloaded from FBIS Website on February 10, 1999.
39. A. D. Sofaer et al., "A Proposal for an International Convention on Cyber Crime and Terrorism," 2000, www.iwar.org.uk/law/resources/cybercrime/stanford/cisac-draft.htm
40. G. W. Bush, Executive Order Establishing Office of Homeland Security, 2001, <http://georgewbush-whitehouse.archives.gov/news/releases/2001/10/20011008-2.html>

NOTES 14 · 33

41. U.S. Department of Homeland Security, “The Secretary’s Five Goals,” 2008, www.dhs.gov/xabout/gc_1207339653379.shtm (URL inactive)
42. M. Chertoff, “Remarks by Secretary Michael Chertoff and President of the Supreme Court of Israel Dorit Beinisch to the Heritage Foundation’s Civil Rights and the War on Terror: Dilemmas and Challenges Event,” April 30, 2008, www.dhs.gov/xnews/speeches/sp_1209741455799.shtm (URL inactive)
43. “FBI: Cyber-Terrorism a Real and Growing Threat to U.S.” *Homeland Security News Wire*, March 5, 2010, www.homelandsecuritynewswire.com/fbi-cyber-terrorism-real-and-growing-threat-us
44. Lucian Constantin, “Report: Flame Part of US-Israeli Cyberattack Campaign against Iran.” *Network World*, June 20, 2012, www.networkworld.com/news/2012/062012-report-flame-part-of-us-israeli-260353.html
45. M. Curtin and J. Dolske, “A Brute-Force Search of DES Keyspace,” 1998, www.interhack.net/pubs/des-key-crack; “Cracking DES: Secrets of Encryption Research, Wiretap Politics & Chip Design—How Federal Agencies Subvert Privacy: Frequently Asked Questions (FAQ) About the Electronic Frontier Foundation’s ‘DES Cracker’ Machine,” Electronic Frontier Foundation, 1998, http://w2.eff.org/Privacy/Crypto/Crypto_misc/DESCracker/19980716_eff_des.faq or <http://tinyurl.com/68thws>; and “RSA Code-Breaking Contest Again Won by Distributed.Net and Electronic Frontier Foundation (EFF): DES Challenge III Broken in Record 22 Hours,” Electronic Frontier Foundation, 1999, http://w2.eff.org/Privacy/Crypto/Crypto_misc/DESCracker/HTML/19990119_deschallenge3.html or <http://tinyurl.com/5n3gqf>
46. E. Rouse, “Psychological Operations/Warfare,” date unknown; www.psywarrior.com/psyhist.html
47. K. Poulsen, “The Backhoe: A Real Cyberthreat,” *Wired*, January 19, 2006, www.wired.com/science/discoveries/news/2006/01/70040; also CGA “CGA DIRT Analysis and Recommendations for Calendar Year 2005,” Common Ground Alliance Damage Information Reporting Tool, 2005, www.commongroundalliance.com/TemplateRedirect.cfm?Template=/ContentManagement/ContentDisplay.cfm&ContentFileID=3269 or <http://tinyurl.com/43obmo>
48. K. Kratovac, “Ship’s Anchor Caused Cut in Internet Cable: Unusual Cuts Led to Disruptions in Services, Slowed Down Businesses,” MSNBC Technology and Science/Internet, February 8, 2008, www.msnbc.msn.com/id/23068571/
49. United States Strategic Command, “U. S. Cyber Command,” Fact Sheet, 2012, www.stratcom.mil/factsheets/Cyber_Command/
50. Scott McNabb, “AFCYBER Takes Part in Second USCYBERCOM Cyber Flag Exercise.” U.S. Air Force Space Command press release, November 21, 2012, updated November 29, 2012, www.afspc.af.mil/news/story.asp?id=123327388
51. E. Bumiller, “Pentagon Expanding Cybersecurity Force to Protect Networks Against Attacks,” *The New York Times*, January 27, 2013, www.nytimes.com/2013/01/28/us/pentagon-to-beef-up-cybersecurity-force-to-counter-attacks.html
52. NATO, “Collaborating against Cyber Threats,” in the Web page “Defending against Cyber Attacks,” NATO Website, June 4, 2013, www.nato.int/cps/en/natolive/75747.htm
53. NATO, “NATO and Cyber Defence,” NATO Website, www.nato.int/cps/en/SID-BB17E53B-F3456C51/natolive/topics_78170.htm

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 15

PENETRATING COMPUTER SYSTEMS AND NETWORKS

**Chey Cobb, Stephen Cobb, M. E. Kabay,
and Tim Crothers**

15.1 MULTIPLE FACTORS INVOLVED IN SYSTEM PENETRATION	15·1	15.3.4 Spying	15·20
15.1.1 System Security: More than a Technical Issue	15·2	15.3.5 Penetration Testing, Toolkits, and Techniques	15·20
15.1.2 Organizational Culture	15·2	15.3.6 Penetration via Websites	15·27
15.1.3 Chapter Organization	15·3	15.3.7 Role of Malware and Botnets	15·31
		15.3.8 Sophisticated Attackers	15·32
15.2 NONTECHNICAL PENETRATION TECHNIQUES	15·3		
15.2.1 Misrepresentation (Social Engineering)	15·3	15.4 POLITICAL AND LEGAL ISSUES	15·36
15.2.2 Incremental Information Leveraging	15·6	15.4.1 Exchange of System Penetration Information	15·36
		15.4.2 Full Disclosure	15·36
		15.4.3 Sources	15·38
		15.4.4 Future of Penetration	15·39
15.3 TECHNICAL PENETRATION TECHNIQUES	15·7		
15.3.1 Data Leakage: A Fundamental Problem	15·7	15.5 SUMMARY	15·40
15.3.2 Intercepting Communications	15·8	15.6 FURTHER READING	15·41
15.3.3 Breaching Access Controls	15·15	15.7 NOTES	15·42

15.1 MULTIPLE FACTORS INVOLVED IN SYSTEM PENETRATION. Although penetrating computer systems and networks may sound like a technical challenge, most information security professionals are aware that systems security has both technical and nontechnical aspects. Both aspects come into play when people attempt to penetrate systems. Both aspects are addressed in this chapter, which is not an instruction guide on how to penetrate systems, but rather a review of the methods and means by which systems penetrations are accomplished.

15 · 2 PENETRATING COMPUTER SYSTEMS AND NETWORKS

15.1.1 System Security: More than a Technical Issue. The primary non-technical factor in system security and resistance to system penetration is human behavior, which can defeat just about any technical security measure. More than anything else, security depends on human beings to understand and carry out security procedures. Consequently, information system (IS) security must be integral to the culture of any organization employing an information system. Without security, systems and networks will not be able resist attempts at penetration.

Often security is represented as a structure of concentric circles. Protection of the central, secured element is then dependent on the barriers imposed by each successive ring. These barriers can be physical or figurative, but the goal of IS security is to protect the integrity, confidentiality, and availability of information processed by the system. This goal is reached using identification, authentication, and authorization. *Identification* is a prerequisite, with each user required to proffer an identifier (ID) that is included in the authorization lists of the system to be accessed. *Authentication* consists of proving that the user really is the person to whom the ID has been assigned. *Authorization* consists of defining what a specific user ID, running specified programs, can legally do on the system. The security perimeter can be penetrated by compromising any of these functions. Chapters 28 and 29 in this *Handbook* discuss identification and authentication in detail.

The trend toward distributed and mobile computers, often utilizing the global networking capability of the Internet, makes it hard to know where to draw these concentric circles of protection. Indeed, the barriers to penetration need to be extended along lines of communication, encompassing end points of the network, which may be geographically dispersed.

15.1.2 Organizational Culture. An organization's general attitude toward security is the key to an effective defense against attack. Security is difficult to sell, especially to an organization that has never experienced a significant problem. (Ironically, the better the defenses, the less evidence there is of their utility.) A basic principle of security is that practitioners must act as if they are paranoid, continuously on guard against attacks from any direction. Many organizations view security precautions as an attack on the integrity of employees. Wearing badges, for example, sometimes is viewed as dehumanizing and offensive. This attitude leads to absurdities, such as having only visitors wear badges. If only visitors wear badges, then taking off the badge automatically reduces the likelihood that a dishonest intruder will be challenged.

Some individual employees also consider security precautions as personally offensive. For example, locking a terminal or workstation when leaving it for a few minutes may be seen as evidence of distrust of other employees. Refusing to allow piggybacking—that is, permitting several colleagues to enter a restricted area on one access card—may be seen as insufferably rude. Where employees are taught to be open and collegial, securing removable computer media and paperwork at night can seem insulting.

These conflicts occur because years of socialization, starting in infancy, are diametrically opposed to the tenets of information security. Politeness in a social context is a disaster in a secure area; for instance, piggybacking into a computer room impairs the accuracy of audit trails kept by the access-control computers. Lending someone a car is kind and generous, but lending someone a user ID and a personal password is a gross violation of responsibility. Chapters 49 and 50 in this *Handbook* discuss psychological aspects of changing corporate culture to support information security.

NONTECHNICAL PENETRATION TECHNIQUES 15 · 3

Carrying out effective security policies and procedures must resolve these conflicts between normal standards of politeness and the standards required in a secure environment. Organizations must foster open discussion of the appropriateness of security procedures, so that employees can voluntarily create a corporate culture conducive to protection of corporate information. Chapters 44, 45, 48, and 51 in this *Handbook* specifically discuss policy issues.

Beyond this, organizations need to be aware of the security posture and attitudes of those with whom they network. These days it is quite possible for one organization's system to be operated, or even owned, by another. And people from many different organizations may be using the same network. A culture of security must permeate all of the organizations that have access to a system, otherwise points of weakness will exist, thus increasing the probability that attempts to penetrate the system will succeed.

15.1.3 Chapter Organization. Section 15.2 looks at methods of tricking people into allowing unauthorized access to systems (Chapter 19 in this *Handbook* explores social engineering in more depth). Section 15.3 examines technical measures for overcoming security barriers and specific techniques (*exploits*) for penetration, while Section 15.4 describes legal and political aspects of system penetration.

15.2 NONTECHNICAL PENETRATION TECHNIQUES. Although the penetration of information systems is often portrayed as the work of the technically adept, many successful penetrations have relied on human factors, such as gullibility and venality. Both are exploited by would-be system penetrators.

15.2.1 Misrepresentation (Social Engineering). Social engineering relies on falsehood. Lies, bribes, and seduction can trick honest or marginally dishonest employees into facilitating a penetration. An attacker might trick an employee into revealing login and authentication codes or even into granting physical access to an otherwise secure site. System penetration can then be accomplished by numerous means, from walking up to an unsecured workstation, to installing Trojan code or a network packet-sniffing device. (Both of these technologies are discussed in more detail later in this chapter.)

15.2.1.1 Lying. Telling lies is a technique often used by persons intent on obtaining unauthorized access to a system. One can obtain valuable information about a system and its defenses by telling lies. Many lies work by playing on the natural human tendency to interpret the world by our internal model of what is most likely. Social psychologists call this model the *schema*. Well-dressed businesspeople who walk briskly and talk assertively are probably what they seem. In a phone conversation, a person who sounds exasperated, impatient, and rude when demanding a new password is probably an exasperated, impatient, and rude employee who has forgotten a password. Unfortunately, many criminals know, sometimes instinctively, how to exploit these interpretations to help get them into secured systems.

Another technique, often used in concert with lying, is to escape notice and avoid suspicion by simply blending in. The way we perceive, or fail to perceive, details are referred to by social psychologists as the figure-ground problem. The normal becomes the background, and the objects of our attention become figures standing out from the ground. The schema influences what is noticed; only deviations from expectation spark figure-ground discrimination. Criminal hackers take advantage of this effect by fading into the background while penetrating security perimeters.

15 · 4 PENETRATING COMPUTER SYSTEMS AND NETWORKS

15.2.1.2 Impersonating Authorized Personnel. Criminal hackers and unscrupulous employees call security personnel, operators, programmers, and administrators to request user IDs, privileges, and even passwords. (This is one reason that the telephone is a poor medium for granting security privileges; if staff members were trained to refuse requests made over the phone, many attempts to penetrate systems could be thwarted.) In sites where employees wear ID badges, intruders have a hard time penetrating physical security by posing as employees. However, physical security in these cases depends on the cooperation of all authorized personnel to challenge everyone who fails to wear a badge. This policy is critically important at entry points. To penetrate such sites physically, criminals must steal or forge badges or work with confederates to obtain real but unauthorized badges.

Sites where physical security includes physical tokens, such as cards for electronic access control, are harder for criminals to penetrate. They must obtain a real token, perhaps by theft or by collusion with an employee. Perimeter security depends on keeping the access codes up to date so that cards belonging to ex-employees are inactivated. Security staff must immediately deactivate all cards reported lost. In addition, it is essential that employees not permit *piggybacking*, the act of allowing another person, possibly unauthorized, to enter a restricted zone along with an authorized person. Too often, an employee, in an act of politeness, will permit others to enter a normally locked door as he or she exits. Once inside a building, criminals can steal valuable information that will allow later penetration of the computer systems from remote locations. This is often accomplished by impersonating third-party personnel.

Even if employees are willing to challenge visitors in business suits, it may not occur to them to interfere with people who look as if they are employees of an authorized support firm. For example, thieves often have succeeded in entering a secured zone by dressing like third-party computer technicians or office cleaners. Few employees will think of checking the credentials of a weary technician wearing grimy overalls, an authentic-looking company badge, a colorful ID card, and a tool belt. When a suitable-looking individual claims to have been called to run diagnostics on a workstation, many nontechnical employees will acquiesce at once, seizing the opportunity to grab a cup of coffee or to chat with colleagues. Minutes later, the thief may have copied sensitive files or installed a sniffing device (e.g., a keystroke recorder or a network packet sniffer). In one case known to one of the authors (MK), a criminal was given a workspace and a network connection in a large bank and allowed to work unmolested and unchallenged for several months on a “secret project.” It was only when an alert security guard realized that no one in the office knew who this person was that she challenged the intruder and broke the scam.

15.2.1.3 Intimidation. A technique related to impersonation of authorized or third-party personnel is intimidation. Someone claiming to be a person in a position of authority displays irritation or anger at delays in granting an unauthorized deviation from policy, such as communicating a password over the phone to a person of unauthenticated identity. The attackers indirectly or directly threaten alarming consequences (e.g., delays of critical repairs, financial losses, disciplinary actions) unless they are granted restricted information or access to secured equipment or facilities.

15.2.1.4 Subversion. People make moral choices constantly. There is always a conscious or unconscious balancing of alternatives. Criminal hackers try to reach

NONTECHNICAL PENETRATION TECHNIQUES 15 · 5

their goals by changing the rules so that dishonesty becomes more acceptable to the victim than honesty.

15.2.1.5 Bribery. A lot of industrial and commercial information has a black market value. The same is true of personally identifiable information that can be used to commit fraud and identity theft. The price of a competitor's engineering plans or customer database may be a year's salary for a computer operator responsible for making backups. There is little likelihood that anyone would notice the subverted operator copying a backup at 3:00 A.M. or a secretary taking an extra compact disc out of the office. Many organizations have failed to install software to prevent a manager sending electronic mail with confidential files to a future employer.

That industrial espionage, with or without state sponsorship, is a thriving business is a fact that is now widely—and sometimes quite openly—acknowledged.¹ Building a corporate environment in which employees legitimately feel themselves to be part of a community is a bulwark against espionage. When respect and a sense of exchange for mutual benefits inform the corporate culture, employees will rebuff spies or even entrap them, but the disgruntled employee whose needs are not addressed is a potential enemy.

15.2.1.6 Seduction. Sometimes criminal hackers and spies have obtained confidential information, including access codes, by tricking employees into believing that they are loved. This lie works well enough to allow access to personal effects, sometimes after false passion or drugs have driven the victim into insensibility. It is not unknown for prostitutes to seduce men from organizations that they and their confederates are seeking to crack. Rifling through customers' wallets can often uncover telltale slips bearing user IDs and passwords.

No one can prevent all such abuse. People who are enthralled by expert manipulators will rarely suspect that they are being used as a wedge through a security perimeter. Along with a general increase in security consciousness, staff members with sensitive codes must become aware of these techniques so that they may be less vulnerable. Perhaps then they will automatically reject a request for confidential information or access codes.

15.2.1.7 Extortion. Criminals can threaten harm if their demands are not met. Threaten someone's family or hold a gun to their head and few will, or should, resist a demand for entry to a secured facility or for a login sequence into a network. Some physical access-control systems include a duress signal that can be used to trigger a silent alarm at the monitoring stations. The duress signal requires a predetermined, deliberate action on the part of the person being coerced into admitting unauthorized personnel. This action may be adding an extra number to the normal pass code, pressing the pound sign (#) twice after entering the code, or entering 4357 (H-E-L-P) into the keypad. The duress signal covertly notifies security that an employee is being forced to do something unwillingly. Security can then take appropriate action.

15.2.1.8 Blackmail. Blackmail is extortion based on the threat of revealing secrets. An employee may be entrapped into revealing confidential data, for example, using techniques just described. Classic blackmail includes seduction followed by pictures in flagrante delicto, which the criminals then threaten to reveal. Sometimes a person can be framed by fabricated evidence; a plausible but rigged image of venality can ruin a career as easily as truth. Healthy respect for individuals and social bonds

15 · 6 PENETRATING COMPUTER SYSTEMS AND NETWORKS

among employees, supervisors, and management can make it difficult for blackmailers to succeed. If employees who are victims of a blackmail attempt feel they can inform management without suffering inappropriately negative consequences, the threat may be mitigated to a certain degree. Perhaps the last, best defense against blackmail is honesty. The exceptionally honest person will reject opportunities that lead to victimization through blackmail and will laugh at fabrications, trusting friends and colleagues to recognize lies when they hear them.

15.2.1.9 Insiders. Many of the world's largest and most daring robberies have, upon examination, turned out to be inside jobs. The same is true of system penetrations. Although many of the just-described techniques can be used to obtain help from the inside, some are made possible by people on the inside who decide, for whatever reason, to aid and abet criminal hackers. For example, a dishonest employee may actively seek to sell access for personal gain. Organizations should try to be alert to this eventuality, but there is very little defense against thoroughly dishonest employees when the only overt act needed to open the gates from the inside is to pass system credentials to an outsider. Chapter 13 in this *Handbook* specifically addresses insider crime.

15.2.1.10 Human Target Range. Organizations should not underestimate the range of targets at which the described techniques may be directed. Although the terms *employees*, *authorized personnel*, and *third-party personnel* are used in the preceding paragraphs, the target range includes all manner of vendors, suppliers, and contractors as well as all levels of employees—from software and hardware vendors, through contract programmers, to soft-drink vendors and cleaning staff. It may even include clients and customers, some of whom possess detailed knowledge of the organization's operations. Employees at every level are likely to be computer literate, although with varying degrees of skill. For example, it is quite possible that someone working as a janitor today knows how to operate a computer skillfully and may even know how to surf hacking sites on the Web and download penetration tools. Indeed, a janitor may have obtained the job specifically with the intent of engaging in industrial espionage, data theft, or sabotage.

In short, anyone who comes into contact with the organization has the potential to provide an attacker with information useful in the preparation and execution of an attack. The human targets of a social engineering attack may not, on an individual basis, possess or divulge critical information, but each may provide clues—pieces of the puzzle—an aggregation of which can lead to successful penetration and compromise of valuable data and resources. Use of this process is a hallmark of some of the most successful criminal hackers. The term *incremental information leveraging* was coined for this use of less valuable data to obtain more valuable data.²

15.2.2 Incremental Information Leveraging. By gathering and shrewdly utilizing small and seemingly insignificant pieces of information, it is possible to gain access to much more valuable information. This technique of incremental information leveraging is a favorite tool of hackers, both criminal and noncriminal. One important benefit of the tool that is particularly appreciated by criminal hackers is the low profile it presents to most forms of detection. By accumulating seemingly innocuous pieces of information over a period of time, and by making intelligent deductions from them, it is possible to penetrate systems to the highest level.

A prime example of this approach is seen in the exploits of Kevin Mitnick, who served almost five years behind bars for breaking into computers, stealing data, and

TECHNICAL PENETRATION TECHNIQUES 15 · 7

abusing electronic communication systems. Illegal acts committed by Mitnick include the 1981 penetration of Computer System for Mainframe Operations (COSMOS), a Pacific Bell facility in downtown Los Angeles. COSMOS was a centralized database used by many U.S. phone companies for controlling basic recordkeeping functions. Mitnick and others talked their way past a security guard and located the COSMOS computer room. They stole lists of computer passwords, operating manuals for the COSMOS system, and combinations to the door locks on nine Pacific Bell central offices. Mitnick later employed knowledge of phone systems and phone company operations to penetrate systems at Digital Equipment Corp. (DEC).

Since his release in January 2000, Mitnick has spoken about information security before Congress and at other public venues. He described social engineering as such a powerful tool that he “rarely had to resort to a technical attack.”³ As to technique, he stated, “I used to do a lot of improvising . . . I would try to learn their internal lingo and tidbits of information that only an employee would know.” In other words, by building up knowledge of the target, using a lot of information that is neither protected nor proprietary, it is possible to gain access to that which is both proprietary and protected. The power of incremental information leveraging is the equivalent of converting a foot in the door into an invitation to come inside.

Protection against incremental information leveraging and all other aspects of social engineering begins with employee awareness. Employees who maintain a healthy skepticism toward any and all requests for information provide a strong line of defense. Another powerful defense mechanism, highlighted by Mitnick, is the use of telephone recording messages, such as “This message may be monitored or recorded for training purposes and quality assurance.” An attacker who hears a message like this may think twice about proceeding with attempts to use voice calls to social engineer information from the target.

15.3 TECHNICAL PENETRATION TECHNIQUES. Technical penetration attacks may build on data obtained from social engineering, or they may be executed on a purely technical basis. Techniques used include eavesdropping, either by listening in on conversations or by trapping data during transmissions, and breaches of access controls (e.g., trying all possible passwords for a user ID or guessing at passwords). Weaknesses in the design and implementation of information systems, such as program bugs and lack of input validation, also may be exploited in technical attacks. Unfortunately, weaknesses of this nature abound in the realm of the Internet, even as more and more organizations increase their Internet connectivity, thus creating more and more potential penetration points.

15.3.1 Data Leakage: A Fundamental Problem. Unfortunately, for information security (INFOSEC) specialists, it is impossible, even in theory, to prevent the unauthorized flow of information from a secured region into an unsecured region. The imperceptible transfer of data without authorization is known as *data leakage*. Technical means alone cannot suppress data leakage.

Consider a tightly secured operating system or security monitor that prevents confidential data from being copied into unsecured files. Workstations are diskless, there are no printers, employees do not take disks into or out of the secured facility, and there are strict restrictions on taking printouts out of the building. These mechanisms should suffice to prevent data leakage.

Not really.

15 · 8 PENETRATING COMPUTER SYSTEMS AND NETWORKS

Anyone with a penchant for mnemonics or with a photographic memory could simply remember information and write it down after leaving the facility. And it is extremely difficult to prevent employees from writing notes on paper and concealing them in their clothing or personal possessions when they leave work. Unless employees are strip-searched, no guard can stop people with crib sheets full of confidential data from walking out of the building. Indeed, this is how Vasili Mitrokhin, head archivist of the KGB's First Chief Directorate, perpetrated the largest breach of KGB security ever, by smuggling thousands of handwritten copies of secret documents out of the KGB headquarters in Moscow in his shoes, socks, and other garments.⁴

Another means of data leakage is steganography, hiding valuable information in plain sight among large quantities of unexceptional information. For example, a corrupt employee determined to send a confederate information about a chemical formula could encode text as numerical equivalents and print these values as, say, the fourth and fifth digits of a set of engineering figures. No one is likely to notice that these numbers contained anything special. The more digitally inclined can use steganography software, freely available on the Internet, to hide data in image files.

The unauthorized transfer of information cannot be absolutely prevented because information can be communicated by anything that can fluctuate. Theoretically, one could transfer data to a confederate by changing the position of a window shade (slow but possible). Or one could send ones and zeroes by the direction of oscillation of a tape reel; or one could send coded information by the choice of music. Even if a building were completely sealed, it would still leak heat outward or transfer heat inward—and *that* would be enough to carry information. In practical terms, system managers can best meet the problem of data leakage by a combination of technical protection and effective management strategies.

It is also important to realize that a significant amount of data leakage occurs through innocuous intending communications from employees. It is common for employees to discuss small aspects of their jobs and work information without realizing the implications and ability of others to collate this information into far more substantial amounts. This has become particularly prevalent with the modern advent of social media. Significant amounts of information about an organization and its internal workings can be derived from employee social media pages. Consequently, social media has become a rich source of data for attackers wishing to target an organization. This approach is particularly common among attackers that target an organization through its employees using techniques such as spear-phishing.

Data loss prevention (DLP), discussed briefly in Chapter 13 in this *Handbook*, is now an established set of techniques, with numerous tools available to enforce restrictions on unauthorized data transfers to storage devices and external sites. Nonetheless, vigilance over human behavior and the effective configuration and real-time or at least frequent analysis of log records remain essential components of effective DLP.

15.3.2 Intercepting Communications. Criminal hackers and dishonest or disgruntled employees can glean access codes and other information useful to their system penetration efforts by monitoring communications. These might be between two workstations on a local area or wide area network, between a remote terminal and a host such as a mainframe, or between a client and a server on the Internet. Attackers can exploit various vulnerabilities of communications technologies. The shift to Transmission Control Protocol/Internet Protocol (TCP/IP)-based Internet communications over the last decade has brought many more communications streams into the target

TECHNICAL PENETRATION TECHNIQUES 15 · 9

range of would-be penetrators. For basic information about data communications, see Chapter 5 in this *Handbook*.

15.3.2.1 Wiretapping. Wiretapping consists of intercepting the data stream on a communications channel (even if that channel is not wire; e.g., fiber optic cable can also be tapped, as can wireless communications, although the latter sometimes are said to be sniffed rather than tapped).

15.3.2.2 Asynchronous Connections. Point-to-point data connections (e.g., using telephone modems or serial devices) have all but disappeared, but they are relatively easy to tap. Physical connection at any point on twisted pair or multiwire cables allows eavesdropping on conversations; if the line is being used for data communications, a monitor is easily attached to display and record all information passing between a node and its host. Asynchronous lines in large installations often pass through patch panels, where taps may not be noticed by busy support staff, as they manage hundreds of legitimate connections. Such communications usually use phone lines for distances beyond a few hundred meters (or about 1,000 feet).

Wiretappers must use modems configured for the correct communications parameters, including speed, parity, number of data bits, and number of stop bits, but these parameters are easy to find out by trial and error.

Countermeasures include:

- Physical shielding of cables and patch panels
- Multiplexing data streams on the same wires
- Encryption of data flowing between nodes and hosts

15.3.2.3 Synchronous Communications. Because synchronous modems are more complex than asynchronous models and because their bandwidths (maximum transmission speeds) are higher, they are less susceptible to attack, but they are not risk-free.

15.3.2.4 Dial-up Phone Lines. Used for both data and voice communications, dial-up lines supplied by local telephone companies and long-distance carriers are vulnerable to wiretapping. Law enforcement authorities and telephone company employees can install taps at central switching. Criminals can tap phone lines within a building at patch panels, within cabling manifolds running in dropped ceilings, below raised floors, or even in drywall. They also can tap at junction boxes where lines join the telephone company's external cables.

The same countermeasures apply to phone lines as to asynchronous or synchronous data communications cables.

15.3.2.5 Leased Lines. Leased lines use the same technology as dial-up (switched) lines, except that the phone company supplies a fixed sequence of connections rather than random switching from one central station to another. There is nothing inherently more secure about a leased line than a switched line; on the contrary, it is easier to tap a leased line at the central switching station because its path is fixed. However, leased lines usually carry high-volume transmissions. The higher the volume of multiplexed data, the more difficult it is for amateur hackers to disentangle the data streams and make sense of them. At the high end of leased line bandwidth

15 · 10 PENETRATING COMPUTER SYSTEMS AND NETWORKS

(e.g., carriers such as T1, T2, etc.), the cost of multiplexing equipment makes interception prohibitively expensive for all but professional or government wiretappers.

Data encryption provides the best defense against wiretapping on leased lines.

15.3.2.6 Long-Distance Transmissions. Dial-up and leased lines carry both short-haul and long-distance transmissions. The latter introduce additional points of vulnerability. Microwave relay towers carry much of the long-distance voice and data communications within a continent. The towers are spaced about 40 kilometers (25 miles) apart; signals spread out noticeably over such distances. Radio receivers at ground level can intercept the signals relayed through a nearby tower, and because microwaves travel in straight lines, rather than following the curvature of the earth, they eventually end up in space, where satellite receivers can collect them. The difficulty for the eavesdropper is that there may be thousands of such signals, including voice and data, at any tower. Sorting out the interesting ones is the challenge. However, given sufficient computing power, such sorting is possible, as is targeting of specific message streams. *Spread-spectrum transmission*, or *frequency hopping*, is an effective countermeasure.

15.3.2.7 Packet-Switching Networks. Packet-switching networks, including historical X.25 carriers such as Telenet, Tymnet, and Datapac, used packet assembler-disassemblers (PADs) to group data into packets addressed from a source to a destination. If data traveled over ordinary phone lines to reach the network, interception could occur anywhere along these segments of the communications link. However, once the data were broken up into packets (whether at the customer side or at the network side), wiretappers had a difficult time making sense of the data stream.

15.3.2.8 Internet Connections. TCP/IP connections are no harder to tap than any others, and they carry an ever-increasing array of data, from e-commerce traffic to television broadcasts and voice communications, the latter using Voice over Internet Protocol (VoIP). (See Chapter 34 in this *Handbook*.) Unless the data stream is encrypted, there are no special impediments to wiretappers. Although the tapping of fiber optic cable requires more specialized equipment than the tapping of copper cables, it is possible.

15.3.2.9 LAN Packet Capture. Local area networks (LANs) are similar to packet-switching networks: both network protocols send information in discrete packages, either over cables or radio waves. Each package has a header containing the address of its sender and of its intended recipient. Packets are transmitted to all nodes on a segment of a LAN. For more information about LANs, see Chapter 25 in this *Handbook*.

Normally a node is restricted to interpreting only those packets that are intended for it alone. However, it is possible to place devices in “promiscuous mode,” overriding this restriction. This can be done with software that surreptitiously converts a device, such as an end user workstation, into a listening device, capturing all packets that reach that node. Of course, network administrators can intentionally create a packet-capturing workstation for legitimate purposes, such as diagnosing network bottlenecks. It is also possible to connect specialized hardware called *LAN monitors* to the network, either with or without permission, for legitimate or illegitimate purposes. Sometimes called *network sniffers*, these devices and programs range from basic freeware to expensive commercial packages that can cost tens of thousands of dollars for a network with

TECHNICAL PENETRATION TECHNIQUES 15 · 11

hundreds of nodes. (The term *sniffer*, although in common use, is a registered trademark, as *Sniffer®*, of Network General Corporation.)

The more sophisticated packet-sniffing programs allow the user to configure profiles for capture; for example, the operator can select packets passing between a host and a system manager's workstation. Such programs allow an observer to view and record everything seen and done on a workstation, including logins or encryption keys sent to a server.

Packet sniffing poses a serious threat to confidentiality of data transmissions through LANs. Most sniffing programs do not announce their presence on the network. Although it may not be apparent to the casual observer that a workstation is performing sniffing, it is possible, as a countermeasure, to scan the network for sniffing devices. Stealthier packet-sniffing technology is constantly improving, and tight physical security may be the best overall deterrence.

LAN users concerned about confidentiality should use LAN protocols that provide end-to-end encryption of the data stream or third-party products to encrypt sensitive files before they are sent through the LAN. Routers that isolate segments of an LAN or WAN (wide area network) can help limit exposure to the threat of sniffers.

15.3.2.10 Optical Fiber. Although optical fibers were once thought to be secure against interception, new developments quickly abolished that hope. An attacker can strip an optical fiber of its outer casing and bend it into a hairpin with a radius of a few millimeters (1/8 inch); from the bend, enough light leaks out to duplicate the data stream. Bryan Betts, writing in *PCWorld*, quoted

Thomas Meier, the CEO of Swiss company Infoguard. . . [who] demonstrated the technique on a fibre carrying a VOIP phone call over Gigabit Ethernet. A section of fibre from inside a junction box was looped into a photodetector called a bend coupler, and the call was recorded and then played back on a laptop.

"People claim optical fibre is harder to tap than copper, but the opposite is true—you don't even have to break the insulation, as you would with copper," Meier said. "You can read through the fibre's cladding with as little as half a dB signal loss."

He claimed that suitable bend couplers can be bought off the shelf—or from eBay—for a few hundred dollars, and connected to the extra fibre that is typically left coiled up in junctions boxes for future splicing needs.

He added that the risk is not imaginary or theoretical—optical taps have been found on police networks in the Netherlands and Germany, and the FBI investigated one discovered on Verizon's network in the United States. Networks used by U.K. and French pharmaceutical companies have also been attacked, probably for industrial espionage, he said.⁵

Luckily, most optical trunk cables carry hundreds or thousands of fibers, making it almost impossible to locate any specific communications channel. (The same is not true of fiber cables used to deliver network connectivity to individual homes and offices.) Equipment for converting optical signals into usable data remains costly, discouraging its use by casual criminal hackers.

15.3.2.11 Wireless Communications. Cable-based communications have the advantage of restricting channel access to at least theoretically visible connections. However, the rapid increase in wireless telecommunications in the last decade of the twentieth century and the first years of the twenty-first has routed increasing

15 · 12 PENETRATING COMPUTER SYSTEMS AND NETWORKS

amounts of information through a broadcast medium in which access—even unauthorized access—may be invisible to users and system administrators.

15.3.2.12 Wireless Phones. Also referred to as cordless phones, conventional wireless phones broadcast their signals, and the traffic they carry can be detected from a distance. Older cordless phones were analog and susceptible to eavesdropping from such basic devices as walkie-talkies and baby monitors. Children sometimes walked around their suburban neighborhoods with a handset from such a phone turned on; once they walked far enough away from their home to lose the signal, any new dial tone belonged to a neighbor, who might be puzzled to discover a call to the Antipodes on the next phone bill. Today's wireless phone models typically use a different set of frequencies, such as 2.4 gigahertz (GHz). These phones generally use frequency-hopping spread spectrum (FHSS) technology to make unauthorized use more difficult, and to impede eavesdropping. FHSS means the signals hop from frequency to frequency across the entire 2.4-GHz spectrum, making tapping their signals harder but by no means impossible. Wireless phones should not be used for confidential voice or data traffic unless encryption is enabled and activated. An added danger lies in hanging up a cordless phone during a conversation, since the cordless phone's base continues to transmit until switched off.

15.3.2.13 Cellular (Mobile) Phones. Early analog cellular (mobile) phone systems had an expectation of privacy equivalent to that of shouting a message through a megaphone from a rooftop. Calls on such phones were easily intercepted using scanners purchased from local electronics stores. Although encryption is possible on the newer digital cell phones that are now widely used, the encryption is not always turned on due to the burden it imposes on the cell company switching equipment. Check with the carrier before assuming that cell calls are encrypted. Also, bear in mind that, although digital cell phone calls are harder to intercept than analog ones, there is a thriving black market in devices that make such interception possible.

As a rule, confidential information should never be conveyed through cellular phones without encrypting the line or the messages first.

Phil Zimmermann, creator of Pretty Good Privacy (PGP, later GPG) in the early 1990s, created a new service in 2012 called Silent Circle to provide encrypted telephony and Internet access to its subscribers.⁶

15.3.2.14 Wireless Networks. The increasingly popular means of networking computers to networks known as WiFi presents many opportunities for interception of communications. In this context, “wireless network” usually means a data network using the 802.11 standard, which comes in a variety of flavors, such as 802.11b, 802.11g (often collectively referred to as WiFi, which stands for “Wireless Fidelity” and is actually a brand name owned by the trade group WiFi Alliance). Typically, this is the sort of local area network created by plugging a wireless access point into an Ethernet network and a WiFi card or adapter into each computer. Most notebook computers now come with a built-in WiFi adapter.

These WLANs, or wireless local area networks, are relatively cheap and easy to create since they do not require network cabling. That helps to explain why more than 200 million WiFi devices were sold in 2006 and more than half of all U.S. companies have been using WLANs to some degree or other since 2002. But cheap WLANs come with hidden costs, namely security. Every WLAN operates, by its very nature, in loose-lips mode. In a sense, the ease of use comes with ease of abuse. They all broadcast

TECHNICAL PENETRATION TECHNIQUES 15 · 13

their traffic into the air, whence it may be overheard by someone with the right set of ears, a legitimate user or a criminal hacker, someone looking for free bandwidth or a war driver. War driving, the practice of driving around town to find wireless access points, is a hobby to some people, and probably not illegal unless done with malicious intent—many notebook computers try to find wireless access points whenever turned on—but bear in mind that laws vary from one country to another. (Those contemplating war driving should check the legal status in their jurisdictions.)

To be a war driver, all that is needed is an old laptop, the right WiFi card, some free software (NetStumbler, e.g.), and an empty Pringles® potato-chip can wired to the WiFi card as an external antenna to boost reception. (The Pringles® can is optional, as is a global positioning system device to mark the location of WiFi access points.) If you drive around with this equipment activated, you will doubtless discover numerous access points in both residential and business districts. If the names of the access points are things like “Linksys” or “netgear,” this is an indication that the owner of the network has not changed the default service set identifier (SSID), which tends to be the brand of the wireless access point, broadcast for the world to see, unless the network owner turns off this feature. The name could also be a person, or place, or company, which helps war drivers figure out whose network they are picking up. (One of the authors detected SCHS near the offices of Sample County Health Services.) A program like NetStumbler will also tell you whether the network is using encryption. There has been a steady rise in the percentage of wireless networks using encryption, but it is far less than 10 percent. According to a survey conducted by AT&T in late 2007, one in six small businesses in America that use wireless technology has taken no precautions against wireless threats, and one-third of small businesses indicated that they were unconcerned about wireless data security.⁷

There is more about wireless network security in Chapter 33 in this *Handbook*, but the point is the relative ease with which networks can be tapped. This means WiFi, whether at home, in the office, or at a hot spot, represents a significant category of data leakage and thus a major avenue for systems penetration.

15.3.2.15 Van Eck Phreaking. This attack is named for Wim Van Eck, a Dutch electronics researcher who in 1985 proved to a number of banks that it was possible to read information from their cathode ray tubes (CRTs) at distances of almost a mile away, using relatively simple technology.⁸ Because many types of electronic equipment emit radio-frequency signals, receivers that capture these signals can be used to reconstruct keystrokes, video displays, and print streams. Using simple, inexpensive wide-band receivers, criminals can detect and use such emissions at distances of tens or hundreds of meters (yards).

Since radio-frequency signals leak easily through single-pane windows, PCs should never be placed in full view of ground-floor windows (and certainly not facing the windows!). Attenuators that tap the window at irregular intervals can be installed to defeat such leakage. A special double-pane window with inert gas between the panes also can lessen the amount of signals leakage.

Other countermeasures include special cladding of hardware, such as computers and printers, to attenuate broadcast signals. This protection often is referred to by the name of the classified government standard for protection of sensitive military systems, TEMPEST. Although TEMPEST was allegedly a classified code word to begin with, it is now sometimes expanded as *Transient ElectroMagnetic Pulse Emission Standard* or *Telecommunications Electronics Material Protected from Emanating Spurious Transmissions*.

15 · 14 PENETRATING COMPUTER SYSTEMS AND NETWORKS

TEMPEST-certified equipment costs many times more than the same equipment without TEMPEST cladding.⁹ A less expensive alternative is to use a special device that emits electromagnetic noise that masks meaningful signals. Yet another approach to protection against this threat is to locate systems within buildings, or rooms within buildings, that have been constructed to TEMPEST standards. There are federal regulations concerning the methods of building sensitive compartmented information facilities (SCIFs), and the testing to obtain a TEMPEST rating is quite stringent. These measures include such things as cladding of all walls and ceilings, cladding of all electrical and network cabling, lead-lined doors, and the absence of any external windows.

15.3.2.16 Trapping Login Information. Criminals can capture identification and authentication codes by inserting Trojan horse programs into the login process on a server host and by using macro facilities to record keystrokes on a client node. More commonly today, however, specifically written malware is used to capture keystrokes and transmit them at intervals back to remote systems.

15.3.2.17 Host-Based Login Trojans. A Trojan horse is a program that looks useful but contains unauthorized, undocumented code for unauthorized functions. The name comes from Greek mythology, in which Odysseus (Ulysses in Latin), weary of the never-ending siege of Troy, sailed his ships out of sight as if he and his warriors were giving up, but left a giant wooden horse at the city gates. Entranced by this magnificent peace offering, the Trojans dragged the great horse into the city. During the Trojans' wild celebrations that night, the soldiers Odysseus had secreted in the belly of the hollow horse let themselves out and opened the gates to their army. The Greeks slaughtered all the inhabitants of the city, and the Trojan War was over.

In February 1994, the Computer Emergency Response Team Coordination Center (CERT-CC) at Carnegie Mellon University in Pittsburgh issued a warning that criminal hackers had inserted Trojan horse login programs in hundreds of UNIX systems on the Internet. The Trojan captured the first 128 bytes of every login and wrote them to a log file that was later read by the criminals. This trick compromised about 10,000 login IDs.

Trojan code might be installed on a computer or terminal used by several people (e.g., on a mainframe terminal in the 1970s or in an Internet café today) so that when someone enters a user ID and password to logon, the system—controlled by the Trojan—displays a message such as “Invalid password, try again,” and the user does so. This time the login is accepted. The victim continues working, unaware that there is anything unusual going on. The Trojan, installed earlier, simulated the normal login procedure, displaying a semblance of the expected screen and dialog. Once the victim entered a password, the Trojan writes the authentication data to a file and then shows a misleading error message. The spoof program then terminates and the regular program is ready for login.

Such a case occurred in April 1993 in a suburb of Hartford, Connecticut. Shoppers noticed a new automated teller machine (ATM) in their mall. At first, the device seemed to work correctly, disbursing a few hundred dollars to bank card users on demand. It quickly changed to a more sinister mode. Users would insert their bank cards and respond as usual to the demand for their personal identification numbers (PIN). At this point, the new ATM would flash a message showing a malfunction and suggesting that the user try an adjacent bank machine. Most people thought nothing of it, but eventually someone realized that the ATM was not posting the usual “Out of Order”

TECHNICAL PENETRATION TECHNIQUES 15 · 15

indicator after these supposed errors. In addition, banks began receiving complaints of a rash of bank card frauds in the immediate area. Investigators discovered that the ATM had no connection to any known bank—that it had been purchased used, along with equipment for manufacturing bank cards. The ATM was a spoof; it was merely collecting the user ID and PIN of every victim for later pickup by the criminals who had installed it without permission in the mall. The criminals were caught after having stolen about \$100,000 over a four-week period using fraudulent bank cards.

Chapter 20 in this *Handbook* includes more details about Trojans and other low-tech attacks.

15.3.2.18 Keystroke Logging. Another threat to identification and authentication codes is the ability to record keystrokes for later playback or editing. Most word processing programs provide macro facilities, so named because of their ability to store and output multiple keystrokes, such as the typing of boilerplate text, with one keystroke. More sophisticated terminate-and-stay-resident (TSR) programs can record sequences of commands, and are sometimes used to demonstrate software or to automate quality assurance tests. This technology can also be used to lay in wait on a workstation and record everything the user does with the mouse and types with the keyboard; such programs are sometimes called *keystroke loggers*. Later, the criminal can harvest the records and pick out the login codes and other valuable information.

There are also hardware implementations of keystroke logging. One is a small device inserted between the keyboard and the computer, capturing what is typed and holding it in nonvolatile memory until it can be retrieved. At the time of this writing in May 2013, such devices were widely available through Internet sales for about US\$40.

By far the most common form of keystroke logging is done by *hooking* the operating system mechanisms for providing keyboard functionality. This is done in software. There are numerous specific technical techniques to accomplish the hooking process depending on the operating system and hardware. Sometimes the software used to accomplish the keystroke logging is stand-alone, but often it is part of a larger piece of software that also provides C2 (command and control) capabilities. It is also quite common for the keystroke logging functionality to be part of a rootkit in order to avoid detection.

It is possible to defeat the attempted reuse and abuse of login credentials captured by any of these methods by switching to one-time passwords generated by microprocessors. One-time passwords are discussed in Chapter 28 of this *Handbook*. However, both key loggers and Trojans can be deployed to gain unauthorized access to data without resorting to the reuse of passwords.

15.3.3 Breaching Access Controls. Criminals and spies use two broad categories of technical attacks to deduce access phone numbers, user IDs, and passwords: brute-force attacks and intelligent guesswork. In addition, there are ways to manipulate people into revealing their access codes; these techniques are discussed in the section on social engineering.

15.3.3.1 Brute-Force Attacks. Brute-force attacks consist of using powerful computers to try all possible codes to locate the correct ones. Brute force is applied to locating modems, network access points, vulnerable Internet servers, user IDs, and passwords.

15 · 16 PENETRATING COMPUTER SYSTEMS AND NETWORKS

15.3.3.2 Demon (War) Dialing. Despite the wholesale shift of data communications to TCP/IP networks and the Internet, modems on dial-up phone lines remain a common, and sometimes forgotten, means of external access to a system. The telephone numbers of any modems connected to hosts or servers or intelligent network peripherals, such as high-end laser printers, are sensitive and should not be posted or broadcast.

Demon dialers are programs that can try every phone number in a numerical range and record whether there is a voice response, a fax line, a modem carrier, or no answer. When phones ring all over an office in numerical order, one at a time, and when there is no one on the line if a phone is picked up, it is undoubtedly the work of someone using a demon dialer. Of course, good demon dialing software accesses the numbers in the target range nonsequentially.

During the heyday of fax machines, some youngsters were reported to have “farmed” entire telephone exchanges during the night, then to have sold the fax numbers for \$1 dollar per number to unscrupulous junk-fax services that sold advertisers access to them.

15.3.3.3 Exhaustive Search. The same approach as demon dialing can find user IDs and passwords after a connection has been made. The attacker uses a program that cycles systematically through all possible user IDs and passwords and records successful attempts. The time required for this attack depends on two factors:

1. The keyspace for the login codes
2. The maximum allowable speed for trying logins

In today’s technical environment, any inexpensive computer can generate login codes far faster than hosts permit login attempts. Processor speed is no longer a rate-limiting factor. Note that this type of attempt to “guess” passwords is different from password cracking, described elsewhere, which operates on captured or stolen copies of encrypted password files.

15.3.3.4 Keyspace. As discussed in Chapter 7 in this *Handbook*, the keyspace for a code is the maximum number of possible strings that meet the rules of the login restrictions. For example, if user passwords consist of exactly six uppercase or lowercase letters or numbers and the passwords are case-sensitive (i.e., uppercase letters are distinguished from lowercase letters), the total number of possible combinations for such passwords is calculated in this way:

- There are 10 digits and 52 upper- or lowercase letters (in the English alphabet) = 62 possible codes for any of six positions.
- If there are no restrictions on repetition, a string of n characters to be taken from a list of r possibilities for each position will generate r^n possible combinations.
- Thus, in our example, there are 62^6 possible sequences of 62 codes taken in groups of six = 56,800,235,584 (more than 56 billion) possible login codes.

If there are restrictions, the keyspace will be reduced accordingly. For example, if the first character of a password of length six must be an uppercase letter instead of being any letter or number, there are only 26 possibilities for that position instead of 62,

TECHNICAL PENETRATION TECHNIQUES 15 · 17

thus reducing the total keyspace to $26 \times 62^5 = 23,819,453,632$ (more than 23 billion) possibilities.

15.3.3.5 Rainbow Tables. Cryptanalysts (including criminal cryptanalysts) can generate all possible one-way hashes for the keyspace of any given password rule; brute-force cracking using the tables can thus be accelerated. There are even Websites distributing such tables freely.¹⁰

15.3.3.6 Login Speed. Generating login codes is not hard. The greatest barrier to brute-force login attacks is interruption in the login whenever the host detects an error. Most operating systems and security monitors allow the administrator to define two types of login delays following errors:

1. A usually brief delay after each failed attempt to enter a correct password
2. A usually long delay after several failed login attempts

Suppose each wrong password entered causes a 1/10th-second delay before the next password can be entered; then for our example involving six repeatable uppercase or lowercase letters or numbers, it would take $5,680,023,558$ seconds = 1,577,784 hours ≈ 180 years to try every possibility.

Suppose, in addition, that after every fifth failed login attempt, the system were to inactivate the user ID or the modem port for three minutes. Such interference would stretch the theoretical time for a brute-force exhaustive attack to around 650 years.

Should the security manager completely inactivate the ID if it is under attack? If the ID is inactivated until the user calls in for help, user IDs become vulnerable to inactivation by malicious hackers. Attackers need merely provide a bad password several times in a row and the unsuspecting *legitimate* user will be locked out of the system until further notice. A widespread attack on multiple user IDs could make the system unavailable to most users. Such a result would be a denial-of-service attack (see Chapter 18 in this *Handbook*).

Should the port be inactivated? If there are only a few ports, shutting them down will make the system unavailable to legitimate users. This drastic response may be inappropriate—indeed, it may satisfy the intentions of criminal hackers. A short delay, perhaps a few minutes, would likely be sufficient to discourage brute-force attacks.

In all of these examples, the illustrations have been based on exhaustive attacks (i.e., trying every possibility). However, if passwords or other codes are chosen randomly, the valid codes will be uniformly distributed throughout the keyspace. On average, then, according to a principle of statistics called the *Central Limit Theorem*, brute-force searches will have to search half the keyspace. For large keyspaces, the difference between a very long time and half of a very long time will be negligible in practice (e.g., 325 years is not significantly different from 650 years if everyone interested will be dead before the code is cracked).

15.3.3.7 Scavenging Random Access Memory. Not all attacks come from outside agents. Criminals with physical access to workstations, malicious software, or authorized users who can use privileged utilities to read main memory, can scavenge memory areas for confidential information such as login IDs and passwords.

On a workstation using a terminal emulator to work with a host, ending a session does not necessarily unload the emulator. Many emulators have a configurable screen

15 · 18 PENETRATING COMPUTER SYSTEMS AND NETWORKS

display buffer, sometimes thousands of lines long. After an authorized user logs off and leaves a terminal, a scavenger can read back many pages of activity, sometimes including confidential information or even login codes. Passwords, however, usually are invisible and therefore not at risk.

If a workstation is part of a client/server system, an application program controlling access may leave residues in random access memory (RAM). A RAM editor, easily available as part of utility packages, can capture and decode such areas as file buffers or input/output (I/O) buffers for communications ports. However, rebooting the workstation after communication is over prevents RAM scavenging by reinitializing memory.

15.3.3.8 Scavenging Cache Files. The same principle can be applied to the various cache and swap files created by the operating system. Cache files are used to keep frequently used data readily available to applications. Swap files store data and code that is moved out of memory onto disk when memory is full. Some operating systems also create hibernation files, writing memory to disk just prior to powering down, and thus enabling a quick resumption of work when the system is powered up. Operating systems may also provide auto-saved recovery files that allow restoration of data after a system error. All of these can be mined for system credentials as well as other valuable data.

15.3.3.9 Scavenging Web History Files. A more recent variation on this scavenging approach is to examine files created by Web browsers. These sometimes contain not only the pages viewed by a user but also the credentials entered to access those pages.

15.3.3.10 Recovering Stored Passwords. Many applications store user passwords locally. Most applications attempt to secure these passwords, unfortunately they are seldom stored with proper encryption methods. It has become common for malicious software to recover these stored passwords and transmit them to remote servers for collection by attackers.

15.3.3.11 Intelligent Guesswork. Users rarely choose random passwords. Much more frequently, passwords are chosen from a subset of all possible strings. Instead of blindly batting at all possible sequences in a keyspace, an attacker can try to reduce the effective keyspace by guessing at more likely selections. Likely selections include canonical passwords, bad passwords, and words selected from a dictionary.

Hardware and software often come from the factory, or out of the box, with user IDs and passwords that are the same for all systems and users.

For example, wireless access points and routers have default user IDs when they ship from the factory. Naturally, these user IDs are set up with the same password. (For example, “admin” on Linksys wireless router devices.) Such systems always include instructions to change the passwords, but, too often, administrators and users neglect to do so. Criminals are familiar with factory presets—most of which are readily discoverable via Google—and exploit them to penetrate systems. The simple routine of changing all canonical passwords prevents hackers from gaining easy access to systems and software.

TECHNICAL PENETRATION TECHNIQUES 15 · 19

15.3.3.12 Stealing. Any element of a computer system has potential value to a thief. Therefore every element including hardware, software, media, files, documentation, and printouts must be safeguarded.

15.3.3.13 Data Scavenging. Criminal hackers have few scruples about using other people's property when they enter computer systems; they have none at all concerning using other people's trash. The term *data scavenging* describes the process that acquires information from throw-away sources. Perhaps the most widely known is Dumpster® diving, sorting through whatever an organization discards. Hard-copy printouts, CD-ROMs, tapes, and other assorted data-bearing media often end up in trash containers where they are easily accessible to Dumpster® divers after hours. In some areas, if one visits an office or industrial park at night, one can see half a dozen people rummaging about, sometimes headfirst in Dumpsters®. Criminal hackers use the information thoughtlessly discarded by naïve office workers as a source of procedures, vocabulary, and proper names that can help them impersonate employees over the phone or even in person. The classic example is a discarded internal phone directory, which can provide a social engineer with valuable data to use when making calls to employees. An employee who hesitates to comply with an attacker's bogus request for information over the phone may well be persuaded if the attacker says something like "I understand your hesitation; if it makes you feel more comfortable you can call me back at extension 2645." If 2645 is a legitimate internal extension, the caller gains considerable credence. Of course, the properly trained employee will hang up and make the call to 2645 rather than take the easy option and say, "I guess that's okay then, here is the information you wanted."

Some printouts contain confidential information that can lead to extortion or system penetration. For example, a thief who steals a list of personally identifiable information about patients with HIV infection could torment the victims and extort money. Every piece of paper, or other media to be discarded, should be evaluated for confidentiality. Unless the information is worthless to everyone, employees should shred paper before disposal or arrange to send paper to a bonded service for destruction. The same applies to CD-ROMs and other media.

15.3.3.14 Discarded Magnetic and Optical Media. Discarded paper poses a threat; discarded magnetic and optical media are a disaster. Many organizations fail to teach employees that the normal commands used to delete files do not remove all trace of them. Either through the use of utility programs or the operating system itself, the original file clusters can be located and any part of the original file that has not yet been overwritten can be regenerated.

Backup tapes, CDs, and DVDs may contain valuable information about the system security structure. For example, in the 1970s and 1980s, system backups on one brand of minicomputer contained the entire directory, complete with every user ID and password *in the clear* on the first tape. Using a simple file copy utility, any user could read these data.

To destroy information on magnetic media, users either must overwrite the medium several times with random data or physically destroy the medium. Degaussers are inadequate unless they meet military specifications, but such units typically cost thousands of dollars.

The problem of readable data is especially troublesome on discarded disk drives or on broken hard disk drives that have been repaired or that are subject to specialized forensic data recovery. Users have received operational, data-laden disk drives as

15 · 20 PENETRATING COMPUTER SYSTEMS AND NETWORKS

replacements for their own broken units. Sometimes the replacement disks have not even been reformatted; they contain entire directories of correspondence, customer databases, and proprietary software. Because it is by definition impossible to overwrite data on a defective disk drive, military security specialists routinely destroy defective hard disks using oxyacetylene torches or purpose-built grinders that reduce hard drives into small chunks.

For more information about secure disposal of magnetic and optical data storage media, see Chapter 57 in this *Handbook*.

15.3.4 Spying. Some techniques used by criminal hackers seem to have been lifted directly from spy novels. For example, laser interferometry can reconstitute vibration patterns from reflected infrared laser beams bounced off windows. Users of such equipment can hear and record conversations in rooms that have external windows that vibrate according to the sounds in the room. For more antispy measures related to physical and facilities security, see Chapters 22 and 23 in this *Handbook*.

Hackers surreptitiously steal people's access codes by watching their fingers as they punch in secret sequences. When pay phones were prevalent, *shoulder surfers* would capture telephone calling-card codes, which they sell to organized crime rings. Codes can be stolen by peering over the shoulders of neighboring callers, or by using binoculars, telescopes, and video cameras to track the buttons pressed by their victims.

Shoulder surfing can occur within installations as well. For example, most users of punch-key locks pay no attention to the visibility of their fingers. Whenever punching in a code, users should guard against observation by unauthorized people. In public places, users should stand up close to the keypad. In fixed installations, facilities managers should cover keypads with opaque sleeves allowing unimpeded access but concealing details of the access codes.

Criminals are adept at surfing both wired and wireless network connections, either by cruising around town with a war-driving setup or hanging out in a target-rich environment such as an airport, train station, coffee shop, or hotel lobby. There is more about wireless hacking in Chapter 33 in this *Handbook*.

One particular type of wired connection that should be used with care is the broadband guest-room connection offered by many hotels. Too often, these are set up without proper security measures, enabling a curious or criminally inclined guest to locate machines belonging to other guests (sometimes by merely clicking the Network Neighborhood icon in Windows Explorer). Employees should be instructed not to plug their company laptops into such connections unless they have a properly configured firewall in place and turned on and are using a virtual private network (VPN) to access corporate systems (see Chapter 32 in this *Handbook*). Even some of the most expensive upscale hotels have been found to suffer from this problem, as illustrated in Exhibit 15.1.

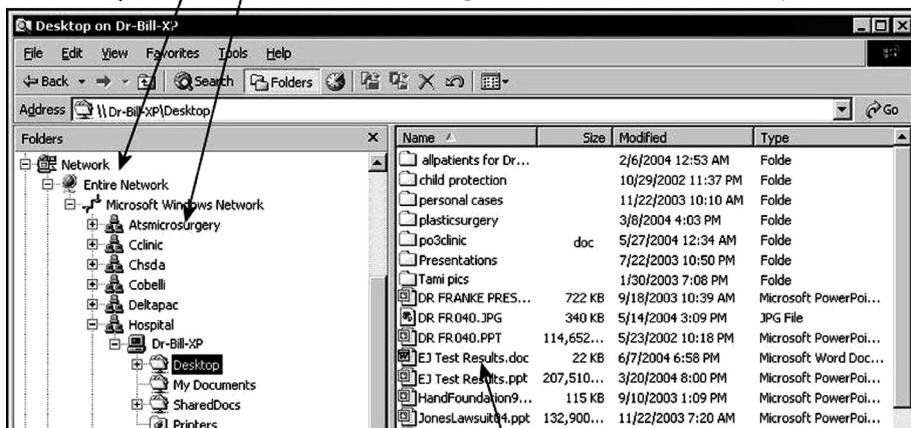
15.3.5 Penetration Testing, Toolkits, and Techniques. Verifying and improving the security of systems by attempting to penetrate them is a well-established practice among security professionals and system administrators. However, although some were practicing this technique earlier, it was not openly discussed prior to 1993. That year, Dan Farmer and Wietse Venema released the pioneering paper entitled "Improving the Security of Your Site by Breaking into It."¹¹ This paper advanced the notion of assessing system security by examining a system through the eyes of a potential intruder. Farmer and Venema showed that scanning for seemingly benign network services can reveal serious weaknesses in any system. Prior to the publication of this important paper, many system administrators were unaware of the extent of vulnerabilities affecting their systems.

TECHNICAL PENETRATION TECHNIQUES 15 · 21

Dr. Bill checks into hotel for medical conference, plugs his laptop into the broadband in his room.

In another room, a different guest also plugs into the hotel broadband. She opens Windows Explorer and double-clicks on the Network icon.

Due to the poor configuration of the hotel network, this guest can see all the other guest machines and is free to explore them.



Dr. Bill has no firewall on his laptop and so all his files are exposed.

This is an actual screenshot taken at a five-star hotel in Boston. To protect privacy, we changed the names of computers, folders, and files. Note that no hacking tools were required to see or to access Dr. Bill's files.

EXHIBIT 15.1 Poorly Configured Hotel Room Internet Connectivity

Farmer and Venema then released a network-testing program called SATAN (Security Analysis Tool for Auditing Networks). Security professionals and system administrators both lauded and were angered by the program. Some system administrators cheered the availability of an all-in-one tool that revealed security holes but did not exploit them. Others questioned the authors' motives in releasing a free and readily available tool that hackers could use to attack networks. While the debate raged on, system administrators and hackers alike began using SATAN to interrogate networks.

15.3.5.1 Common Tools. Since that time, hundreds of penetration toolkits have appeared; they are commonly referred to as *scanners*. Today one can find innumerable freeware tools or invest in one of the commercial tools. Scanners vary in complexity and reliability. However, a majority of the tools employ the same basic functions to test a network: query ports on the target machines and record the response or lack of response.

Used in the proper manner, these tools can be effective in discovering and recording vast amounts of data about a network of computer systems and revealing security holes in the network. Many scanner packages also include packet-sniffing applications, described earlier. Administrators can use this information to reduce the number of systems that can be compromised.

A wide variety of basic network tools may be used in any penetration test. These tools may include mundane programs, such as PING, FINGER, TRACEROUTE, and NSLOOKUP. However, most serious penetration tools make use of an automated vulnerability analysis tool consisting of a series of port queries and a database of known vulnerabilities. Some tools also attempt to exploit identified vulnerabilities in order to eliminate false positives. Once the vulnerabilities have been found, it is

15 · 22 PENETRATING COMPUTER SYSTEMS AND NETWORKS

remarkably easy to obtain “exploits” or programs with which to launch an attack against the susceptible machines. All the tools make use of the basic operations of the TCP suite of protocols. Although these protocols will operate on different port numbers, they all share a common structure of a three-way handshake. All TCP protocols look for a connection attempt (connect), a synchronization and acknowledgment exchange (SYN/ACK), various conditions (FLAGS), and a close port request (FIN). Therefore, the operation of port scanners is quite similar. They attempt to find open, or listening, ports on a machine, ask for a connection, and then log the results to a file to be compared to the internal database. The scanner will display the results of the scans by listing the open ports and the services that appear to be running. At this point the various programs differ. Some attempt an in-depth analysis of the possible security holes associated with the ports and services, along with the appropriate security measures to secure them. The more malicious scanners also include automated scripts for exploitation of the vulnerabilities thus revealed.

Of the freeware tools, Nessus, Netcat, and nmap are probably the best known, although there always seems to be a new flavor of the month among both hackers and system administrators. SATAN is still available, as are its spin-offs, SAINT and SARA. A fair amount of skill is required to use these tools as they are fairly sophisticated, and some of the scans can overload a system and cause it to hang or crash.

15.3.5.2 Common Scans. As previously mentioned, most of the scanners/sniffers available will run through the same basic routines in order to develop a picture of a machine as a whole. The purpose is to determine what is running on the machine and what its role is in the network. The next sections describe most basic scans and their results.

15.3.5.2.1 TCP Connect. This is the most basic form of TCP scanning. This system call is provided by the operating system. If the port is listening, the attempt at a connection will proceed. This scan does not require root or supervisor privileges. However, this scan is easily detectable, as many connection and termination requests will be shown in the host’s system logs.

15.3.5.2.2 TCP SYN. This scan does not open a full connection and is sometimes referred to as a *half-open scan* because a full handshake never completes. A SYN scan starts by sending a SYN packet. Any open ports should respond with an SYN|ACK. However, the scanner sends an RST (reset) instead of an ACK, which terminates the connection. Fewer systems log this type of scan. Ports that are closed will respond to the initial SYN with an RST instead of an ACK, which reveals that the port is closed.

15.3.5.2.3 Stealth Scans. Also referred to as Stealth FIN, Xmas Tree, or Null scans, the stealth scan is used because some firewalls and intrusion detection systems watch for SYNs to restricted ports. The stealth scan attempts to bypass these systems without creating a log of the attempt. The scan is based on the fact that closed ports should respond to a request with an RST and open ports should just drop the packet without logging the attempt.

15.3.5.2.4 UDP Scans. There are many popular User Datagram Protocol (UDP) holes to exploit, such as an rpcbind hole or a Trojan program, such as cDc Back Orifice, which installs itself on a UDP port. The scanner will send a 0-byte UDP packet to each port. If the host returns a “port unreachable” message, that port is considered closed.

TECHNICAL PENETRATION TECHNIQUES 15 · 23

This method can be time consuming because most UNIX hosts limit the rate of Internet Control Message Protocol (ICMP) errors. Some scanners detect the allowable rate on UNIX systems and slow the scan down, so as not to flood the target with messages.

15.3.5.2.5 IP Protocol Scans. This method is used to determine which Internet protocols are supported on a host. Raw IP packets without any protocol header are sent to each specified protocol on the target machine. If an ICMP unreachable message is received, then the protocol is not in use. Otherwise, it assumed to be open. Some hosts (AIX, HP-UX, Digital UNIX) and firewalls may not send protocol unreachable message, so all protocols appear to be open.

15.3.5.2.6 ACK Scan. This advanced method usually is used to map out firewall rule-sets. In particular, it can help determine whether a firewall is stateful or just a simple packet filter that blocks incoming SYN packets. This scan type sends an ACK packet with random-looking acknowledgment/sequence numbers to the ports specified. If an RST comes back, the port is classified as “unfiltered.” If nothing comes back, or if an ICMP unreachable is returned, the port is classified as “filtered.”

15.3.5.2.7 RPC Scan. This method takes all the TCP/UDP ports found open and then floods them with SunRPC program NULL commands in an attempt to determine whether they are RPC ports and, if so, what program and version number they return.

15.3.5.2.8 FTP Bounce. This scan looks like it is an FTP proxy server within the network (or trusted domain). It could eventually connect to an FTP server behind a firewall. Once the FTP server has been found, scanning of ports normally blocked from the outside can be made from the internal FTP server. Of course, reading and writing to directories can be checked from this server as well.

15.3.5.2.9 Ping Sweeps. This scan uses Ping (ICMP echo request) to find hosts that are up. It can also look for subnet-directed broadcast addresses on the network. These are IP addresses that can be reached externally. Ping sweeps often are used to try to “map” the network as a whole.

15.3.5.2.10 Operating System Fingerprinting. As many security holes are dependent on the operating system of the host, this scan attempts to identify which operating system is running, based on a number of suppositions. It uses various techniques to detect subtleties in the underlying operating system (OS) network stack of the computers being scanned. The data gathered are used to create a “fingerprint” that is compared to the scanner’s database of known fingerprints. If an unknown fingerprint is found, attackers can check Websites and newsgroups where information about fingerprints is freely traded, to discover what a particular OS might be. Once the OS has been identified, it is quite easy to find exploits by simply using a search engine on the Web. OS fingerprinting is unnecessary if the OS can be discovered by reading the banners. For example, if one was to telnet to a machine, the response could be:

```
badgny~> telnet abcd.efg.com
Trying 163.143.103.12 ...
Connected to abcd.efg.com
Escape character is '^]'.
HP-UX hpx 8.10.01 A 9000/715 (tttyp2)
login:
```

15 · 24 PENETRATING COMPUTER SYSTEMS AND NETWORKS

The banner, which was included in the default configuration, simply indicates that the OS is HP-UX. A good system administrator will turn off the banners on all services that have them.

15.3.5.2.11 Reverse Ident Scanning. This scan usually is used to see if a Web server on the network is running as root. If the *identd* daemon is running on the target machine, then a TCP Ident request will cause the daemon to return the username that “owns” the process. Therefore, if this request is sent to port 80 (hex), and the return user is root, then that server can be used for an attack on the system. This scan requires a full TCP connection to the port in question before it will return the username.

15.3.5.2.12 Scanning One’s Own Port Configuration. Steve Gibson of Gibson Research Corporation makes available a free port scanner that identifies open ports among the first 1056 TCP ports on any system within less than a minute.¹² Any open port can then be controlled using an appropriate firewall setting. Ideally, all ports on workstations are nonresponsive to pings.

15.3.5.3 Basic Exploits. Using the results of a scanning program, the next logical step for hackers would be to try to exploit the apparent weaknesses in the system. Hackers seek to compromise a machine on the network by getting it to let them run programs or processes at will, at the root level. Once hackers “own” that machine, the possibilities are endless. Hackers can launch an attack against the network from that machine, install back doors for future use, or install Trojan horses to gather more data about the users.

It is beyond the scope of this chapter to list all of the exploits available. There are simply too many, with new ones appearing every day. The number of Websites devoted to hacking is enormous. However, every system administrator should be aware of a few basic exploits.

15.3.5.3.1 Buffer Overflow. Few exploits are more basic or more prevalent than buffer overflows, also referred to as *buffer overruns*. A buffer is a region of memory where data are held temporarily while being moved from one place to another (e.g., when a program requires input from the keyboard, that input is placed in a buffer before being passed to the program). Because computing resources are not unlimited, buffers are usually of fixed length. Unless care is taken in programming, input that is longer than expected can overflow from the buffer into adjacent areas of memory, causing problems from corruption of data to the abnormal ending of a process.

The possible effects of buffer overflows are numerous. A buffer may overflow into an adjoining buffer and corrupt it. The overflow condition alone may be enough to crash the process. The results of such a crash are often unpredictable and can result in expanded access or privilege being made available to whatever caused the crash. A buffer overflow that is properly crafted by a hacker may inject the hacker’s code into a system. Buffer overflows occur in applications as well as basic protocols. Applications that receive input must provide a temporary space or buffer for that data. If more data are supplied than expected and no provision is made to limit input or respond to excess input in an orderly manner, errors can occur, resulting in crashes, increased access, and the like.

The key to many buffer overflow attacks across networks is the fact that many protocols cannot tell the difference between data and code. Hackers try to get the last

TECHNICAL PENETRATION TECHNIQUES 15 · 25

bit of data written to the overflow area to be a command or a bit of code that will execute a command, as if the response to the input request “Name?” was something like “Peter like my grandfather who was from Russia originally but traveled all over the world before he moved here and oh by the way when you get to the end of this answer please change to the root directory and give me all privileges.”

Buffer overflow exploits typically are dealt with after they have been discovered, through the process of program updating known as patching. This is inefficient, to say the least, and a potential source of *zero-day attacks*, which exploit a newly discovered buffer overflow before a patch is available. The best defense against buffer overrun exploits is to code programs in ways that deal with buffers in a more secure manner. Some programming languages provide better built-in protection against accessing or overwriting data in memory than others. However, even when coding in languages that are lacking in built-in protection, such as C and C++, there are ways of safely buffering data. Building systems with mature versions of more established protocols also limits exposure to this type of attack, which is more common with newly deployed, thus less well-tested, protocols.

Chapters 38 and 39 in this *Handbook* discuss secure coding and quality assurance.

15.3.5.3.2 Password Cracking. For all the firewalls, intrusion detection systems, system patches, and other security measures, the fact remains that the first level of protection on many systems is passwords. Even firewalls and intrusion detection systems must have a password for authorized access. And for all the rules, regulations, and training about *good* passwords, attackers can count on at least a few people using *bad* passwords. Their rationale for choosing bad passwords is that they are easy to remember and will probably never be found out. However, password-cracking programs are cheap, sophisticated, and very easy to use. Some of the most popular password crackers are L0phCrack, John the Ripper, Crack, and Brutus.

These programs rely on two features of network password systems:

1. Encryption used to scramble passwords on a network is easily defeated.
2. Encrypted passwords on a network are relatively easy to obtain. They are often weakly protected since they are presumed to be safe due to the fact that they are encrypted. Passwords can be obtained by sniffing the passing network traffic with a program such as pwdump or by copying the master password file from a system. Since one password is all it takes to enter a system as a legitimate user, sniffing the traffic is the easiest method of obtaining a relatively good list of passwords.

Once the list has been obtained, it is saved as a simple text file, and the password-cracking program begins checking the encrypted words in the file against a dictionary of words that have previously been encrypted with the same algorithm. Whenever a match is found between an encrypted string in the file and a word in the encrypted dictionary, the cracking program displays and records the plaintext of the encrypted dictionary word. Thus the password is revealed.

In addition to checking ordinary dictionary words, some password crackers check for both uppercase and lowercase letters, numbers before and after a word, and numbers used in lieu of vowels within a word. The speed at which these programs operate, even on a basic desktop or laptop computer, is impressive, and it is entirely

15 · 26 PENETRATING COMPUTER SYSTEMS AND NETWORKS

possible to obtain cracked passwords within seconds. Indeed, a useful security awareness exercise is to demonstrate such a program to employees: The first passwords to be cracked will be the weakest ones, and that can serve as a warning to users who choose such words.

A good security officer will ensure that passwords on a network are checked regularly with a password cracker or by implementing one of the many strong-password enforcers. Password enforcers augment the password program by comparing the passwords chosen by the user to the rules set by the enforcer.

A simple online tool from Steve Gibson's GRC allows one to calculate the keyspace and estimate cracking time for sample passwords.¹³ Users should not enter their actual password but only one using the same rules; for example, if one's actual password were *Dk3*(4n\$p2*, one could enter *Eh5&%9g#t8* to arrive at exactly the same analysis. In the example presented here, the calculations result in a keyspace of 6.05×10^{19} with a cracking time of 1 week if a massively parallel array of processors supported 10^{15} guesses per second. The password *evanescent porridge* would take the same processors 14.32 billion centuries to include by brute-force cracking of the 4.50×10^{33} keyspace.

15.3.5.3.3 Rootkits. Rootkits are one of the many tools available to hackers to disguise the fact that a machine has been “rooted.” A rootkit is not used to crack into a system but rather to ensure that a cracked system remains available to the intruder. Rootkits are comprised of a suite of utilities that are installed on the victim machine. The utilities start by modifying the most basic and commonly used programs so that suspicious activity is cloaked. For example, a rootkit often changes simple commands such as “ls” (list files). A modified “ls” from a rootkit will not display files or directories that the intruder wants to keep hidden.

Rootkits are extremely difficult to discover since the commands and programs appear to work as before. Often a rootkit is found because something did not “feel right” to the system administrator. Since rootkits vary greatly in the programs they change, one cannot tell which programs have been changed and which have not. Without a cryptographically secure signature of every system binary, an administrator cannot be certain to have found the entire rootkit.

Some of the common utilities included in a rootkit are:

- Trojan horse utilities
- Back doors that allow the hacker to enter the system at will
- Log-wiping utilities that erase the attacker's access record from system log files
- Packet sniffers that capture network traffic for the attacker

15.3.5.3.4 Trojan Code. As described earlier in the context of compromised login procedures, Trojan code is something other than it appears to be. In this case, the Trojans are the changed programs in a rootkit that allow an intruder's tracks to be hidden or allow the program to gather more information as it sits silently in the background. Local programs that are Trojaned often include “chfn,” “chsh,” “login,” and “passwd.” In each case, if the rootkit password is entered in the appropriate place, a root shell is spawned.

15.3.5.3.5 Back Doors. Back door utilities often are tied to programs that have been Trojaned. They are used to gain entry to a system when other methods fail. Even if

TECHNICAL PENETRATION TECHNIQUES 15 · 27

a system administrator has discovered an intrusion and has changed all the usernames and passwords, there is a good chance that he or she does not know that the back doors exist. To use the back door, the hacker needs only to know the correct port to connect to the compromised machine and to enter a password or command where one is not usually entered.

For example, `inetd`, the network super daemon, is often Trojaned. The daemon will listen on an unusual port (`rfe`, port 5002 by default in Rootkit IV for Linux). If the correct password is given after connection, a root shell is spawned and bound to the port.

The function `rshd` can be similarly Trojaned so that a root shell is spawned when the rootkit password is given as the username and thus `rsh [hostname] -l [rootkit password]` will obtain access to the compromised machine.

15.3.6 Penetration via Websites. A vast new network territory opened up in the final decade of the twentieth century, partially devoted to commerce and largely driven by attempts to make money from a technology that originally had been developed for military and academic purposes. The relentless growth of this network surpassed 145 million registered domains at the time of writing in May 2013 as reported by the Whois Source (www.whois.sc/internet-statistics), one of the most reliable sources of statistics about the Internet.

Worldwide Internet penetration, measured as the number of Internet users as a percentage of total population, was over one third by June 2012, with North America, surpassing 75 percent.¹⁴ Asia had almost 4 billion users at that time. Not surprisingly, with so many machines in one network and over 2.4 billion users (twice what was reported in the 2009 Fifth Edition of this *Handbook*), this new territory is the primary playground for hackers, from the merely curious to the seriously criminal. The Web presents a target-rich environment for people seeking unauthorized access to other people's information systems. There are several reasons for this; chief among them is the fact that many organizations, both commercial and governmental—including the military—have external, public Websites that are connected, in some way, to internal, private networks. This connection provides a system penetration path that can be exploited in many different ways, as outlined in this section.

For more detailed analysis of Website security, see Chapter 21 in this *Handbook*.

15.3.6.1 Web System Architecture. Standard practice when placing a commercial Website on the Internet is to screen it from hostile activity, typically using a router with access-control lists (ACLs) or a firewall, or both. However, unless the Web site is of the basic, “brochure-ware” kind, which simply exists to provide information on a read-only basis, the site has to allow for user input of data. Input is required for something as simple as a guest book entry or an information request form; more complex applications such as online shopping have more complex input requirements.

A typical method of processing input is the Common Gateway Interface (CGI). This is a standard way for a Web server to pass user input to an application program and to receive a response, which can then be forwarded to the user. For example, when a user fills out a form on a Web page and submits it, the Web server typically passes the form information to a small application. This application processes the data and may send back a confirmation message. This method is named *CGI*, and it forms part of the Hypertext Transfer Protocol (HTTP). Because it is a consistent method, applications written to employ it can be used regardless of the operating system of the server on which it is deployed. Further adding to the popularity of CGI is the fact that it works

15 · 28 PENETRATING COMPUTER SYSTEMS AND NETWORKS

with a number of different languages including C, C++, Java, and PERL. The term “CGI” is used generically for Web server code written in any of these languages.

Any system that receives input must provide a path within the system through which that input can flow. This is referred to as an “allowed path.” The necessity of allowed paths, combined with the ability to exploit them for penetration attacks, led the system security expert David Brussin (author of Chapter 26 in this *Handbook*) to coin the term “allowed path vulnerabilities” for this class of vulnerability. The two leading categories of exploits that employ allowed paths are input validation exploits, which are akin to buffer overflow exploits and often employed to abuse CGIs, and file system exploits, which abuse the Web server operating system and services running on the server.

Of course, there are other ways to abuse Websites as well. Denial-of-service attacks can be used to inhibit legitimate access and thus compromise availability. (Chapter 18 in this *Handbook* covers DoS attacks.) Not every attack is aimed at further penetration of systems; the Web pages themselves may be the target of attack, as in a defacement, an unauthorized change to Web pages. Defacement often is committed to embarrass the Website owner, publicize a protest message, or enhance the reputation of the criminal hacker who is performing the defacement. But in terms of penetration, the primary goal of attacks on Websites is to compromise internal networks that might be connected to the Web server. Exploitation of allowed path vulnerabilities is probably the most common form of such attacks.

15.3.6.2 Input Validation Exploits. Whenever an allowed path is created to accommodate user input, the possibility exists that it will be abused. Such abuse can lead to unauthorized access to the system. This section describes a range of penetration methods using this approach, all of which somehow employ invalid input, or input that is:

Not expected by the receiving application on the server.

Not “allowed” according to the rules by which the receiving application is operating.

15.3.6.3 Unexpected Input Attacks. How is it possible to submit unexpected input to a server? The answer lies in the architecture of the Internet and the paradoxical nature of the client system that is accessing the server. The typical Web client is a client in name only. It is often a powerful machine, capable of being a server in its own right, and very difficult for any other server to control, due to the inherent peer-to-peer nature of the huge network that is the Internet. All nodes of the Internet are considered hosts. And, of course, many of these hosts are outside the physical control of the organizations hosting those machines that are acting as servers. This fact has serious implications for security.

Unless the server can install tightly controlled application code on the client, and restrict user input to that code, the server must rely on the coding most commonly used to implement Web client-server interaction, Hypertext Markup Language (HTML) and Hypertext Transfer Protocol Daemon (HTTPD). Both are complex and relatively immature. For example, they do not automatically identify the source of the input. Consider an HTML form on a Web page, designed to be presented to a visitor to the Website who fills in the fields and then clicks a button to submit the form. There is nothing on the client side to control the user’s input. So instead of entering a first name in the First Name field, a user might enter a long string of random characters. Unless the application processing this data field performs extensive input validation, the effects of

TECHNICAL PENETRATION TECHNIQUES 15 · 29

such action can be unpredictable, even more so if the user includes control characters. Similarly, if the Web server itself is not designed to validate page requests, a user might cause problems by submitting a bogus Universal Resource Locator (URL).

The problem is even more severe than this. Unless the Web server application is specifically written to defeat the following abuse, it can be used to cause all sorts of problems that potentially lead to successful penetration. Suppose that, instead of simply filling out a form, the user creates a local copy of the page containing the form, then alters the source code of the form, saves the file, and submits it to the Website instead of the original page. This does not violate the basic protocols of the Web but clearly provides considerable penetration potential. Many Websites are still vulnerable to attacks of this type.

15.3.6.4 Overflow Attacks. As described earlier in the context of hidden form fields, it is possible to gain access to Web servers by supplying more input than expected. Such an attack is possible with any field on a user-submitted form, not just a hidden field. The defense is to build extensive error checking into the application that processes the form.

Overflow attacks, which were described in general terms earlier, also can be directed at applications or services running on the Web server. For example, in June 2001, CERT announced a remotely exploitable buffer overflow in one of the Internet Server Application Programming Interface (ISAPI) extensions installed with most versions of Microsoft Internet Information Server 4.0 and 5.0, specifically the Internet/Indexing Service Application Programming Interface extension, IDQ.DLL. An intruder exploiting this vulnerability may be able to execute arbitrary code in the local system security context, giving the attacker complete control of the victim's system.

15.3.6.5 File-System Exploits. Another category of attack against Websites exploits problems with the file system of the Web server itself. Ever since Web servers started appearing on the Internet, there has been a constant procession of vulnerability announcements arising from file system issues. The main ones are presented here, but the possibility of others appearing is high due to a lack of what David Brussin has called *vulnerability class analysis*. A vulnerability class is a type of problem, such as buffer overflow or file system access control. Developers of Web servers, and many other applications, are often averse to, or resource-constrained from, the elimination of vulnerabilities as a class, being focused instead on the hole-by-hole fixing of specific instances of the vulnerability as they arise. This phenomenon is largely a result of the rapid pace at which the Web has been developed and deployed, driven by powerful commercial forces.

15.3.6.6 Dot Dot, Slash, and Other Characters. Persons responsible for the security of information systems that employ Web servers must be alert for new vulnerabilities. Web server software has proven particularly susceptible to certain categories of vulnerability that tend to recur in new versions. Whenever these vulnerabilities are discovered, attackers quickly exploit them. Typically, software vendors issue patches to solve the problem, but systems remain susceptible until patched, and attackers use automated tools to scan the Internet for servers that are still susceptible. For example, in April 2001, a flaw was discovered in versions of Microsoft Internet Information Server (IIS) in use at that time. This flaw made it possible for remote users to list directory contents, view files, delete files, and execute arbitrary commands.

15 · 30 PENETRATING COMPUTER SYSTEMS AND NETWORKS

In other words, if a Web server running IIS was connected to the Internet, anybody using the Internet could potentially copy, move, or delete files on the Web server. With this level of access, it was possible to use the Web server to gain access to connected networks unless strong internetwork access controls were in place. In alerts that were issued to warn users of this software, exploitation of this vulnerability was described as “trivial.” In fact, a large number of sites were penetrated because of it, and many suffered defacement (i.e., unauthorized changes to the appearance of their Websites). In other cases, attackers downloaded sensitive customer data and uploaded and installed back door software.

Particularly worrying about this vulnerability was the fact that it was essentially a recurrence of the so-called dot dot directory traversal attack, which was possible on a lot of early Web servers. These servers would, upon request, read “..” directories in URLs, which are unpublished parent directories of the published directory. Thus, attackers were able to back out to the Web root directory and then to other parts of the server’s directory structure. This technique basically allowed attackers to navigate the file system at will. Many Web servers, including IIS, began to incorporate security measures to prevent the “dot dot” attack, denying all queries to URLs that contain too many leading slashes or “..” characters.

The vulnerability published in April 2001 involved bypassing these restrictions by simply substituting a Unicode translation of a “/” or “.:.” Attackers found that, by appending the “..” and a Unicode slash or backslash after a virtual directory with execute permissions, it was possible to execute arbitrary commands. Attackers could execute any command via a specially crafted HTTP query. The frequency with which “old” vulnerabilities reappear in new software should serve as a warning to information security professionals not to assume that “new” is the same as “improved.” Indeed, all software needs to be treated with a healthy degree of skepticism and a fair amount of heavy testing prior to deployment.

15.3.6.7 Metacharacters. Dots and slashes used in field system references are closely related to metacharacters, which also can be used to attack Web systems. A metacharacter is a special character in a program or data field that provides information about other characters, for example, how to process the characters that follow the metacharacter. Users of DOS or UNIX are probably familiar with the wildcard character, a metacharacter that can represent either any one character or any string of characters. If used inappropriately, for example, in user-supplied data, metacharacters can cause errors that result in unintended consequences, including privileged access.

15.3.6.8 Server-Side Includes. Server-side includes (SSIs) are special commands in HTML that the Web server executes as it parses an HTML file. SSIs were developed originally to make it easier to include a common file, called an *include file*, inside many different files; examples include files containing a logo or text files consisting of the page, date, author, and so on. This capability was expanded to enable server information, such as the date and time, to be included automatically in a file. Eventually, several different types of include commands were provided on Web servers: *config*, *include*, *echo*, *fsize*, *flastmod*, *exec*. The last of these, *exec*, is quite powerful, but it is also a security risk, as it gives to the user of the Web client permission to execute code. A number of attacks are possible when *exec* is permitted within an inadequately protected directory.

Analogous to SSIs are ASPs (Active Server Pages) and JavaServer Pages, as well as Hypertext Processors (PHP). All are technologies that facilitate dynamic page building

TECHNICAL PENETRATION TECHNIQUES 15 · 31

and allow execution of codelike instructions within an HTML page. All should be employed with special attention to security implications and strict adherence to secure Web application programming methods. Security professionals will note that many of the weaknesses in these newer technologies are simply old exploits reborn. Anything that offers “greater interactive functionality” may increase the likelihood of new vulnerabilities simply because, in the absence of rigorous software quality assurance, changes are so often associated with new bugs.

15.3.7 Role of Malware and Botnets. Several new twists on system penetration via Websites and Internet email have emerged in recent years, starting with the evolution of computer viruses and worms. Some worms and viruses are designed to install Trojan code. The term “malware,” derived from “malicious” and “software,” is now widely used to refer to the entire category of software coded with malicious intent, including viruses, worms, and Trojans. In fact, viruses and worms, which are dealt with in more detail in Chapter 16 in this *Handbook*, are themselves a form of system penetration; after all, the creators of virus and worm code are getting their code onto systems that they are not authorized to use, so one may consider those machines to have been penetrated. Some criminal hackers have combined different elements of malware to penetrate computers used to surf the Web. Those compromised machines are then used, in turn, to spread malware and compromise additional machines. The goal may be gathering of user names and passwords (helpful for yet more system penetrations) or financial data used to perpetrate fraud and identity theft.

The two main components of this penetration strategy are drive-by downloads and botnets. A drive-by download attempts to compromise machines used to visit a malicious Website, taking advantage of either user gullibility or vulnerabilities in their Web browsers to install, without explicit permission, unsolicited code, typically a Trojan of some sort. A botnet is a collection of bots, host computers that can be controlled remotely (i.e., robotically) through Trojan code installed on those machines, either via a drive-by download or other means, such as a virus or worm. Here is how researchers at Google described the phenomenon in a landmark 2007 report:

[C]omputer users have become the target of an underground economy that infects hosts with malware or adware for financial gain. Unfortunately, even a single visit to an infected Website enables the attacker to detect vulnerabilities in the user’s applications and force the download of a multitude of malware binaries. Frequently, this malware allows the adversary to gain full control of the compromised systems leading to the ex-filtration of sensitive information or installation of utilities that facilitate remote control of the host.¹⁵

What makes the Google report a landmark is not the existence of drive-by exploits, which have been on the rise for several years, but their prevalence. Because Google maintains a massive repository of Web pages in order to operate its search engine, it is in a fairly unique position when it comes to analyzing the content of the Web as a whole. The solidly researched finding that at least 1 in 10 of all Web pages contained some form of malware should be a wakeup call to IT departments everywhere. Are your users surfing to Facebook pages? Are they aware that something as seemingly harmless as visiting their favorite singer’s page on the Web might cause malicious code to be downloaded onto their computer from a server in China?

That sort of attack started to become commonplace in 2007, as documented by Roger Thompson of Exploit Prevention Labs, with Alicia Keys being one of a number of artists targeted by criminal hackers. One exploit used in these attacks installed a

15 · 32 PENETRATING COMPUTER SYSTEMS AND NETWORKS

proxy network bot, known as a *flux bot*. The goal of a flux bot is to obscure the location of phishing sites by using constantly changing proxy servers, thus making it harder for banks and other institutions targeted by the phishing scam to shut it down. (There is more about phishing scams in Chapter 20 in this *Handbook*.)

According to the Google report, there are four main methods by which Web pages are turned into malware infection vectors: advertising, third-party widgets, user-contributed content, and Web server security. This implies that company Websites not only need to be firmly secured, but that all advertising and user-supplied content should be validated, as should any widgets that are employed or distributed by the site.

15.3.8 Sophisticated Attackers. In recent years several groups of attackers have been very successful using a very different approach to compromise security systems than has been used historically. These groups have been termed “Advanced Persistent Threats” by the United States Air Force in 2006. While many of these groups are nation state–sponsored these same techniques are being employed by groups with purely financial gain as well.

It is important to understand up front that these attacks, regardless of specific source group, are targeted. In the case of the nation state–sponsored groups, the goal is information theft. This information theft often takes the form of email, but also can range much deeper to research and development information or other intellectual property. The financially motivated groups focus on financial institutions for the purposes of financial gain.

As of this writing, there are many known groups active in the world. Each of these groups uses a little different specific components in their overall process, but the overall process is the same and can be broken down into some distinct stages.

The first stage is reconnaissance, planning, and preparation. This is where they select specific individuals within the target organization. They also plan out their attack based on the results of the reconnaissance. Finally they set up specific command and control servers to be used against the target and build the malicious software to be used as part of the attack.

The second stage is obtaining foothold and persistence. In this stage their goal is to get onto the target network. This is the point they deploy their attack (most often through spear-phishes) and gain remote access to the target network. Since gaining a foothold is one of the more challenging parts, they also seek to ensure they can maintain that foothold on the network.

The third stage is mapping the internal network and locating the desired data. At this point, they seek to understand the target environment and find their ultimate target of either specific data or systems on which they can perpetrate their financial impact.

The fourth and final stage is exfiltration or exploitation. This is the point at which they transmit the collected information outside of the target environment (typically nation-state actor goal) or execute their financial attack (in the case of the cybercrime actors).

15.3.8.1 Stage One—Reconnaissance, Planning, and Preparation.

The attackers often begin with selecting a handful of individuals at the target organization. They employ open source intelligence sources like the company Website, Facebook, and LinkedIn to gather information and identify recipients at the target organization. Other sources, such as press releases from companies and news articles about those companies, provide a rich resource for attackers to zero in on individuals close to the ultimate target of the attackers.

TECHNICAL PENETRATION TECHNIQUES 15 · 33

Once individuals have been identified an appropriate delivery can be created. The most common delivery means is through spear-phishes. A spear-phish is a targeted phish (as compared to a normal phish that is sent to thousands or millions of recipients). These spear-phishes are often simple and straightforward but can be quite sophisticated.

While spear-phishing is the most common means of delivery attackers will also compromise Websites commonly used by their targets as a means of delivery. This method of delivery is referred to as watering hole attacks. The term *watering hole* refers to the hunting method whereby a hunter waits for their prey at a watering hole, knowing that eventually the target will come for a drink. Watering hole attacks work by modifying a compromised Website's code so that when users view the Website an exploit is delivered to the victim browser. This exploit is in turn used to install malicious code on the system.

Least commonly, the sophisticated attackers have been found to use malicious code placed on USB keys. These are either left where target users will find them or even mailed directly to the targets. Exploit code is used on the USB device so that when users access the device the exploit code will install malicious code onto the user system.

Also as part of stage one the attackers will configure C2 (command and control) servers to be used over the course of the attack. These C2 servers are often compromised hosts at legitimate organizations but can also be hosting servers obtained legitimately. In particular, attackers prefer servers at places such as large universities, as these are unlikely to be blocked by common prevention mechanisms like Web proxy filters commonly employed by organizations today.

The final major aspect of the preparation is the creation of the backdoors to be employed against the target organization. The attackers generally have several variants of their backdoor with slightly different characteristics. Some groups will even use commonly available RATs (Remote Access Trojans) such as Poison Ivy and configure them to use the appropriate C2s. Many of the attackers have tools that will weaponize legitimate PDF or office documents. Weaponization is the process of turning a legitimate file into a malicious one. If, for instance, the individuals share a particular industry in common, the attackers might find a relevant industry conference in the near future and download an actual agenda PDF file from the conference Website. They might then weaponize that PDF or simply use information from the file to add to their spear-phish to add realism and improve their chances for the targeted users to fall for the spear-phish.

15.3.8.2 Stage Two—Foothold and Persistence. The end result of stage one is an RAT being installed on one or more target systems. This Trojan provides covert backdoor access to the victim systems. Now that the attacker has a foothold on the target network, they want to maintain that access. Attackers want access from more than a single system, since a single host might not be running when they want access. The attackers also understand that eventually the RAT will be found or removed. The attackers also want to increase their access to ensure they can get to whatever their final objective is.

To accomplish all of this, the attackers will usually follow a consistent process. They start with dumping the password hashes from the local host cache. By default, Microsoft Windows (the predominant operating system in use by users at most organizations) keeps a cache of the last 10 user IDs and passwords used to log into a computer. This is done so that if the computer is not connected to the corporate environment and thus can't access the domain controllers, the user will still have the ability to log into the system locally. Tools like pwdump, Windows credential editor, and fgdump

15 · 34 PENETRATING COMPUTER SYSTEMS AND NETWORKS

will display the user IDs and passwords in the local system cache. These can then be cracked using rainbow tables or cracking software back at the attacker's location.

The second step is to look at the host network connections to determine what computers the host is communicating with. The logic behind this is simple—any credentials in the local computer will very likely work on other hosts the system is interacting with. Remember that the purpose of user IDs is to provide authentication to other systems for the purpose of using them. Using the cracked password for the locally cached user IDs, the attacker connects to the subsequent systems and proceeds to extract the user IDs and password hashes from that host. The process of moving from host to host like this is termed *lateral movement*. It is especially important to note that the attackers are not using malware to accomplish this access. They are using legitimate user IDs and passwords. Ultimately, the attacker will continue this process until they find a system with an administrative account that will give them access across the entire environment. It is rare for them to have to traverse more than six to ten systems in order to find one with a network or domain administrator account cached. Once they have an administrative account they are free to use it to move throughout the network.

In addition to obtaining local user IDs and passwords from system caches, the attackers also map the target environment at this stage. Mapping the network is done through simple querying of Active Directory. The simplest way for the attackers to gain information about the environment is with a series of simple commands:

```
net group "domain computers" /domain
net group "domain users" /domain
net group "domain controllers" /domain
net group "domain admins" /domain
```

If you aren't familiar with those commands, you should try them. The only requirement is to do them from a computer that is a member of the local domain. No administrative credentials are necessary. The result is a list of all computers, users, domain controllers, and domain admin accounts in the environment. This information is returned immediately and is far more useful for attackers' purposes than running a tool like nmap. A further benefit is that the use of these commands is very, very difficult to detect as compared to nmap or similar tools that are very noisy. When the attackers want to gain more thorough information about an environment they will typically use dsquery and dsget. These are two command line tools provided by Microsoft to query active directory. Here are the contents of a batch script that was recovered from a recent incursion showing the use of these tools.

```
dsquery user -limit 0 | dsget user -samid -display -title
    -email -dept -office -company -c >1
dsquery user -limit 0 | dsget user -samid -pwdneverexpires
    -acctexpires -loscr -profile -hmdir -hmdrv -c >2
dsquery user -limit 0 | dsget user -c > 3
dsquery group -limit 0 | dsget group -c >g
dsquery subnet -limit 0 | dsget subnet -c >sn
dsquery site -limit 0 | dsget site -c >s
dsquery computer -limit 0 | dsget computer -c>c
dsquery ou -limit 0 | dsget ou -c>o
```

Attackers understand that users leave organizations and change roles. To ensure they have plenty of user IDs at their disposal it is a high priority for them to dump the entire domain credentials. This usually occurs within minutes of them achieving an administrative account with sufficient credentials. Their tool of choice for this is the

TECHNICAL PENETRATION TECHNIQUES 15 · 35

well-known pwdump. They simply run pwdump against one of the domain controllers using a recovered domain administrator account. Once the hashes are obtained, they can crack them at their leisure back at their location.

Another priority of the attackers is to install additional mechanisms for gaining access to the environment. To accomplish this they will use a variety of methods. If a company employs a single-factor VPN, then the attackers have unlimited access to the environment using the stolen user credentials. They will also install several additional backdoors. These are distributed throughout the environment to make it difficult for security personnel to find them all. The attackers understand that malware is the easiest thing to spot in an environment, however, and so they will only install copies on a handful of the overall hosts they compromise. As a further precaution to detection, they will usually configure the extra backdoors to sleep for extended periods of time. This tactic results in it being very challenging for security staff to find and eliminate all copies of their backdoors. Finally, in the last couple years, they have been observed also deploying Web shells in corporate DMZs as a further level of backup access to an environment. Web shells were very popular with attackers in the late 1990s and early 2000s but have been rare for the last several years, and thus are not commonly looked for. The Web shells being employed by the attackers are very small and do nothing until remotely accessed, making them very difficult to uncover as well.

As if the use of so many techniques in order to maintain access to a company was not enough, most groups will also monitor their installed backdoors in real time. If they observe an organization actively removing them they will quickly jump onto compromised hosts and install entirely different backdoor versions with very different characteristics so as to maintain access and remain undetected. It is important to realize that the majority of the backdoors are proprietary, not commodity, and very rarely found with standard antivirus. The net result for the attackers is a long-time presence on an environment, during which they can move around and access whatever they desire.

15.3.8.3 Stage Three—Information Theft. The final stage and aspect of the advanced attackers is where they find and exfiltrate the information they see or execute the financial transfers in the case of the cybercrime-motivated actors. All of the efforts of the attackers are for the ultimate purpose of stealing something.

After establishing long-term persistence in an environment, the attackers begin moving from host to host looking for the information they seek. Mostly this is done by command line, but some groups will also tunnel RDP protocol through their C2 infrastructure. If an organization has some other form of remote control, such as VNC, the attackers will gladly utilize that as well with the legitimate credentials they compromised. Through the simple expedient of the Windows ‘dir’ command, they look in the documents of each host, seeking the data they desire. This is facilitated by the deep organizational knowledge they acquired in earlier stages.

Often the majority of their theft is simply e-mail. When you consider the amount of business activity details that are conducted with e-mail, you realize an attacker can gain incredible competitive advantage from information contained in e-mail.

When the attackers find the data they seek, they will compress and encrypt it, then transmit it back to the C2 infrastructure. This process is known as *exfiltration*. Since sophisticated attackers understand organizations have logging measures in place, the transfer is done to intermediary systems, rather than directly to their location. Subsequently they can move it from the intermediary systems to the final destination without risk of the compromised organization being able to trace it.

15 · 36 PENETRATING COMPUTER SYSTEMS AND NETWORKS

Most often RAR is used to package the data for exfiltration. RAR is particularly useful for attackers because of its unique characteristic among compression formats: Data in an RAR file can be recovered even if the RAR is incomplete. If exfiltration is interrupted midtransmission, the attackers will still be able to extract whatever portion of the data made it out. Encryption is used to make it more difficult for a compromised company to figure out what data was stolen. Most groups also prefer to use a media file extension for their RAR files as a further means of obfuscation. The renamed files may have a .mpg or .avi extension but they are simply normal encrypted RAR files.

For transmission of the data the attackers typically use Secure Sockets Layer (SSL). This tool, combined with the compressed and encrypted contents, makes it very difficult for companies to detect that anything malicious is occurring. The sophisticated attackers are so successful overall because the majority of their activities are done by hiding in the large volume and noise of legitimate activity. Consequently, they've all been observed using services like dropbox, box.net, and other cloud providers for exfiltration as well, given the significant adoption rate of these types of services legitimately by our users. Detecting malicious activity to a cloud provider from legitimate activity is very difficult indeed.

15.4 POLITICAL AND LEGAL ISSUES. Thanks to the World Wide Web, the Internet has become a self-documenting phenomenon. One can use the Internet to find out everything one wants to know about the Internet, including how to penetrate information systems that employ Internet technology. However, the penetration information available on the Internet is not restricted to Internet systems, and the Internet is not the only source of penetration information. Furthermore, the very availability of penetration information is fraught with political and legal issues. These are discussed briefly in this final section of the chapter.

15.4.1 Exchange of System Penetration Information. The sharing of system penetration information—that is, information that could facilitate illegal penetration of an information system—is the subject of a long, heated, and ongoing debate. This debate encompasses both practical and ethical aspects of the issue. Although complete coverage is not possible within the confines of this chapter, we do review the question of full disclosure, along with some of the sources for penetration information.

15.4.2 Full Disclosure. How should we handle known vulnerabilities and potentially damaging computer viruses? Should we publish full details, conceal some details, or suppress any publication that would allow exploitation until patches or updates are available from manufacturers? Is there a case for handling some viruses and exploits differently from others?

Over the last two decades, formal venues for full disclosure of system or network vulnerabilities and exploits have evolved (e.g., BugTraq). Support for full disclosure of such details, down to the source code or script level, from professional, honest security experts (the Good Guys) is based on subsets of several key beliefs:

- The Bad Guys know about the vulnerabilities anyway.
- If they do not know about it already, they will soon with or without the posted details.
- Knowing the details helps the Good Guys more than the Bad Guys.

POLITICAL AND LEGAL ISSUES 15 · 37

- Effective security cannot be based on obscurity.
- Making vulnerabilities public may force vendors to improve the security of their products.

Since those who would use vulnerabilities and exploits for gain or harm often learn of them before system administrators and security experts, it makes sense—so the argument goes—for the Good Guys to spread the knowledge where it can do some good. One might also argue that, because cryptographic techniques are routinely exposed to public scrutiny by experts to detect vulnerabilities and to avoid future failures, other aspects of security should also be made public.

As for putting pressure on manufacturers, one colleague describes an incident that illustrates the frustrations that sometimes underlie full disclosure. He informed an important software supplier of a major security vulnerability. The product manager ignored him for a month. At that point, patience gone, the colleague informed the product manager that he had exactly one more day in which to produce a patch; otherwise, he said, he would publish the hole in full in the appropriate Usenet group. A patch was forthcoming within one hour.

Why would anyone object to full disclosure of detailed viral code and exploits? The arguments are that:

- Nobody except researchers needs to know the details of viruses or even of specific exploits.
- Publishing in full gives credibility to ill-intentioned Bad Guys who do the same.
- Full disclosure makes impressionable youths more susceptible to the view that illegal computer abuse is acceptable.

How, exactly, does publishing the details of a new virus help system administrators? In one view, such details should be exchanged only among colleagues who have developed trust in each other's integrity and who have the technical competence to provide fixes. For example, a "zoo" of computer viruses serves this function, with access limited to legitimate virus researchers who sign a code of ethics that forbids casually distributing viruses to anyone who wants samples. Opponents of this stance see the attitude as arrogant and elitist. Furthermore, some virus and worm code is written in a form that is easily read by anyone who receives a copy (a large population, considering that in 1999, the Melissa virus is thought to have infected over a million computers within a matter of days).

There is a danger in publishing exploits where any person with access to the Web can use them for automated attacks on Websites. This gives naïve people the impression that it is okay to publish any attack code, regardless of consequences. (Note that the consequences are often relatively minor—the author of the Melissa virus, which is estimated to have caused at least \$80 million in damages, served only 20 months in a federal prison and paid a fine of only \$5,000.) What is the difference, then, between publishing vulnerabilities or exploits and actually creating attack tools? Was creating and publishing BackOrifice a morally neutral or even a useful act? BackOrifice is a tool that is explicitly designed to install itself surreptitiously on systems and then hide in memory, using stealth techniques modeled on what some viruses use. Is this a contribution to security?

There are no simple or uncontested answers to these questions, but in some cases technology has answered them for us. Before 1995, when macro viruses first appeared

15 · 38 PENETRATING COMPUTER SYSTEMS AND NETWORKS

“in the wild,” most viruses were written in assembly language or machine code. This limited the number of people who could read them or understand them to persons familiar with assembly language and machine code. However, anyone who receives a macro virus has the text editor tools needed to read its code. The tools required to develop and test macro viruses are provided in mainstream applications such as Microsoft Word, which has tens of millions of users worldwide, a significant percentage of whom have, or can easily acquire, the rudimentary programming skills required to develop their own viruses. The rapid spread of the original Word Concept virus in the summer of 1995 pretty much ensured that anyone who wanted a copy of it could get one (and many of those who did not want a copy got one anyway). The idea of keeping the content of the virus secret was a nonstarter.

Faced with vendor inability or unwillingness to fix security vulnerabilities in a timely manner, some users and experts can be expected to turn to full disclosure of security information, even though many of them may deplore using such tactics. However, they will always run the risk of making the wrong call when it comes to the effect of such disclosures, which are inherently unpredictable. Indeed, the release of the Word Concept virus may have been motivated by a desire to make Microsoft change the way its office applications handle macros. Other office applications, such as Word Perfect and Lotus 1-2-3, were designed to keep macro code separate from document content, making a malicious document much harder to create.

15.4.3 Sources. There are many sources for information about how to penetrate systems. The motives behind these sources range from highly ethical to downright criminal. Chapter 74 in this *Handbook* discusses training and certification in penetration testing for legitimate, authorized purposes.

15.4.3.1 Online Sources. There is no small irony in the fact that much of what a person needs to know about how to penetrate information systems is made available by information systems. Fortunately, the inverse is also true, as one of this chapter’s authors has observed: “The best weapon with which to protect information is information.”¹⁶ For this reason, security professionals need to know what sources of penetration information are available.

Today there are thousands of Websites that

- Document security holes in different versions of operating systems
- Distribute hacking tools and discuss how to use them
- Catalog default credentials for network hardware
- Teach malicious code writing, including how to make viruses and worms
- List license codes to enable activation of pirated software
- Buy, sell, and trade system access codes, stolen credit cards, stolen identity data, and networks of compromised hosts (botnets)
- Provide a forum and meeting place for those seeking to penetrate systems

These sites have assumed the mantle of earlier online communication channels, such as bulletin boards and Usenet groups, where legitimate sharing of information occurred alongside the exchange of illegal information, pirated code, and so on. Some Websites are moderated and so maintain certain ethical standards. Others follow the Internet tradition of “anything goes.” System administrators and employees should

POLITICAL AND LEGAL ISSUES 15 · 39

never participate in discussions on these Websites with company email addresses. In a popular strategy, hackers wait for platform-specific vulnerabilities to be announced, then search the Internet for messages from people using that platform, look at their email addresses to see where they work, and attack systems at those companies in the hope that patches are not yet installed.

15.4.3.2 Publications. Over the years many publications have specialized in hacking. Often these provided information on how to penetrate systems. Some publications were primarily electronic, such as *Phrack*, while others have been print based, such as *2600*, which is now widely distributed through conventional magazine channels as well as through paid subscriptions.

15.4.3.3 Hacker Support Groups. Numerous groups of people exist to share hacking information, ranging from relatively stable entities with their own publications (e.g., *2600* and *cDc*) to annual conventions with thousands of participants (e.g., DefCon). Although some members of these groups may have committed criminal acts, and some participants at hacker conventions actually have been convicted of such and served time, there is usually a diverse mix of elements with different motivations in these groups and meetings. DefCon, for example, draws not only people who openly advocate unauthorized security testing of other people's systems and networks, but also law enforcement personnel and legitimate security experts. Some participants go to DefCon specifically to convince young people not to break laws while trying to learn about security.

Many security professionals would prefer that the line between white-hat hacking and black-hat hacking be clearer and more sharply enforced; however, some companies overlook past transgressions in order to gain the perceived value of the hackers' technical security expertise. Indeed, in recent years, investors have even cooperated with groups of hackers in founding security consulting companies, complete with some employees who continue to use hacker handles. The fact remains that some of the best technical training in the field is provided by people who gained their expertise in criminal (or quasi-criminal) hacking, but who now help defend against such activity.

15.4.4 Future of Penetration. Trends over the last 15 years strongly suggest that attempts to penetrate information systems will not decrease any time soon. These factors have been in play for some time:

- The declining cost of, and increased access to, penetration technology—from software and hardware used to crack passwords and encryption to eavesdropping and interception devices
- The continuing practice of fielding inadequately tested systems, built with immature technology and with insufficient attention to security
- The increased availability of automated hacking tools with easy-to-use interfaces
- The continuing allure, and portrayal in popular culture, of hacking as a “cool” activity, without adequate reflection on its legality or consideration of its morality

Additionally, these factors have emerged strongly in recent years:

- The very real opportunity to make money from penetrating systems, given a thriving market in purloined personal data, compromised hosts (bots), and exploits

15 · 40 PENETRATING COMPUTER SYSTEMS AND NETWORKS

- The increased interest and involvement of organized crime in system penetration
- The rise of transnational terrorist organizations that are increasingly computer-literate and may be inclined to penetrate systems belonging to entities or countries to which they are opposed¹⁷

Although some companies and government agencies are actively pursuing improved responses to these threats, others are not. Through lack of concern, resources, or time to address these trends, many entities are at increasing risk. Each new technology brings new threats, but new threats typically are discounted as scaremongering by vendors who offer defenses against them. For many companies and government agencies, the enthusiasm to reap the benefits of new technology overrules the warnings about risks inherent in its deployment. When those risks finally manifest themselves in ways that threaten the organization, the reaction is usually to buy a technical fix, while the root causes of vulnerability, namely human behavior and employee awareness of security issues, fail to receive the attention and resources they deserve.

New applications of computer technology, such as implanted medical devices and control systems for automobiles and autonomous vehicles, are providing fertile ground for experimentation by research scientists and criminal hackers, who are finding many vulnerabilities.¹⁸

15.5 SUMMARY

- Penetration of information systems is possible by means of a wide range of methods, some of which are very hard to defend against.
- Those responsible for securing systems have to defend against this wide range of penetration methods.
- Making sure all defenses against all attacks are effective all of the times is a lot harder than finding a single point of failure within those defenses.
- Although all systems do not need to defend against all types of attack equally, the cost of even the more exotic attack strategies is constantly falling, expanding the range of possible attackers.
- The cheapest and most effective attacks are often nontechnical, exploiting human frailty rather than weaknesses in the technology.
- Experienced criminal hackers tend to favor the nontechnical attack over the technical; and the best defense, employee awareness, is also nontechnical.
- Systems can be attacked at the client, at the server, or at the connection between the two.
- Both wired and wireless systems are highly susceptible to eavesdropping and interception.
- Many systems today are built with immature and insecure technology, making them susceptible to a wide range of attacks.
- New attacks come to light with alarming but predictable regularity.
- Many of these new attacks are old attacks reborn, due to a lack of vulnerability class analysis. As a result of economic pressures, faulty reasoning, and insufficient desire for security, vulnerabilities are fixed one instance at a time rather than one class at a time.

FURTHER READING 15 · 41

- Allowed path attacks against Websites are consistently the most effective strategy for system penetration whenever a system is Web connected or Web enabled.
- Penetration testing, both by internal staff and by objective external experts, always should precede system deployment.
- Given the inevitability of penetration attempts and the high probability of their eventual success, systems should be designed to survive attacks, limiting the scope of compromise from any single point of failure.
- Penetration of systems will continue to fascinate the curious and tempt them to break the law by illegally accessing systems. The potential gains from system penetration, in terms of money, power, competitive advantage, and notoriety, will continue to motivate those to whom laws and morality are not effective deterrents.
- Penetration of systems will become increasingly automated and simplified, further widening the range of possible attackers.
- Human nature, not technology, is the key to defense against penetration attempts. Only by raising society's ethical standards and educating employees to understand the willingness of others to behave unethically can the occurrence of criminal hacking into information systems be significantly reduced.

15.6 FURTHER READING**Websites**

CERIAS Hotlist: www.cerias.purdue.edu//hotlist

INFOSEC and INFOWAR Portal: www.infowar.com

SANS InfoSec Reading Room—Penetration Testing: www.sans.org/reading_room/whitepapers/testing

SearchSecurity.com: <http://searchsecurity.techtarget.com>

SecurityFocus: www.securityfocus.com

Web Application Security Consortium: www.webappsec.org

BooksAllen, Lee. *Advanced Penetration Testing for Highly-Secured Environments: The Ultimate Security Guide*. Packt Publishing, 2012.

Chappell, Laura. *Wireshark® 101: Essential Skills for Network Analysis*. Laura Chappell University, 2013.

Engebretson, Patrick. *The Basics of Hacking and Penetration Testing: Ethical Hacking and Penetration Testing Made Easy*. Syngress, 2011.

Faircloth, Jeremy. *Penetration Tester's Open Source Toolkit*, 3rd ed. Syngress, 2011.

Fialka, J. J. *War by Other Means: Economic Espionage in America*. New York: W. W. Norton, 1999.

Goodell, J. *The Cyberthief and the Samurai: The True Story of Kevin Mitnick—and the Man Who Hunted Him Down*. New York: Dell, 1996.

Kennedy, David. Jim O'Gorman, Devon Kearns, and Mati Aharoni. *Metasploit: The Penetration Tester's Guide*. No Starch Press, 2011.

Litchfield, D., Anley, C., Heasman, J., and Grindlay, B. *The Database Hacker's Handbook: Defending Database Servers*. Hoboken, NJ: John Wiley & Sons, 2005.

McClure, S., Scambray, J., and Kurtz, G. *Hacking Exposed: Network Security Secrets & Solutions*, 7th ed. New York: McGraw-Hill Osborne Media, 2012.

McGraw, G. *Software Security: Building Security In*. New York: Addison-Wesley Professional, 2006.

15 · 42 PENETRATING COMPUTER SYSTEMS AND NETWORKS

- O'Connor, T. J. *Violent Python: A Cookbook for Hackers, Forensic Analysts, Penetration Testers and Security Engineers*. Syngress, 2012.
- Scambray, J., Shema, M., and Sima, C. *Hacking Exposed: Web Applications*, 2nd ed. New York: McGraw-Hill Osborne, 2006.
- Schwartzau, W. *Pearl Harbor Dot Com*. Seminole, FL: InterPact Press, 2001.
- Shimomura, T., and J. Markoff. *Takedown: The Pursuit and Capture of Kevin Mitnick, America's Most Wanted Computer Outlaw—by the Man Who Did It*. New York: Hyperion, 1996.
- Slatalla, M., and J. Quittner. *Masters of Deception: The Gang that Ruled Cyberspace*. New York: HarperCollins, 1995.
- Sterling, B. *The Hacker Crackdown: Law and Disorder on the Electronic Frontier*. New York: Bantam Doubleday Dell, 1992.
- Stoll, C. *The Cuckoo's Egg: Tracking a Spy through the Maze of Computer Espionage*. New York: Pocket Books/Simon & Schuster, 1989.
- Stuttard, D., and Pinto, M. *The Web Application Hacker's Handbook: Discovering and Exploiting Security Flaws*. Hoboken, NJ: John Wiley & Sons, 2007.

15.7 NOTES

1. “Report on the Existence of a Global System for Intercepting Private and Commercial Communications (ECHELON Interception System),” PE 305.391, July 11, 2001, <http://cryptome.org/echelon-ep-fin.htm>
2. S. Cobb, *The Stephen Cobb Guide to PC & LAN Security* (New York: McGraw-Hill, 1992).
3. Robert Lemos, “Mitnick Teaches ‘Social Engineering’,” News.com, published on ZDNet News, July 17, 2000, http://news.zdnet.com/2100-9595_22-522261.html (URL inactive)
4. “Defector Smuggled Out Copies of the ‘Crown Jewels’ of Soviet Espionage,” *The Times (London)*, September 12, 1999.
5. Bryan Betts, “Optical Nets Easier to Hack than Copper,” *PCWorld*. April 27, 2007, www.pcworld.com/article/131306/article.html
6. Madeleine Acey, “‘Scrambler Software’ Will Protect Phone Calls from Prying Ears,” *CNN | Edge of Discovery*, August 17, 2012, www.cnn.com/2012/08/11/tech/silent-circle-encryption
7. AT&T Press Release, Dallas, TX, December 6, 2007, www.att.com/rss
8. Wim van Eck, “Electromagnetic Radiation from Video Display Units: An Eavesdropping Risk?” Cryptome, 1985, <http://cryptome.org/emr.pdf>
9. NSA CSS, “TEMPEST Level I,” NSA, January 15, 2009, www.nsa.gov/applications/ia/tempest/TEMPESTLevel1.cfm
10. Distributed Rainbow Table Project, “Free Rainbow Tables,” Free Rainbow Tables, May 8, 2013, <https://www.freerainbowtables.com>
11. Dan Farmer and Wietse Venema, “Improving the Security of Your Site by Breaking into It,” available on <http://nsi.org/library/compsec/farmer.txt>
12. Steve Gibson, *Shields UP!!* 2013, <https://www.grc.com/x/ne.dll?bh0bkyd2>
13. S. Gibson, “How Big Is Your Haystack ... and How Well Hidden Is YOUR Needle?” March 28, 2012, <https://www.grc.com/haystack.htm>
14. “According to Internet World Stats,” www.internetworldstats.com/stats.htm

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 16

MALICIOUS CODE

Robert Guess and Eric Salveggio

16.1 INTRODUCTION	16·1		
16.2 MALICIOUS CODE THREAT MODEL	16·2		
16.2.1 Self-Replicating Code	16·2	16.4.1 Signature-Based Malicious Code Detection	16·11
16.2.2 Actors: Origin of Malicious Code Threats	16·2	16.4.2 Network-Based Malicious Code Detection	16·11
16.2.3 Actors: Structured Threats	16·3	16.4.3 Behavioral Malicious Code Detection	16·12
16.2.4 Actors: Unstructured Threats	16·3	16.4.4 Heuristic Malicious Code Detection	16·12
16.2.5 Access versus Action: Vector versus Payload	16·3		
16.3 SURVEY OF MALICIOUS CODE	16·4	16.5 PREVENTION OF MALICIOUS CODE ATTACKS	16·12
16.3.1 Viruses	16·4	16.5.1 Defense in Depth	16·12
16.3.2 Worms	16·7	16.5.2 Operational Controls for Malicious Code	16·12
16.3.3 Trojans	16·8	16.5.3 Human Controls for Malicious Code	16·13
16.3.4 Spyware	16·9	16.5.4 Technical Controls for Malicious Code	16·13
16.3.5 Rootkits	16·10		
16.3.6 IRC Bots	16·10	16.6 CONCLUSION	16·14
16.3.7 Malicious Mobile Code	16·11	16.7 FURTHER READING	16·14
16.4 DETECTION OF MALICIOUS CODE	16·11	16.8 NOTES	16·15

16.1 INTRODUCTION. Malicious logic (or code) is “hardware, software, or firmware that is intentionally included in a system for an unauthorized purpose.”¹ In this chapter, we enumerate the common types of malicious code, sources of malicious code, methods of malicious code replication, and methods of malicious code detection.

A 2011 study of 200 small and medium businesses (SMBs) with up to 249 employees reported that “... two in five SMBs know with certainty that they have suffered some sort of security breach as a result of employees navigating to Web sites that host malware, infected downloads or have been corrupted by malicious code.”²

16 · 2 MALICIOUS CODE

Common types of malicious code include viruses, worms, Trojan horses, spyware, rootkits, and bots. Emerging malicious code threats include kleptographic code, cryptoviruses, and hardware-based rootkits. Present-day malicious code threats do not always fit into neat categories, resulting in confusion when discussing the topic. It is not possible to classify all code as being *good* code or *malicious* code. Absent the *mens rea*, or criminal intent of the author or user, code is neither good nor bad. Authors develop code to achieve some goal or fulfill some purpose just as users run code to achieve some goal or purpose. It is therefore the context of use and the intent of the wielder that determines whether code is malicious.

16.2 MALICIOUS CODE THREAT MODEL. In a threat model or profile, an *actor* uses *access* to target an *asset*, with an *action*, to yield an *outcome*.³ To understand the scope of the problem (and possible prevention), it is useful to study malicious code threats in this model. *Actors* may be structured or unstructured threats posed by individuals, organizations, or nation-states. *Access* is some allowed physical or logical path to the targeted *asset*. The execution of malicious code or logic is an action used to yield the desired *outcome*. This outcome could be intelligence (such as theft of trade secrets), surveillance, reconnaissance, disruption of operations, destruction of assets, publicity for some cause, or negative publicity for the victim.

16.2.1 Self-Replicating Code. Self-replicating code is not inherently malicious. A good deal of artificial intelligence research focuses on iterative self-replication. Hewlett-Packard and others⁴ have researched how techniques of self-replication could yield beneficial software such as code-patching worms. John Von Neumann proposed the basic concept of self-replicating code in his 1949 paper, “Theory of Self-Reproducing Automata.”⁵ In 1961 at Bell Labs, Douglas McIlroy, Victor Vyssotsky, and Robert H. Morris played a game called Darwin in which two opposing programs (memory worms) would enter a system and only one would leave. In an interview with one of the authors, Robert H. Morris the former Chief Scientist of the National Security Agency (NSA), stated:

We had this notion of putting two programs in a machine that would fight with each other and one would win and the other would die. Basically, the notion and the program were written by McIlroy and Vyssotsky and I was just on the side. But then, a few days later, I had a good idea and simply won the game and everyone else gave up ... just because I happened to hit upon a very good idea for writing it.⁶

Although it would be desirable to live in a world free of war, the ever-increasing integration of technology and warfare necessitates the development of offensive information warfare capabilities. Malicious code implants can serve useful intelligence, surveillance, and reconnaissance capabilities in national security and law enforcement operations. In addition, by researching and developing malicious code techniques, practitioners can better prepare for defense against developing threats.

16.2.2 Actors: Origin of Malicious Code Threats. Prior to discussing the specific types of malicious code, it is worth understanding the origin and source of malicious code attacks. Malicious code may originate from structured or unstructured threats. Structured threats include nation-states, corporate criminals, and organized crime. Unstructured threats include rogue actors such as individual intruders and so-called script kiddies.

MALICIOUS CODE THREAT MODEL 16 · 3

16.2.3 Actors: Structured Threats. Structured threats are organized, well funded, and may operate with a long-term strategic view. Structured malicious code threats tend to have intelligence, surveillance, and reconnaissance capabilities among their primary functions. Structured threats may use these capabilities to engage in industrial espionage, adversarial information operations, or serious ongoing fraud and theft.

Organized crime is responsible for 90 percent of malicious code threats.⁷ Extortionists target the online gambling industry with threats of distributed denial of service (DDoS) attacks launched from compromised personal computers. Criminals use compromised systems to engage in pump-and-dump stock fraud with one media outlet reporting annual gains by such groups approaching \$1 billion a year.⁸ Malicious software-for-hire incidents are on the rise, indicating a maturing marketplace for malicious code. In 2006, an Israeli court sentenced a couple to prison time and ordered them to pay fines for authoring and placing malicious code used to spy on corporations and individuals by members of telecommunications and trading firms.⁹

Coordinated, systematic attacks originating from China routinely target defense and research and development facilities. In the 1999 text *Unrestricted Warfare*, two Chinese air force officers called for a transformation of warfare that would involve ongoing technological attacks on western assets.¹⁰ A Chinese military white paper released in 2006 called for a strategy whereby China could win an “informationized war”¹¹ against the West that is “high-paced, high-technology, and digitized.” Although Chinese officials publicly deny sponsoring such attacks, the lack of law enforcement action indicates, at the very least, toleration for electronic attacks on Western assets.

16.2.4 Actors: Unstructured Threats. Unstructured threats include rogue actors not acting in concert or coordination with larger entities. Although serious attacks may result from unstructured threats, they do not pose the same long-term challenges as structured threats. In testimony before the United States Congress, former NSA director Kenneth Minihan stated:

The unstructured threat is random and relatively limited. It consists of adversaries with limited funds and organization and short-term goals. While it poses a threat to system operations, national security is not targeted. This is the most obvious threat today. The structured threat is considerably more methodical and well-supported. While the unstructured threat is the most obvious threat today, for national security purposes we are concerned primarily with the structured threat, since that poses the most significant risk.¹²

16.2.5 Access versus Action: Vector versus Payload. Malicious code attacks involve a vector and a payload. In biology, a vector is an agent that transfers (potentially harmful) material (code) from one location to another. In computer attacks, a vector is an avenue of *access*, such as an allowed path via physical access or via the network. Physical access occurs via internal personnel or others with access to the premises. Access via the network may occur via an allowed path to a Web server, an allowed path from a malicious Web server to a Web client, to a user via an email attachment, or to some other software process through an accessible port. The payload is a function (*action*) placed on the system to achieve some end. Payloads may include additional malicious logic, remote access software (rootkits and the like), or remote control (robot or bot) software to achieve some objective (spamming, DDoS attacks, and so forth).

16 · 4 MALICIOUS CODE

16.3 SURVEY OF MALICIOUS CODE

16.3.1 Viruses. In biology, a *virus* is “[a]n infectious organism that is usually submicroscopic, can multiply only inside certain living host cells, and is now understood to be a non-cellular structure lacking any intrinsic metabolism and usually comprising a DNA or RNA core inside a protein coat.”¹³ In computer terms, a virus is self-replicating code that requires a host executable or document and the aid of a human to replicate. Joe Dellinger created the first virus for the Apple disk operating system in 1982 while a student at Texas A&M.¹⁴ Fred Cohen created the first VAX computer virus in 1983 as a part of his doctoral research.¹⁵ Len Adleman, the A in RSA, first applied the biologic metaphor virus to describe the results.

Types of viruses include boot sector, file infector, macro virus, logic bomb, cross-site scripting viruses (really a form of worm), polymorphic viruses, and cryptoviruses. However, modern malicious code threats tend not to fall so neatly into categories.

16.3.1.1 Boot Sector Viruses. When users boot a computer from digital media, they allow code present in the boot sector of the media to run on the microprocessor with very little intermediation. Boot sector viruses in the 1980s and 1990s used this mechanism to spread via infected removable media such as (now extinct) floppy disks. If a user errantly booted from an infected floppy disk, compact disc (CD), DVD, or flash drive, the virus would copy itself to the boot sector of the hard drive. Thereafter, the malicious program code would copy itself to each medium inserted into the system. Although reportedly still in existence, boot sector viruses comprise very few modern malicious code threats. *Virus Bulletin* reports, “The last boot sector viruses fell off the WildList in early 2006, but various types continue to appear in our prevalence reports and a batch of laptops infected with ‘Stoned.Angelina’ was released in Germany and Denmark in mid-2007. More recently, [T]rojans have been observed using similar techniques to plant rootkits to hide their activities.”¹⁶

16.3.1.2 File Infector Viruses. File infector viruses inserted themselves into programs present on the host system. Whenever a user ran the host program (usually an .EXE or .COM file), the malicious code used the opportunity to insert itself into random access memory (RAM) and then replicate to other files on the system. Although probably still in existence, file infector viruses comprise relatively few modern malicious code threats.

16.3.1.3 Macro Viruses. Macro viruses spread via the macro definition languages used by some applications. The most widely abused is the Visual Basic for Applications (VBA) scripting language that Microsoft developed to allow the automation of functions inside the Office product suite. Any feature that allows for automation is a likely point of attack. Developers should carefully weigh security and ease of use when including such features in products. The first macro virus to target Microsoft Word¹⁷ appeared in the wild in 1995. Since that time, macro viruses have comprised a significant number of successful attacks.

16.3.1.4 Logic Bombs. A logic bomb is a form of malicious code function sometimes built into viruses that wait for some sequence of events to activate such as disappearance of an employee record from the human resources database or a particular

SURVEY OF MALICIOUS CODE 16 · 5

date and time. Kabay gave some examples of early logic bombs in a 2002 Network World article:

In 1985, a disgruntled computer security officer at an insurance brokerage firm in Texas set up a complex series of Job Control Language (JCL) and RPG (an old programming language) programs described later as “tripwires and time bombs.” For example, a routine data retrieval function was modified to cause the IBM System/38 midrange computer to power down. Another routine was programmed to erase random sections of main memory, change its own name, and reset itself to execute a month later.

In 1992, a computer programmer was fined \$5,000 for leaving a logic bomb at General Dynamics. His intention was to return after his program had erased critical data and get paid lots of money to fix the problem.

Time bombs are a subclass of logic bombs that “explode” at a certain time. Some of the first viruses, written in the 1980s, were time bombs. For example, the infamous “Friday the 13th” virus was a time bomb; it duplicated itself every Friday and on the 13th of the month, causing system slowdown. In addition, on every Friday the 13th it also corrupted all available disks. The Michelangelo virus from the early 1990s—one of the first viruses to make it into public consciousness because of news coverage—tried to damage hard disk directories on the 6th of March. The Win32.Kriz.3862 virus, discovered in 1999, detonates on Christmas day; its payload includes massive overwriting of data on all data storage units and also damage to the BIOS.

In 2000, a Stamford, Conn., man was indicted in New York State Supreme Court in Manhattan on charges of unauthorized modifications to a computer system and grand larceny. The defendant worked for Deutsche Morgan Grenfell starting in 1996 as a programmer. By the end of 1996, he became a securities trader. The indictment charged that he inserted a programmatic time bomb into a risk model on which he worked as a programmer; the trigger date was July 2000. The unauthorized code was discovered by other programmers, who apparently had to spend months repairing the program because of the unauthorized changes the defendant allegedly inserted.¹⁸

In 2002 an employee of UBS PaineWebber planted a logic bomb in his employer’s servers as part of an attempted stock manipulation scheme.¹⁹ The perpetrator, Roger Duronio, purchased numerous put option stock contracts allowing him to sell UBS stock at a fixed, high price. On March 4, 2002, at 9:30 A.M., the logic bomb began deleting files on over 1,000 computers, causing a reported \$3 million in damage. Duronio thought that the effects of the logic bomb would bring down the stock value of UBS, allowing him to make a large profit from his stock options. Duronio’s plot failed and his profit did not materialize. Federal authorities later charged him with securities and computer fraud. In 2006 a judge sentenced Duronio to 97 months in prison for the attack.²⁰

In March 2013, a logic bomb overwrote data on hard drives in South Korean banks and broadcasters. Kim Zetter of WIRED wrote,

The logic bomb dictated the date and time the malware would begin erasing data from machines to coordinate the destruction across multiple victims, according to Richard Henderson, a threat researcher for FortiGuard Labs based in Vancouver, the research division of the security firm Fortinet.

The attack, which struck machines on March 20, wiped the hard drives and master boot record of at least three banks and two media companies simultaneously. The attacks reportedly put some ATMs out of operation, preventing South Koreans from withdrawing cash from them.

The malware consisted of four files, including one called AgentBase.exe that triggered the wiping. Contained within that file was a hex string (4DAD4678) indicating the date and time the attack was to begin—March 20, 2013 at 2 pm local time (2013-3-20 14:00:00). As soon as the internal clock on the machine hit 14:00:01, the wiper was triggered to overwrite the hard drive and master boot record on Microsoft Windows machines and then reboot the system.²¹

16 · 6 MALICIOUS CODE

16.3.1.5 Cross-Site Scripting Viruses (or Worms). Cross-site scripting (XSS) viruses (or worms) replicate via flawed Web application servers and client code. A cross-site scripting exploit involving the social networking site MySpace and flaws in Microsoft Internet Explorer occurred in 2005 whereby a user named Samy amassed over 1 million friends overnight using a JavaScript insertion bug.²² The Los Angeles Superior Court sentenced the author, Samy Kamkar, to three years' probation and 90 days of community service for his actions.

In an excellent review of XSS, Sherif Koussa wrote,

Cross-site scripting is an attack that targets the application users. The simplicity of cross-site scripting is also why it is so powerful. If the application is vulnerable to cross-site scripting, the developer is not in charge of what runs on the user's browser anymore ... the attacker is. Cross-site scripting could be used in attacks like authentication hijacking and session hijacking. The power of cross-site scripting manifests itself even more when combined with cross-site request forgery.²³

He then cites the case of the Twitter attack below:

In 2009, Mikey Mooney, then 17 years old, claimed responsibility for attacking Twitter using XSS techniques: "... at least four separate variants of the original StalkDaily.com XSS worm hit the popular micro-blogging site Twitter, automatically hijacking accounts and advertising the author's web site by posting tweets on behalf of the account holders, by exploiting cross site scripting flaws at the site." Writer Dancho Danchev provided more details about the attack in his ZDNet article.²⁴

16.3.1.6 Polymorphic Viruses. Polymorphic code modifies itself to evade detection. The virus code may accomplish this by dynamically reassembling itself to modify the underlying structure while retaining overall functionality. In other methods, the virus is encrypted, encoded, or packed, and a stub loader decrypts, decodes, or unpacks the virus at run time. UPX, the Ultimate Packer for eXecutables, is currently the most widely used packer format.

In the SOPHOS Security Threat Report 2013, the authors write:

Polymorphism is not a new idea—malware authors have been using it for 20 years. Simply stated, polymorphic code changes its appearance in an attempt to avoid detection, without changing its behavior or goals. If a program looks different enough, attackers hope, antivirus software might miss it. Or the antivirus software might be forced to generate too many false positives, leading users to disable it.

In a polymorphic attack, code is typically encrypted to appear meaningless and paired with a decryptor that translates it back into a form that can be executed. Each time it's decrypted, a mutation engine changes its syntax, semantics, or both.

For instance, Windows malware authors have often used structured exception handling to obfuscate control flow and make it tougher to perform static analysis of programs before they run.

Traditional polymorphic viruses are self-contained and must contain the mutation engine in order to replicate. Sophos and other security companies have become adept at detecting these forms of malware. With access to the mutation engine, it's easier to analyze its behavior.

Today attackers are rapidly moving to web-distributed malware relying on server-side polymorphism (SSP). Now, the mutation engine and associated tools are hosted entirely on the server. Criminals can use these tools to create diverse file content on the fly. Recipients of this content (whether it is a Windows .exe, Adobe PDF, JavaScript, or anything else) see only one example of what the engine can create. They don't get to see the engine itself.²⁵

SURVEY OF MALICIOUS CODE 16 · 7

16.3.1.7 Cryptographic Viruses. Cryptoviruses use encryption to encrypt data, making it inaccessible to the user. Although some viruses use symmetric algorithms for self-encryption to evade detection, such algorithms are inappropriate for data encrypting viruses, as any copy of the virus is going to yield a copy of the symmetric key. For this reason, proper cryptoviruses use asymmetric techniques. The Gpcode virus encrypts files by extension (.xls, .doc, and the like) and leaves behind a text message prompting the user to email a given address to pay a ransom for access to their files. The next generation of cryptoviruses will likely use hybrid cryptosystems.

In a 2012 report on cryptoviruses, researchers developed proof-of-concept models of viruses using the public-key cryptosystem (PKC) to conceal themselves from scanners and then activate in response to an encrypted key or ticket. Their list of applications of cryptography in malware includes:

- Resist reverse engineering
- Improve anonymity of communications from controllers to the malware
- More effective data theft and denial-of-service attacks
- Remote-control back doors for extortion²⁶

16.3.2 Worms. The term *worm* refers to any form of self-replicating code that does not integrate into executable code. Common worm vectors include vulnerable services, email, instant messaging applications, and open file shares. Many worms use multiple vectors. For example, the Nimda worm spread via email, Web server vulnerabilities, hosted malicious Webpages, and open file shares.²⁷

The first large-scale Internet worm infection was the Morris (aka Internet) Worm released by Robert T. Morris on November 2, 1988, while he was a student at Cornell University. The worm exploited flaws in the Sendmail and finger services to replicate and infected nearly 10 percent of Internet hosts. Although he claimed that this was an experiment gone awry, Morris became the first person convicted under the Computer Fraud and Abuse Act.

In 2001, Nicholas Weaver coined the term “Warhol Worms”²⁸ for those worms that had the capability of propagating very quickly (taking a cue from Warhol’s famous quip that in the future everyone would have 15 minutes of fame). The SQL Slammer worm became the first such worm when it infected approximately 90 percent of vulnerable systems in 10 minutes.²⁹ The Slammer worm replicated due to a flaw in the Microsoft SQL database server. Because this was installed along with other components like Visual Studio, many people did not even know that they were running an SQL server. Slammer was exceedingly effective, because a Microsoft service pack downgraded a previously patched dynamic link library (dll) to an older, vulnerable version. Specifically, Microsoft issued four updates to *ssnetlib.dll* in 2002. Unfortunately, in October, hotfix Q317748 downgraded *ssnetlib.dll* to a vulnerable version, ironically making those who most faithfully applied patches and hotfixes the most vulnerable to Slammer.³⁰

The SQL Slammer worm is an interesting case study because it was a vector without a payload. The 376-byte worm ran in memory without touching the hard disk, leading some to believe that Slammer was an experiment in propagation techniques. Since it was based on the User Datagram Protocol (UDP) and traveled in a single packet, the overhead was minimal and the propagation rate was greater than any previous malicious code threat.³¹ The fact that the author(s) released it on a Saturday is also curious. It is unknown why a malicious attacker would deliberately release such a rabid

16 · 8 MALICIOUS CODE

virus on a day of the week that would lessen the overall business impact. If the author had released the worm on a peak business day such as a Tuesday, the damage caused would have been much more severe.

Worms are used to distribute other forms of malware such as Trojans, spyware, rootkits, and remote command and control channels called *bots*. An indication of the overall maturity of malicious code attacks is the Bagel worm, which has been in production at least since 2004. Bagel, like many other mail worms, arrives as an email with a malicious attachment. When the user runs the attachment, the payload is executed, which does a number of things depending on the variant. Later variants of Bagel install an open application framework that the remote attackers may update and extend. The harvested systems deliver spam (unsolicited commercial email), harvest additional addresses, and act as a staging point for other attacks. The author(s) appear to follow a sophisticated software development methodology and testing process.³² In 2007, attackers released 30,000 new variants in a six-week time span.³³ In April 2012, the Flashback worm “infected more than 650,000 Mac OS X systems using a vulnerability in Apple’s version of Java.”³⁴ Analysts with F-Secure wrote,

Flashback is the most advanced OS X malware we’ve ever seen. It boasts a series of firsts for its kind. It was both the first to be VMware-aware and the first to disable XProtect, OS X’s built-in malware protection program. Both these features were removed from later variants (the former presumably to avoid heuristic detections, and the latter presumably once the authors realized it was unnecessary, as XProtect was not designed to protect against non-quarantine files). Their removal indicates that Flashback is actively being reviewed and improved by its authors.

Another interesting first is Flashback’s exploitation of an unpatched vulnerability in the Java distribution of OS X, which allowed it to infect more than 650,000 Macs around the world. . . . This made Flashback roughly as common for Macs as Conficker was for Windows. . . . This means Flashback is not only the most advanced, but also the most successful OS X malware we’ve seen so far.

Flashback’s infection strategy is explicitly designed to select unprotected systems and will not infect a machine if certain security software or analysis tools are found. This implies that Flashback’s authors are targeting less security-conscious users, at the expense of the total number of potential targets. This turns out to be an effective strategy, as security researchers had difficulties getting sufficient samples from users. It took a mistake on the part of Flashback’s author to alert users to the presence of an infection and subsequently, to lead to the mass discovery of the malware.³⁵

16.3.3 Trojans. A Trojan horse application, like the horse of Greek mythology, carries both an overt function and a covert function. Although attackers may use worms as one possible propagation vector, Trojans require that a user run the malicious program in order to be effective. Trojans tend to use some form of social engineering or manipulation to persuade the user to run the program. Email worms may appear to originate from a known associate and thereby trick the user. Other Trojans may take the form of games, free offers, pictures of popular celebrities, or files on peer-to-peer file sharing services. One study of the Limewire peer-to-peer file sharing service found that “68% of all downloadable responses containing executable, archival, and Microsoft Office file extensions”³⁶ contained malware. Queries for movies were most likely to hold malicious code.

The covert function of Trojan horse application is typically some form of a remote access Trojan, keylogger, dialer, IRC bot, or rootkit. Remote access Trojans provide full remote access to the system. The developers of the Bo2 K Trojan bill it as “the

SURVEY OF MALICIOUS CODE 16 · 9

most powerful network administration tool available for the Microsoft environment” and insist upon the fact that it is functionally no different from other remote access solutions.³⁷ Keylogging Trojans record keystrokes and periodically upload the data to a remote user. Attackers carried out the Windows 2000 source code theft using credentials stolen by a QAZ Trojan installed on a remote workers computer. A dialer is a form of Trojan that silently dials remote toll numbers and runs up a large telephone bill for the victim. Internet relay chat (IRC) bots act as autonomous IRC clients. Early IRC bots provided technical support and channel monitoring features but are now widely used for malicious purposes such as command and control capabilities.

In 2013, Kaspersky Labs reported on “a Trojan build specifically for Android smartphones.”³⁸ sean Gallagher wrote,

On March 25, the e-mail account of a Tibetan activist was hacked and then used to distribute Android malware to the activist’s contact list. The e-mail’s lure was a statement on the recent conference organized by the World Uyghur Congress. . . If the targets opened the attachment . . . they received malware packaged in an Android APK file.

When opened, the Trojan installs an app called “Conference” on the Android devices’ desktops. If the app is launched, it displays a fake message from the chairman of the WUC—while sending back a message to a command and control server to report its successful installation. The malware provides a backdoor to the device via SMS messages sent by the server. On command, it returns the phone’s contact lists, call logs, data about the smartphone, its geo-location data, and any SMS messages stored on it to a server via a Web POST upload.

The server itself is running on a Chinese-language configured Windows Server 2003 machine sitting in a data center in Los Angeles. In addition to providing an upload point for the data stolen from Android devices, it also hosts more Android malware in its home page and provides a public Web interface (in Chinese) that allows direct control over phones that have been infected with the malware. While the server itself is at an IP address registered to a company called Emagine Concept, a domain pointed at the machine is registered to Shanghai Meicheng Technology Information Development Co., Ltd., a Chinese company with a contact in Beijing.³⁹

This Trojan thus distributes spyware, the next topic.

16.3.4 Spyware. The term *spyware* refers to any software that collects user information without consent. Common varieties of spyware collect information on Web usage, serve advertising content (pop-ups), log keystrokes, engage in click fraud, or monitor program usage and licensing. Unauthorized access or exceeding authority on a computer system is a violation of the Computer Fraud and Abuse Act. Although some spyware may be illegal, some developers insist that they are engaged in a legitimate business activity. Such spyware provides an end user license agreement (EULA) allowing the user to opt out of installation. Since the user must grant permission for the software to install, it is almost certainly legal. However, whether this is ethical and conscionable is a different matter. Sony attracted significant negative publicity for their use of spyware technology to limit the ability of listeners to copy music compact discs. Due to this activity, Sony faced legal charges in multiple states and in 2005 settled a class-action lawsuit. Informed consent is the rule of thumb in any monitoring system. Organizations should endeavor to be forthright in these matters in order to maintain a positive public profile and to avoid legal challenges.

In May 2013, news surfaced about the use of spyware to monitor the communications of a Bahrain human-rights activist. In a court hearing, details arose describing

16 · 10 MALICIOUS CODE

how Dr Ala'a Shehabi was monitored through infection of her systems through a phishing attack:

According to her witness statement, a few weeks after her arrest Shehabi received a series of emails, the first purportedly from Kahil Marzou who was the deputy head of Bahrain's main opposition party, including one containing a virus. Other emails that claimed to be from an Al Jazeera journalist were also infected. Research found that the emails contained a product called FinSpy, distributed by a British company, Gamma International.

The witness statement claims that when a person's computer is infected with FinSpy, "it allows access to emails, social media messaging, and Skype calls, as well as copying the files saved on the hard disk. These products also enable whoever is doing the targeting to commandeer and remotely operate microphones and cameras on computers and mobile phones."⁴⁰

16.3.5 Rootkits. A rootkit consists of a set of tools for covertly compromising a system and maintaining administrative (root) access for the intruder. Present-day rootkits compromise a system at the application, library, kernel, hypervisor, or hardware level. Early application-level rootkits targeted UNIX-like systems and replaced standard system utilities (netstat, ps, and the like) with versions that omitted information on the intruder, such as open ports, running processes, open files, and other activity. API-level rootkits modify or patch the system call table to redirect system calls (like an API-level monkey-in-the-middle attack). Kernel-level rootkits run as device drivers and dynamically load into the kernel; compromising the integrity of the core of the operating system to filter information presented to users via application-level processes.

Nearly all modern operating systems make use of virtualization for memory and process management. An operating system normally runs in ring 0 of the microprocessor; the most privileged level. A hypervisor is a layer of code between the operating system and the hardware that fools the operating system into believing that it is running in ring 0. The hypervisor monitors and arbitrates exchanges between virtual machines and the real hardware. Hypervisor design is a key part of the trusted computing framework, and microprocessor designers like Intel and AMD are including enhanced hardware-based virtualization capabilities in their products. An emerging generation of rootkit technology uses this framework to create potentially undetectable code.^{41,42}

Other potential rootkit threats come from insiders, hardware designers, and manufacturers. An attacker with physical access could insert a peripheral component into a computer to create a logically undetectable rootkit. A hardware manufacturer or designer could include rootkit-like capabilities in a microcircuit. The 2005 sale of IBM's Personal Computer Division to Chinese manufacturer Lenovo caused enough concern in the United States to prompt a national security review of the transaction.⁴³ The review by the House Committee on Foreign Investment in the United States (CFIUS) eventually approved the sale despite concerns expressed by Reps. Henry Hyde and Don Manzullo of Illinois and Rep. Duncan Hunter of California.⁴⁴

16.3.6 IRC Bots. IRC bots are autonomous agents that make use of Internet Relay Chat (IRC) to offer interactive services like channel monitoring, support, information services, and games. Greg Lindhal wrote the first IRC bot GM (game master), which led users through the text-based role-playing game "Hunt the Wumpus." In 1999, the PrettyPark worm became the first worm to use IRC as a remote control channel. Infected systems would check into an IRC server and channel to download updates and upload stolen data. Attackers create malicious bots to carry out DDoS attacks (DDoS

DETECTION OF MALICIOUS CODE 16 · 11

bots), send unsolicited commercial email (SpamBots), and engage in exploitation, theft, or fraud.

Examples of current bot frameworks are GTbot, SDbot, Agobot, Goabot, Randex, Spybot, and Phatbot. The current generation of bot technologies includes keylogging, port scanning, exploitation, packet sniffing, process hiding, and adware fraud capabilities. A single network of these compromised systems may reach more than 100,000 bots. Researchers in 2004 estimated that there were over one million bots currently connected to the Internet.⁴⁵ Bot herders rent out these massive distributed networks systems to criminal organizations. In 2006, a California court sentenced one such bot herder, Jeanson Ancheta, to five years in prison for conspiring to violate the Computer Fraud and Abuse Act, conspiring to violate the CAN-SPAM Act, causing damage to computers used by the federal government in national defense, and accessing protected computers without authorization, in order to commit fraud.⁴⁶

16.3.7 Malicious Mobile Code. Web servers may host pages containing malicious mobile code. ActiveX controls, Java applets, JavaScript, Adobe Flash animations, and any other type of dynamic executing code can download to a user's system and run in their context with all associated privileges. Malicious Web servers are one method for dropping Trojans and bots onto computers. An attacker may use a spam or phishing attack to encourage users to click on a link embedded within an email. Another vector is to use a similar domain name to a legitimate organization. Yet another vector is for the site to offer information in some matter. Once indexed in the major search engines, this content will draw a number of users to the site. See Chapter 17 in this *Handbook* for further details on malicious mobile code.

16.4 DETECTION OF MALICIOUS CODE. Common methods of detecting malicious code include signature-based, network-based, and behavioral heuristic techniques. However, as far back as in the 1984 work *Computer Viruses—Theory and Experiments*, Fred Cohen demonstrated that the only way one could prevent all possible viral code would be through isolation.⁴⁷ While there are numerous methods of detecting malicious code, no one technique works on every variety or in every circumstance, and there is always some method for successfully evading detection.

16.4.1 Signature-Based Malicious Code Detection. Some of the oldest methods of malicious code detection are signature-based methods that utilize known strings or patterns in the code. Signature-based methods are easy to implement and impose very low overhead on the system but are just as easily evaded. Polymorphism and metamorphism⁴⁸ are two methods by which malware may change form over time and thereby evade signature-based detection. Detection systems can use hash functions to fingerprint a given binary program or code fragment. However, the hash will fail to match on any modified version, making this method reliable but not always useful. Overall, signature-based methods are not terribly reliable, although signatures are one useful metric in more complex heuristics.

16.4.2 Network-Based Malicious Code Detection. Network-based methods of malicious code detection look for network artifacts associated with malicious code, such as a connection to a server as in the case of a keylogging Trojan or Internet Relay Chat (IRC) bot. Network anomaly detection works well but is expensive and poorly understood by many security practitioners. A simple method for

16 · 12 MALICIOUS CODE

network-based detection is to analyze network flow (netflow) data in comparison with a statistical database of known good traffic. However, malicious code that acts normal can bypass such methods of detection. Predicting what is normal for a given environment may be difficult for an outsider, but analyzing certain activity like domain name system, email, and Web usage can produce generic models for evasion.

16.4.3 Behavioral Malicious Code Detection. Behavioral methods of detecting malicious code analyze the actions of running software to look for illegitimate activity. This could be opening a port, connecting to a remote host, or modifying the system call table or other memory areas. Behavioral methods of detection will fail if the malware acts normal or if the malicious code can target and exploit the detection system itself. A slow but effective method of detecting potentially malicious code is to use a virtual machine approach, whereby the system allows code to execute in a sandboxed virtual machine. This allows for a full functional analysis of the code, but the current throughput of such systems makes them not useful for today's high-speed production environments.

16.4.4 Heuristic Malicious Code Detection. Heuristics that are more complex may use both statistical and behavioral models to determine a relative score to a normal corpus (statistical database in this case) of legitimate behavior. Bayesian analysis is widely used in the detection of spam. However, this method can also detect new variants of existing malicious code with a high degree of accuracy. N-gram analysis is a form of frequency analysis borrowed from natural language processing that can model software and fingerprint data types. This method is useful for detecting executable code embedded within other data objects. In short, there is no single, monolithic way to detect all malicious code. Spinellis demonstrated that reliable detection of malicious code is NP-complete,⁴⁹ meaning that the dilemma is not solvable in polynomial time. However, the technologies discussed, used in concert, are relatively effective in detecting many common malicious code threats.

16.5 PREVENTION OF MALICIOUS CODE ATTACKS

16.5.1 Defense in Depth. As the problem of malicious code is demonstrably NP-complete (not solvable), a single antivirus program cannot protect against all malicious code threats. One strategy, called *defense in depth*, uses operational, human, and technical controls. Another strategy is to build networks and applications that only function one correct way. Such *orthogonal* networks and applications are a rarity as they are expensive and difficult to design, and many firms will resist imposing rigid limitations on the enterprise.

16.5.2 Operational Controls for Malicious Code. All organizations must create written policies and procedures regarding the introduction of program code into the operating environment. Policies should define what persons the firm allows to install programs, acceptable use for Internet access, acceptable use for email systems, and what to do if users suspect the compromise of a system or user. Organizations should subject all new employees to some level of background investigation. This should include, at a minimum, a criminal records search, verification of all references and credentials, and a credit report. Employers should ensure that the prospective employees are forthright and truthful on their application. If employees lie prior to employment, the odds are that they will lie later as well.

PREVENTION OF MALICIOUS CODE ATTACKS 16 · 13

16.5.3 Human Controls for Malicious Code. All users (including executives) should receive training on the policies and procedures of the organization. Due to the evolving nature of the malicious code threat, the organization should update and refresh this training annually at the very least. The training sessions should introduce the prevalent types of threats, how to detect the threats, and proper response. Currently, this training should include identification of advance-fee fraud (also known as Nigerian 419 frauds), social engineering attempts, and detection of malicious attachments. Users should notify the help desk or other entity if they encounter any anomalous system or user behavior.

16.5.4 Technical Controls for Malicious Code. For more detailed information on antivirus technology, see Chapter 41 in this *Handbook*.

16.5.4.1 Implementing Antivirus Systems. Implement an antivirus (A/V) solution that suits the operational environment. This solution should include both network-based and host-based systems. Network-based systems function inline like a gateway, while host-based systems run on host end points. These systems should come from different vendors, as there is less benefit when using the same software on the network and hosts. A diverse detection and containment strategy will help firms avoid the pitfalls associated with active malcode exploitation of antivirus software, as occurred in 2004⁵⁰ and again in 2006.⁵¹ The last thing that organizations should permit when designing an A/V strategy is for antivirus software to serve as the vector for an active malcode attack.

The A/V solution should provide a mechanism for dynamically applying updates and should do so daily, at the least. Email systems may require a separate inline appliance to detect email-borne malicious code as well as spam, fraud, and phishing attacks. By using a diversity of detection approaches, organizations can attain higher rates of detection than by using a single stand-alone product.

16.5.4.2 Host Configuration Controls and Security. Host configuration can mitigate many malware threats prior to emergence. Implement a form of automatic updates (patches, etc.) that supports the operating environment. Many technical malicious-code threats (e.g., worms) target well-known and patched flaws. Eliminate all noncritical software and services. This will help minimize threats that target the ever-burgeoning code complexity that security professionals face. In one study, programmers trained in the capability maturity model for secure software development continued to make 4.5 errors per 1,000 lines of code.⁵² For an operating system like Microsoft Vista, which one estimate puts at 50 million lines of code,⁵³ extrapolating this statistic means that there are likely at least 225,000 code errors. Disabling or removing as much of this code as possible is a reasonable preventive control. If the environment permits, remove all Web browsers. If the environment does not allow this, lock down the browser configuration and use a secure Web proxy.

16.5.4.3 Network-Based Security Controls. A layered defense of routers, firewalls, proxies, and switched virtual local area networks (VLANs) can mitigate malicious code propagation. At the router, filter all inbound bogus network addresses (BOGONs), Request for Comment (RFC) 1918 addresses, and spoofed internal addresses per RFCs 2267⁵⁴ and 3704.⁵⁵ If the enterprise faces specific threats from certain nations, filter the network blocks allocated to those nations at the border router.

16 · 14 MALICIOUS CODE

Use the current best practices in firewall configuration. (These change rapidly.) Use a secure, authenticated Web proxy, and force all clients to use the proxy. (No user should be able to access the Web directly.) Disable all unused LAN access ports, or consider using 802.1x authentication to do so automatically. Segment the network into functional workgroups using VLANs or physically separate switch architecture. If possible, implement access control lists between VLANs or switched segments. The goal is to make the network as orthogonal as possible. If the network can function only one right way, the organization will avoid many automated malicious code threats entirely.

16.5.4.4 Network Monitoring. Prevention of malicious code threats is ideal, but detection is critical. Most malicious code attacks will have some artifact, whether host based or network based. To detect these artifacts, organizations should establish a security information management system that aggregates device logs, server logs, host logs, intrusion detection system alerts, and network flow data. Network flow data is a useful tool for detecting anomalous network activity. A network flow is a 5-tuple consisting of a source address, destination address, source port, destination port, and protocol with an associated time stamp. Any malicious network activity is going to have an associated flow. Stealthy malware, however, will attempt to act normal in order to evade detection. Detecting anomalous activity requires that operators understand what is normal for the environment. Without a historical statistical database or extensive experience, this can be difficult. Network anomaly detection (NAD) uses a statistical modeling methodology⁵⁶ to use this data to detect statistical outliers. This is extraordinarily effective for detecting new malicious code threats as well as other forms of anomalous behavior. However, a lack of understanding appears to limit the adoption of NAD. This is one area where many organizations could improve their practice of information security controls.

16.6 CONCLUSION. Malicious code threats are as numerous as the variety of nonmalicious code. Prevention of all malicious code is not possible as the problem is demonstrably NP-complete.⁵⁷ However, a strategy of defense in depth that uses operational, human, and technical controls can be relatively effective. Used properly, the current generation of technical controls available to organizations is effective in stopping a majority of malicious code threats. However, the trusted insider typically has the access needed to turn a threat into a reality of operational risk. Current trends in malicious code threats indicate a continued pattern of organized criminal involvement and international espionage that continues to target the weakest link in the security chain: the human being.

16.7 FURTHER READING

- Blunden, Bill. *The Rootkit Arsenal: Escape and Evasion in the Dark Corners of the System*, 2nd ed. Jones & Bartlett Learning, 2012.
- Butler, Jamie and Greg Hoglund. *Rootkits: Subverting the Windows Kernel*. Addison-Wesley, 2005.
- Harley, David., et al. *Avien Malware Defense Guide for the Enterprise*. Burlington, MA: Syngress, 2007.
- Harley, David. *OS X Exploits and Defense*. Burlington, MA: Syngress, 2008.
- Sikorski, Michael. *Practical Malware Analysis: The Hands-On Guide to Dissecting Malicious Software*. No Starch Press, 2012.

NOTES 16 · 15

- Szor, Peter. *The Art of Computer Virus Research and Defense*. Cupertino, CA: Symantec Press, 2005.
- Young, Adam, and Moti Yung. *Malicious Cryptography: Exposing Cryptovirology*. Hoboken, NJ: John Wiley & Sons, 2004.

16.8 NOTES

1. NCSC, “Glossary of Computer Security Terms,” Version 1, 1988, NCSC-TG-004 in the Rainbow Series, www.fas.org/irp/nsa/rainbow/tg004.htm
2. GFI, “GFI Software Survey: 40% of SMBs Have Suffered a Security Breach Due to Unsafe Web Surfing,” *EON*. October 12, 2011, <http://eon.businesswire.com/news/eon/20111012005707/en/web-reputation/web-security-survey/small-business-web-monitoring-software>
3. C. Alberts and A. Dorofee, “OCTAVE Threat Profiles,” 2001, www.cert.org/archive/pdf/OCTAVEthreatProfiles.pdf
4. R. Lemos, “Good Worms Back on the Agenda,” *Security Focus*, January 27, 2006, www.securityfocus.com/news/11373
5. John von Neumann, “Theory of Self-Reproducing Automata,” *Part 1: Transcripts of Lectures Given at the University of Illinois, Dec. 1949*, ed. A. W. Burks (University of Illinois, 1966). Urbana, IL.
6. R. H. Morris, Recorded interview with Robert Guess, July 2005.
7. D. Llet, “Antivirus Firm Says Organized Crime Growing Online,” *ZDnet News*, December 9, 2004, http://news.zdnet.com/2100-1009_22-5486201.html (URL inactive).
8. J. Hoskyn, “SEC Targets Pump and Dump Scammers,” *IT Week*, March 9, 2007, www.itweek.co.uk/vnunet/news/2185164/sec-targets-pump-dump-stock (site discontinued).
9. SOPHOS, “Court hands hefty fine and jail sentence to Israeli spyware couple, reports Sophos,” *SOPHOS | Press Releases*, March 26, 2006, www.sophos.com/en-us/press-office/press-releases/2006/03/israelspyduo.aspx (URL inactive).
10. Quio Liang and Wang Xiangsui, *Unrestricted Warfare* (Beijing: PLA Literature and Arts Publishing House, 1999), www.terrorism.com/documents/TRC Analysis/unrestricted.pdf (Access Permission required).
11. J. Rogin, “Cyber Officials: Chinese Hackers Attack ‘Anything and Everything,’” February 13, 2007, www.fcw.com/article97658-02-13-07-Web&printLayout (URL inactive).
12. Kenneth A. Minihan, “Prepared Statement before the Senate Governmental Affairs Committee,” June 14, 1998, www.senate.gov/~gov_affairs/62498minihan.htm (URL inactive).
13. *Oxford English Dictionary*, 2007.
14. J. Dellingar, Usenet post, December 2, 1987, http://yarchive.net/risks/early_virus.html
15. F. Cohen, “Experiments with Computer Viruses,” 1984, www.all.net/books/virus/part5.html
16. Virus Bulletin, “Boot sector virus,” *Virus Bulletin | Glossary*. 2013, www.virusbtn.com/resources/glossary/boot_sector_virus.xml
17. Microsoft, “What to Do If You Have a Macro Virus,” October 4, 2002, <http://support.microsoft.com/kb/181080/en-us>

16 · 16 MALICIOUS CODE

18. M. E. Kabay, "Logic Bombs, Part 1." *Network World | Security Newsletter*. August 21, 2002, www.networkworld.com/newsletters/sec/2002/01514405.html
19. U.S. Department of Justice, United States Attorney District of New Jersey, "Disgruntled UBS PaineWebber Employee Charged with Allegedly Unleashing 'Logic Bomb' on Company Computers," December 17, 2002.
20. U.S. Department of Justice Press Release, "Former UBS Computer Systems Manager Gets 97 Months for Unleashing 'Logic Bomb' on Company Network," 2006, www.usdoj.gov/usao/nj/press/files/pdffiles/duro1213rel.pdf (URL inactive).
21. Kim Zetter, "Logic Bomb Set Off South Korea Cyberattack," *Wired*, March 21, 2013, www.wired.com/threatlevel/2013/03/logic-bomb-south-korea-attack
22. E. Lai, "Teen Uses Worm to Boost Ratings on MySpace.com," *ComputerWorld*, October 17, 2005, www.computerworld.com/securitytopics/security/holes/story/0,10801,105484,00.html
23. Sheif Koussa, "Security Code Review Techniques: Cross-Site Scripting Edition," *Developer Connection | Jonathan Rozenblit*, March 5, 2013, <http://blogs.msdn.com/b/cdndevs/archive/2013/03/05/security-code-review-techniques-cross-site-scripting-edition.aspx>
24. Dancho Danchev, "Twitter Hit by Multiple Variants of XSS Worm," *ZDNET | Security*, April 14, 2009, www.zdnet.com/blog/security/twitter-hit-by-multiple-variants-of-xss-worm/3125
25. SOPHOS, "Security Threat Report 2013: New Platforms and Changing Threats," *SOPHOS*, February 1, 2013, [www.sophossecuritythreatreport2013.pdf](http://www.sophos.com/en-us/mediabinary/PDFs/other/sophossecuritythreatreport2013.pdf)
26. Shafiqul Abidin, Rajeev Kumar, and Varun Tiwari. "Review Report on Cryptovirology and Cryptography," *International Journal of Scientific & Engineering Research*, 11, 2012: 1–4.
27. CERT. "Advisory CA-2001-26 Nimda Worm." Original release date: September 18, 2001, Revised: September 25, 2001. Source: CERT/CC, www.cert.org/advisories/CA-2001-26.html
28. N. Weaver, "Warhol Worms: The Potential for Very Fast Internet Plagues," 2001, www.iwar.org.uk/comsec/resources/worms/warhol-worm.htm
29. D. Moore, V. Paxson, S. Savage, C. Shannon, S. Stanford, and N. Weaver, "The Spread of the Sapphire/Slammer Worm," CAIDA, 2003, www.caida.org/publications/papers/2003/sapphire/sapphire.html
30. S. Berinato, "Patch and Pray," *CSO Online*. August 14, 2003, www.csoonline.com.au/index.php?id=1337625166;fp;8;fpid;5
31. T. Vogt, "Simulating and Optimising Worm Propagation Algorithms" September 29, 2003, <http://web.lemuria.org/security/WormPropagation.pdf>
32. "Year of the Beagle: The Beagle Worm History, Part III," *Infection Vectors Website*, February 2005, www.infectionvectors.com/vectors/year_of_the_beagle.htm
33. CommTouch, "Malware Outbreak Trend Report: Bagle/Beagle," March 6, 2007, www.commtouch.com/documents/Bagle-Worm_MOTR.pdf (URL inactive).
34. Brian Krebs, "Who Wrote the Flashback OS X Worm?" *Krebs on Security*, April 3, 2013, <http://krebsonsecurity.com/2013/04/who-wrote-the-flashback-os-x-worm>
35. Broderick Ian Aquilino, "Flashback OS X Malware," *F-Secure*, October 5, 2012. www.f-secure.com/weblog/archives/Aquilino-VB2012.pdf

NOTES 16 · 17

36. A. Kalafut, A. Acharya, and M. Gupta, "A Study of Malware in PeertoPeer Networks," 2006, www.imconf.net/imc-2006/papers/p33-kalafut.pdf (site discontinued).
37. BO2 k Website, "A Note on Product Legitimacy and Security," ND, http://bo2 k.sourceforge.net/docs/bo2 k_legitimacy.html
38. Costin Raiu, Kurt Baumgartner, and Denis, "Android Trojan Found in Targeted Attack," *SECURELIST*, March 26, 2013, www.securelist.com/en/blog/208194186/Android_Trojan_Found_in_Targeted_Attack
39. Sean Gallagher, "First targeted attack to use Android malware discovered: Kaspersky uncovers trojan spread by 'spear-phish' to Tibet activists." *ars technica*, March 26, 2013, <http://arstechnica.com/security/2013/03/first-targeted-attack-to-use-android-malware-discovered/>
40. Jamie Doward, "UK company's spyware 'used against Bahrain activist', court papers claim: Human rights groups hope email evidence can force review of export controls on surveillance equipment," *The Guardian*, May 11, 2013, www.guardian.co.uk/world/2013/may/12/uk-company-spyware-bahrain-claim
41. S. King and P. Chen, "SubVirt: Implementing Malware with Virtual Machines," 2006. www.eecs.umich.edu/virtual/papers/king06.pdf
42. J. Rutkowska, "Subverting the Vista Kernel for Fun and Profit," 2006, www.blackhat.com/presentations/bh-usa-06/BH-US-06-Rutkowska.pdf
43. A. Wolfe, "U.S. To Review IBM-Lenovo Sale," *CMP*, 2005, www.crn.com/hardware/59100372
44. IBM, "Committee on Foreign Investment in U.S. completes review of Lenovo-IBM deal," *Information Week*, 2005, www.informationweek.com/news/hardware/showarticle.jhtml?articleID=162400445 (URL inactive).
45. P. Bacher, T. Holz, M. Kotter, and G. Wicherksi, "Know Your Enemy: Tracking Botnets," *The Honeynet Project* Website, 2005, www.honeynet.org/papers/bots
46. U.S. Department of Justice, "'Botherder' Dealt Record Prison Sentence for Selling and Spreading Malicious Computer Code," 2006, www.cybercrime.gov/anchetaSent.htm (URL inactive).
47. F. Cohen, "Computer Viruses—Theory and Experiments," 1984, www.all.net/books/virus/index.html
48. M. Christodorescu, S. Jha, S. Seshia, D. Song, and R. Bryant, "Semantics-Aware Malware Detection," 2005, www.eecs.berkeley.edu/~sseshia/pubdir/oakland05.pdf
49. D. Spinellis, "Reliable Identification of Bounded-length Viruses is NP-complete," 2003, www.dmst.aueb.gr/dds/pubs/jrnl/2002-ieeeit-npvirus/html/npvirus.html
50. D. Fisher, "Fast-Moving Worm Crashes Computers," *eWeek* Website, March 22, 2004, www.eweek.com/article2/0,1895,1551998,00.asp
51. T. Gray, "New Worm Attacks Through Symantec Antivirus App," *Newsfactor* Website, December 18, 2006, www.newsfactor.com/news/Worm-Attacks-Through-Antivirus-Flaw/story.xhtml?story_i=d121000E3SRW1 (URL inactive).
52. W. Humphrey, "Three Dimensions of Process Improvement Part II: The Personal Process," 1998, www.stsc.hill.af.mil/crossTalk/frames.asp?uri=1998/03/dimensions.asp
53. S. Manes, "Dim Vista," *Forbes Magazine*, February 26, 2007, www.forbes.com/free_forbes/2007/0226/050.html (URL inactive).

16 · 18 MALICIOUS CODE

54. P. Ferguson and D. Senie, “RFC 2267: Network Ingress Filtering: Defeating Denial of Service Attacks which employ IP Source Address Spoofing,” *Internet Engineering Task Force*, 1998, www.ietf.org/rfc/rfc2267.txt
55. F. Baker and P. Savola, “RFC 3704: Ingress Filtering for Multihomed Networks,” *Internet Engineering Task Force*, 2004, www.ietf.org/rfc/rfc3704.txt
56. National Security Agency, “Network Anomaly Detection Algorithm,” n.d., www.nsa.gov/techtrans/techt00029.cfm (URL inactive).
57. B. Aditya Prakash, Lada Adamic, Theodore Iwashyna, Hanghang Tong, and Christos Faloutsos. “Fractional Immunization in Networks.” *Carnegie Mellon University | Computer Science*, January 24, 2013, www.cs.cmu.edu/~badityap/papers/smrtalloc-sdm13.pdf

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 17

MOBILE CODE

Robert Gezelter

17.1 INTRODUCTION	17·1	17.3.3 Operating Environment Summary	17·12
17.1.1 Mobile Code from the World Wide Web	17·3		
17.1.2 Motivations and Goals	17·5		
17.1.3 Design and Implementation Errors	17·5		
17.2 SIGNED CODE	17·6	17.4 DISCUSSION	17·12
17.2.1 Authenticode	17·6	17.4.1 Asymmetric, and Transitive or Derivative, Trust	17·12
17.2.2 Fundamental Limitations of Signed Code	17·6	17.4.2 Misappropriation and Subversion	17·13
17.2.3 Specific Problems with the ActiveX Security Model	17·7	17.4.3 Multidimensional Threat	17·13
17.2.4 Case Studies	17·8	17.4.4 Client Responsibilities	17·14
		17.4.5 Server Responsibilities	17·15
17.3 RESTRICTED OPERATING ENVIRONMENTS	17·11	17.5 SUMMARY	17·16
17.3.1 Java	17·11	17.6 FURTHER READING	17·16
17.3.2 Virtual Machines	17·12	17.7 NOTES	17·17

17.1 INTRODUCTION. Computing has been transformed in the last twenty years. In these past two decades, access to the World Wide Web has gone from a rare desktop novelty to the ubiquitous; now the majority of people carry mobile devices with full wireless Internet access. The online world has blossomed over the past two decades. The term mobile device is no longer an ironic oxymoron. Twenty years ago, today's dumb phone was the size of a suitcase; today, a far more powerful smart phone easily fits in a child's palm.

Today's portable devices are far more powerful than the mainframes of yesteryear. Their capabilities present a security and information assurance challenge. The consumer-friendly model of easy-to-add enhancements to these devices complicates the security ecosystem.

The World Wide Web began with simple, static pages. Applications soon started to produce pages on demand. The evolution of Web-related technologies has grown to include significant reliance on JavaScript (also known as ECMAScript¹) through the use of AJAX and other techniques, as well as the use of Java and other plug-ins.

17 · 2 MOBILE CODE

The challenge is that such code is often delivered by outside agencies as it is needed (e.g., JavaScript elements are downloaded from the supplier or other Web server as pages are loaded; apps for mobile devices are installed on the fly at little or no nominal cost when desired by the owner of the device). Thus, there is often little or no opportunity to vet the behaviors of such code, or even to maintain a semblance of control. One version of the code may be well-behaved; the next version may enable a compromise of the device, whether inadvertent or deliberate. The Apple App Store used with iPhones®, iPods®, and iPads® has an approval and vetting process, but the process has not been totally without flaw.² As an example, it is not uncommon to reference external JavaScript elements using unencrypted hypertext transfer protocol (HTTP), rendering the Web page inherently vulnerable to masquerading, man-in-the-middle, domain name impersonation, and other attacks. If nothing else, the integrity of an externally supplied element is dependent on the integrity and security of every element in the network, a daunting reality.

At the most basic level, mobile code is a set of instructions that are delivered to a remote computing device for dynamic execution or installation on the device. Such prepackaged components are a boon to applications on desktops and mobile devices; however, the inherent power of the approach comes with significant hazards. The hazards associated mobile code stem from its inherent ability to do more than just display characters on the remote display.

It is this dynamic nature of mobile code that causes policy and implementation difficulties. A blanket prohibition on mobile code is secure, but that prohibition would prevent users of the dynamic Web from performing their tasks and deny users of mobile devices much of the device's utility. The tension among integrity, flexibility, and dynamism is at the heart of the issue.

The advent of mass-market appliances, from personal devices (e.g., smart phones, tablets) to out-of-the-box appliances (e.g., file servers) makes blanket prohibitions often infeasible. Such devices are delivered with prepackaged or downloaded software, which can then be enhanced via the downloading of entire new versions, or extensions known as *apps*. In both cases, the software elements developed by the manufacturer or third parties are prepackaged, with little or no possibility for customer organization testing. Although IT consumerization has benefits, one of the significant drawbacks is the loss of user and corporate visibility into what software is in use, and what access the software has to information.^{3,4,5,6}

Smart phones and tablets are emblematic of IT consumerization. Without the enhancements provided by apps these devices provide a faint shadow of their true capability. However, a rogue or weakly defended app can be downloaded by an unsuspecting user, compromising sensitive data. Apps are supplied as prepackaged, binary installation files, typically delivered over wireless communications pathways (e.g., Wi-Fi or the cellular communications network).

Several definitions, as used by United States military forces but applicable to all, are useful in considering the content of this chapter:

Enclave. An information system environment that is end to end under the control of a single authority and has a uniform security policy, including personnel and physical security. Local and remote elements that access resources within an enclave must satisfy the policy of the enclave.

Mobile code. Software obtained from remote systems outside the enclave boundary, transferred across a network, and then downloaded and executed on a local system without explicit installation or execution by the recipient. Mobile code is a powerful software tool that enhances cross-platform capabilities, sharing of resources, and Web-based solutions. Its use

INTRODUCTION 17 · 3

is widespread and increasing in both commercial and government applications. . . . Mobile code, unfortunately, has the potential to severely degrade . . . operations if improperly used or controlled.

Malicious mobile code. Mobile code software modules designed, employed, distributed, or activated with the intention of compromising the performance or security of information systems, increasing access to those systems, providing the unauthorized disclosure of information, corrupting information, denying service, or stealing resources.⁷

17.1.1 Mobile Code from the World Wide Web. In discussions of the World Wide Web, the phrase *mobile code* generally refers to executable code, other than HTML and related languages (e.g., Extensible Markup Language, or XML), delivered electronically over a communications network (e.g., Web, email) for execution on the client's computer. The most common packaging technologies for mobile code are apps, ActiveX, Java, and JavaScript.

Mobile code can directly perform covert functions on a client system, accessing information or altering the operation or persistent state of the system; it can create accidental or deliberate vulnerabilities that can be exploited at a later time.

The widespread use of mail clients that support executable email content (e.g., HTML email, with either embedded or referenced program code), has become a significant source of vulnerability.

Pop-ups can also be a source of vulnerability in many ways; in particular, they can access other WWW sites and pages and may give rise to log entries that can cause legal problems. The experience of Julie Amero, a Norwich, Connecticut, substitute teacher, is as troubling as it is cautionary. Ms. Amero was accused⁸ and convicted of using a classroom computer to access inappropriate material. The worst charges were subsequently withdrawn when it became clear that the classroom computer at the center of the incident had been infected with malware days before the incident, the antivirus subscription had lapsed, and the testimony of willful access was demonstrated to be factually incorrect. However, Ms. Amero was subjected to an extended prosecution and pilloried in the press. Settling the case required Ms. Amero to plead guilty to a misdemeanor, surrender her teaching license, and pay a fine.^{9,10,11}

Although malicious software such as viruses, worms, and Trojan horse programs written in compiled, interpreted, or scripting languages such as Visual Basic also might be considered mobile code, these pests are not generally labeled as such; this chapter deals only with *apps*, ActiveX controls, Java applets, and JavaScript programs. See Chapter 16 in this *Handbook* for details of other kinds of malicious code.

Today's increasingly dynamic and rich Internet applications rely on AJAX,¹² WebSockets,¹³ and other technologies, which in turn rely on JavaScript and other mobile code technologies. These applications increase the scope of the threat while also making it more difficult to ban the use of the vulnerable technologies.

Although systems or applications crashes are the most spectacular problems with mobile code, they are not the most dangerous. Persistent, silent, and covert access or modification of client system data pose far more serious threats. For example, some Web sites have covertly obtained email addresses from users' browsers, resulting later in unwanted commercial email. Similar attacks have been made against mobile devices phonebooks and contact lists. As the motivations for such attacks have migrated from pranks to criminal and nation-state-sponsored espionage, the importance of undetected covert access to data has correspondingly increased.

In the past, antimalware tactics have relied heavily on widespread distribution of threats. The expected emergence of designer mobile code, specifically targeted to a

17 · 4 MOBILE CODE

particular system or a small number of systems, is a dangerous trend.¹⁴ These targeted attacks often slip through signature scanning systems; widespread distribution and publicity are not part of their intended use.

Fraudulent electronic mail attacks are generally referred to as *phishing*, reportedly derived from the combination of phreaking (as in phone phreaking, bypassing long-distance billing systems) and fishing. With some sense of humor, targeted phishing is referred to as *spearphishing*; attacks against senior managers (e.g., CEOs) are referred to as *whaling* or *harpooning*. The nomenclature may be humorous, but well-executed harpooning attacks can be difficult to detect or prevent. Such attacks have been found at the center of several important compromises (e.g., the 2011 spear phishing attack against RSA¹⁵) at major corporations. High-value targets are not only senior executives; individuals with access to or authority over high-value information are equally at risk.

One of the earliest documented episodes of this type is the 2004 use of Trojan horse software, discussed in §2.10.5 in this *Handbook*, for extensive data theft. This case was not an isolated incident. The trend of targeted attacks against senior personnel, rather than random attacks, has accelerated.¹⁶

Investigative agencies have also entered the fray. In July 2007, an affidavit filed by the Federal Bureau of Investigation in connection with a series of bomb threats described the use of spyware to infiltrate a suspect's computer and return information to investigators.¹⁷

Mobile code presents a complex set of issues to information technology (IT) professionals. Allowing mobile code into an operating environment compromises any enclave; however, even commercial off-the-shelf (COTS) programs breach the integrity of an enclave. The differences between mobile code and other forms of code are primarily in the way these external programs are obtained, installed, documented, and controlled. In an enterprise computing environment, COTS acquisition normally involves conscious and explicit evaluation of the costs and benefits of installing a particular named and documented program. In contrast, mobile code programs are installed largely without notification to the user, and generally without documentation, change control, or review of any kind. Unless client-system firewalls are set to reject mobile code automatically (which is increasingly difficult, if not a technical and practical impossibility), system administrators cannot be certain exactly which software has been executed on client machines under their nominal control. The use of Secure Sockets Layer (SSL)-based technologies such as HTTPS to otherwise secure connections used by applications also has the side effect of preventing the detection of mobile code at the firewall level. Although such control is often illusory, due to user circumvention of restrictions on installation of unauthorized software, the use of mobile code, installed by external Web sites, seriously compromises any remaining control over employee software configurations.

Mobile code has also been used in some cases to enforce proprietary rights in content, as was the case in a 2005 affair involving Sony Music.¹⁸ The Sony Music case involved software contained on music CD-ROMs; precisely the same effect could have occurred with a downloaded file or Web page. The covert installation of software of any kind is a serious hazard to integrity, security, and privacy. Malfunction or misuse of such software would likely fit within the criminal statutes defining illegal, unauthorized alteration of systems. The attorney general of Texas,¹⁹ as well as private class actions in New York²⁰ and California,²¹ all filed cases against Sony. All of these actions were settled by Sony BMG in December 2006.²² In January 2007, the U.S. Federal

INTRODUCTION 17 · 5

Trade Commission announced a settlement with Sony BMG on the charges of installing software without permission.²³ Investigations were also opened by the attorneys general of Massachusetts and Florida, as well as overseas in Italy and Canada.

There are also reports that the Sony Root Kit was exploited by the Backdoor.IRC.Snyd.A exploit²⁴ and others to hide files from malware scans. The widespread nature of this induced vulnerability should also give pause. A widespread vulnerability provides an ecological niche ready for exploitation.

17.1.2 Motivations and Goals. The motivations and goals of malware propagators have continued an evolutionary trend from the unintentionally destructive to the vengeful, vindictive, and criminal. Unsurprisingly, the actors have migrated from rogue individual pranksters to criminals and nation states and proxies.

In the beginning, many incidents were randomly damaging, or pranks with unintended side effects. This is no longer the case. Now malevolent mobile code is often code with a purpose. That purpose may be embarrassment, it may be blackmail, it may be corporate espionage, or it may be out-and-out theft. In a different dimension, the goal may be the subordination of otherwise innocent computer resources for a criminal enterprise against unrelated third parties.

The change in goals also has a dramatic impact on counter strategies. When the goal was mass publicity, the same infection was widespread, and scanning technologies could be used to identify known threats. When the goal is no longer publicity, publicity and widespread infection are maladaptive. Covert infection is then a far more attractive strategy than mass distribution. Custom mobile code designed to achieve selective covert infections is unlikely to quickly appear in the crosshairs of scanning software. This follows the evolutionary trajectory common in the biological world, where pathogens tend to mutate into less fatal forms over time. It is maladaptive for a parasite, which is what most malware is, to fatally damage its host. The downside of this effect is that the chronic infections with no apparent side effects are often overlooked.

Going forward, technologies and operational routines that make it difficult for unauthorized code to take up residence in or compromise the persistent state of the system are far more desirable counterstrategies than approaches based on scanning for known infections.

17.1.3 Design and Implementation Errors. Design and implementation errors take a variety of forms. The simplest cases involve software that malfunctions on a constant predictable basis. More pernicious and more dangerous are those errors that silently compromise the strict containment of a multiuser environment. Errors in such prophylactic layers, known as brick walls or sandboxes, compromise the integrity of the protection scheme. In the worst cases, they permit unfettered access to system-level resources by unprivileged user programs.

Design and implementation errors can occur within any program or procedure, and mobile code is no exception. *Sandboxes* (nonprivileged, restricted operating environments) are intended to prevent unauthorized operations. *Authentication* determines which organization takes responsibility for such errors.

This chapter looks at a security model based on authentication of mobile code and then examines how restricted operating environments help to limit damage from harmful mobile code.

17 · 6 MOBILE CODE

These concerns are appropriate to both widely distributed and targeted attacks. The challenge of targeted attacks lies in their small population; targeted attacks are unlikely to appear on the radar of general distribution scanning programs.

17.2 SIGNED CODE. Authentication technologies are designed to ensure that information supplied by an organization has not been altered without authorization. The technology used to implement this assurance is based on use of the Public Key Infrastructure (PKI), as discussed in detail in Chapter 37 in this *Handbook*. Code authenticated using PKI-based mechanisms often is referred to as *signed*. Signed code generally is immune to unauthorized modification; however, a signature guarantees integrity only from the point in time that the code is signed; the signing process does not imply safety or quality of the code prior to the point of signing.

Once signed, a file cannot be altered without invalidating the original signature; without the cooperation of someone holding access to the private key associated with the creating organization's X.509 certificate, the signature cannot realistically be tweaked to look legitimate. (An X.509 certificate, digitally signed by an authorized user, authenticates the binding between a user's name and the user's public key.) However, looking below the surface, such precautions do not address a variety of vulnerabilities:

- Unauthorized access to private (signing) keys
- Access to the code base prior to signing
- Fraudulent certificates
- Design and implementation errors

17.2.1 Authenticode. Microsoft's Authenticode technology is an example of an authentication-based approach.²⁵ Developers wishing to distribute code obtain an appropriate digital certificate from a Certification Authority (CA) and use the digital certificate to sign the code. The signature is then checked by the client system each time that the code is executed.

Authenticode relies on several components:

- PKI and the X.509 certificates issued by a Certification Authority.
- Limited access to the private keys associated with the issuing organization's X.509 certificate. In Microsoft terminology, the term Software Publishing Certificate or SPC refers to a PKCS #7 object, which in turn contains a collection of X.509 certificates used to sign code.
- The integrity of the processes used by the CA to ensure that requests for X.509 certificates are legitimate.

Authenticode does not address issues relating to the safety or accuracy of signed code, merely that it is authentic and unaltered since signing. For example, signing does not provide any guard against employee malfeasance.

17.2.2 Fundamental Limitations of Signed Code. Signing technologies, regardless of the context (e.g., email, applets, and archives), do not directly address questions of accuracy or correctness; they merely address questions of origin. The

SIGNED CODE 17 · 7

biggest danger in signing schemes is the all-or-nothing approach taken to trust. Signed items are presumed to be trustworthy to the fullest extent of the requesting user's authority. The signed item can perform any operation that the user would be permitted to execute. There is no concept of *partial trust*. In an attorney's words, such an acceptance would be a *general power of attorney*. In the words of the CERT/Coordination Center (CERT/CC)-sponsored *Security in ActiveX Workshop*: "A digital signature does not, however, provide any guarantee of benevolence or competence."²⁶

At the same time, the inherent power and apparent legitimacy of a digital signature place a heavy burden on signers and the higher levels of the PKI to ensure the integrity of the mechanisms and secrets.

The key to the integrity of signed code is the signing process and the process that generates the object to be signed; the security of the secret keys required for its implementation determines the degree of trust in attribution of the signed code. In the truest sense, the private keys associated with the X.509 certificate represent the keys to the kingdom, as valuable as a signature chop in the Far East or a facsimile signature plate for a bank account.

On a practical level, accepting code signed by an organization is an explicit acceptance that the signing organization has good controls on the use of its signing keys. Organizations that take security seriously, segregating access to privileged accounts and controlling access to systems, are well positioned to manage the procedures for signing code.

Thus, the procedures and systems used for signing code should be treated with the same caution as is used for the aforementioned signing plates or the maximum security cryptographic facilities familiar to those in the national security area.

Unfortunately, despite years of publicity about the dangers of shared passwords and accounts, many IT installations continue to use shared accounts and passwords. There is little reason to assume that the secrets relating to PKI are better protected, despite extensive recommendations that those details be well guarded.

17.2.3 Specific Problems with the ActiveX Security Model. The CERT/CC workshop on Security in ActiveX summarized the security issues in three major areas: importing and installing controls, running controls, and the use of controls by scripts.²⁷ The key findings from this report are as follows:

- Importing and Installing Controls
 - As discussed, the sole basis for trusting a signed control is its presumed origin. However, the originator of the code may have incorporated a design flaw in the control or may not have done adequate software quality assurance to prevent serious bugs.
 - A trusting user may install a signed control that contains a vulnerability making it useful for attackers simply because it is signed.
 - On Windows systems with multiple users, once a control has been permitted by one user, it remains available for all users, even if their security stances differ.
- Running Controls
 - An ActiveX control has no limitations on what it can do on the client machine, and it runs with the same privileges as those of the user process that initiated the control.

17 · 8 MOBILE CODE

- Although ActiveX security measures are available in Internet Explorer, other client software may run controls without necessarily implementing such security. Internet Explorer security levels tend to be all or nothing, making it difficult to allow a specific control without allowing all controls of that type. Remote activation of controls can bypass normal security perimeters such as those imposed by firewalls.
- There is no basis for deciding whether a particular control is safe to execute or not, in any particular context.
- Scripting Concerns
 - Lacking a general basis for limiting the actions of controls, ActiveX programmers must effectively determine their own precautions to prevent harmful actions. It is difficult enough to develop a good set of boundaries on program activity, even if one uses a general model such as the sandbox described later; it is extremely difficult to see how individual developers can be expected or, more importantly, trusted to create their own equivalent of the sandbox for each individual control. In light of these hazards, the authors of the CERT/CC report stated that “there is a large number of potential failure points.”

17.2.4 Case Studies. Several security breaches or demonstrations mediated through ActiveX have occurred since the introduction of this technology in the mid-1990s.

17.2.4.1 Internet Exploder. In 1996, Fred McLain wrote Internet Exploder, an ActiveX control designed to illustrate the broad degree of trust conferred on an ActiveX control by virtue of its having been signed. Exploder, when downloaded for execution by Internet Explorer, will shut down the browser’s computer (the equivalent of the *Shut down | Shut down* sequence from the Start menu on a Windows system). This operation is operationally disruptive but does not actually corrupt the system. McLain notes in his frequently asked questions (FAQ) on Exploder that it is easy to build destructive or malicious controls.²⁸

Exploder raises an important question: Who and what are the limits on trust when using signed code? In normal commercial matters, there is a large difference between an inauthentic signature, a forgery, and a properly signed but unpayable check. In software, the difference between an inauthentic control and a dangerous one is far less clear.

17.2.4.2 Chaos Computer Club Demonstration. On January 27, 1997, a German television program showed members of the Chaos Computer Club demonstrating how they could use an ActiveX control to steal money from a bank account. The control, available on the Web, was written to subvert the popular accounting package Quicken using what amounts to an electronic version of tailgating. A victim need merely visit a site and download the ActiveX control in question; the control then automatically checked to see if Quicken was installed. If so, the control ordered Quicken to issue a transfer order to be saved in its list of pending transfers. The next time the victim connected to the appropriate bank and sent all pending transfer orders to the bank, all the transfers would be executed as a single transaction. The user’s personal identification number (PIN) and transaction authorization number (TAN) would apply to all the transfers, including the fraudulent one in the pile of orders. Most

SIGNED CODE 17 · 9

victims would be unaware of the theft until they received their next statement—if then.²⁹

Dan Wallach of Princeton University, commenting on this case, wrote:

When you accept an ActiveX control, you're allowing completely arbitrary code to rummage around your machine and do anything it pleases. That same code could make extremely expensive phone calls, to 900 numbers or over long distances, with your modem; it can read, write, and delete any file on your computer; it can install Trojan horses and viruses. All without any of the subterfuge and hackery required to do it with Java. ActiveX hands away the keys to your computer.³⁰

Responding to criticisms of the ActiveX security model, Bob Atkinson, architect and primary implementer of Authenticode, wrote a lengthy essay explaining his point of view. Among the key points:

- Microsoft never claimed that it would certify the safety of other people's code.
- Authentication is designed solely to permit identification of the culprits *after* malicious code is detected.
- Explorer-based distribution of software is no more risky than conventional purchases through software retailers.³¹

Subsequent correspondence in the RISKS Forum chastised Mr. Atkinson for omitting several other key points, such as:

- Interactions among ActiveX controls can violate system security even though individual controls appear harmless.
- There is no precedent in fact for laying liability at the feet of software developers even when you can find them.
- Under attack, evidence of digital signature is likely to evaporate from the system being damaged.
- Latency of execution of harmful payloads will complicate identification of the source of damage.
- Malice is not as important a threat from code as incompetence.
- Microsoft has a history of including security-threatening options, such as automatic execution of macros in Word, without offering any way of turning off the feature.
- A Web site can invoke an ActiveX control that is located on a different site or that already has been downloaded from another site, and can pass, by means of that control, unexpected arguments that could cause harm.³²

Wallach's criticisms of ActiveX are generally applicable to signed code technologies, not peculiar to Microsoft's ActiveX.

17.2.4.3 Certificates Obtained by Imposters. In January 2001, VeriSign issued two *Class 3 Digital Certificates* for signing ActiveX controls and other code to someone impersonating a Microsoft employee. As a result, users receiving code signed using these certificates would receive a request for acceptance or rejection of a certificate apparently signed by Microsoft on January 30 or 31, 2001. As Russ Cooper

17 · 10 MOBILE CODE

commented on the NTBUGTRAQ Usenet group when the news came out in March 2001:

The fact that unless you actually check the date on the Certificate you won't know whether or not its [sic] one you can trust is a Bad Thing(tm)[sic], as obviously not everyone (read: next to nobody) is going to check every Certificate they get presented with.

You gotta wonder how VeriSign's issuance mechanism could be so poorly designed and/or implemented to let something like this happen.

Meanwhile, Microsoft are [sic] working on a patch that will stick its finger in this dam.

Basically, VeriSign Code-Signing Certificates do not employ a Certificate Revocation List (CRL) feature called CDP, or CRL Distribution Point, which causes the Certificate to be checked for revocation each time it's read. Even if you have CRL turned on in IE, VeriSign Code-Signing Certificates aren't checked.

Microsoft's update is going to shim in some mechanism which causes some/all Code-Signing Certificates to check some local file/registry key for a CRL, which will (at least initially) contain the details of these Certificates. Assuming this works as advertised, any attempt to trust the mis-issued Certificates should fail.³³

Roger Thompson, Chief Technical Officer for Exploit Prevention Labs, explained that the imposters' motives would determine how bad the results would be from the fraudulent certificates. "If it was someone with a purpose in mind, then six weeks is a long time to do something," he said. "If the job was to install a sniffer, then there could be a zillion backdoors as a result of it." Published reports indicated that the failure of authentication occurred due to a flaw in the issuing process at VeriSign: The certificates were issued *before* receiving verification by email that the official customer contact authorized the certificates. This case was the first detected failure of authentication in over 500,000 certificates issued by VeriSign.³⁴

Some recent cases involving stolen or forged certificates include the following:

- In November 2011, antimalware vendor F-Secure's Mikko Hypponen reported "finding a piece of malicious software that was cryptographically signed by a forged Adobe certificate originating with Government of Malaysia: Malaysian Agricultural Research and Development Institute ..."³⁵
- In September 2012, Adobe reported that they "... received two malicious utilities that appeared to be digitally signed using a valid Adobe code signing certificate."³⁶ The company explained, "Sophisticated threat actors use malicious utilities like the signed samples during highly targeted attacks for privilege escalation and lateral movement within an environment following an initial machine compromise."
- In February 2013, Bit9 Systems revealed that "Due to an operational oversight within Bit9, we failed to install our own product on a handful of computers within our network. As a result, a malicious third party was able to illegally gain temporary access to one of our digital code-signing certificates that they then used to illegitimately sign malware. There is no indication that this was the result of an issue with our product. Our investigation also shows that our product was not compromised."³⁷

17.2.4.4 Certificate Forgery. Certificate forgery is another hazard. Certificates are not magical; they are merely documented cryptographically signed by a recognized Certification Authority. The cryptographic techniques used to sign certificates

RESTRICTED OPERATING ENVIRONMENTS 17 · 11

are subject to attack. A research result demonstrating the hazard of MD-5 collisions leading to forged certificates was originally published in 2008.³⁸ For this reason, many authorities have recommended that MD-5 signed certificates be withdrawn from use. In August 2011, an MD-5 signed fraudulent certificate for *.google.com was reported by Iranian Internet users.³⁹

17.3 RESTRICTED OPERATING ENVIRONMENTS. From a Web perspective, the term *sandbox* defines what could be referred to as a restricted operating environment. Restricted operating environments are not new; they have existed for nearly 50 years in the form of multiuser operating systems, including MULTICS, OS/360 and its descendants, OpenVMS, UNIX, and others.

In simple terms, a restricted, or nonprivileged, operating environment prohibits normal users and their programs from executing operations that can compromise the overall system. In such an environment, normal users are prohibited from executing operations such as HALT that directly affect hardware. User programs are prevented from executing instructions that can compromise the operating system memory allocation and processor state and from accessing or modifying files belonging to the operating system or to other users. Implemented and managed carefully, such systems are highly effective at protecting information and data from unauthorized modification and access. The National Computer Security Center (NCSC) *Orange Book* contains criteria for classifying and evaluating trusted systems.⁴⁰

The strengths and weaknesses of protected systems are well understood. Permitting ordinary users unrestricted access to system files compromises the integrity of the system. Privileged users (i.e., those with legitimate access to system files and physical hardware) must be careful that the programs they run do not compromise the operating system. Most protected systems contain a collection of freestanding programs that implement useful system functions requiring some form of privilege to operate. Often these programs have been the source of security vulnerabilities. This is the underlying reasoning behind the universal recommendation that programs not run as root or Administrator or with similar privileges unless absolutely necessary.

17.3.1 Java. Java is a language developed by Sun Microsystems for platform independent execution of code, typically within the context of a Web browser. The basic Java environment includes a Java Virtual Machine (JVM) and a set of supporting software referred to as the Java Run Time Environment. Applets downloaded via the World Wide Web (intranet or Internet) have strict limitations on their ability to access system resources. In particular, these restrictions prevent the execution of external commands and read or write access to files.

The Java environment does provide for *signed applets* that are permitted wider access to files. Dynamically downloaded applets also are restricted to initiating connections to the system that supplied them, theoretically limiting some types of third-party attacks.

In concept, the Java approach, which also includes other validity tests on the JVM pseudocode, should be adequate to ensure security. However, the collection of trusted applets found locally on the client system and signed downloaded applets represent ways in which the security system can be subverted. Without signature, the Java approach is also vulnerable to attack by Domain Name System (DNS) spoofing.

Multiuser protection and virtual machine protection schemes also are totally dependent on the integrity of the code that separates the nonprivileged users from privileged, system-compromising, operations. Java has not been an exception to this rule. In 1996,

17 · 12 MOBILE CODE

a bug in the Java environment contained in Netscape Navigator Version 2.0 permitted connections to arbitrary Universal Resource Locators (URLs).⁴¹ Later, in 2000, errors were discovered in the code that protected various resources.⁴² Although the Java environment is less widely exploitable than ActiveX, vulnerabilities continue to be uncovered. In 2007, at least two vulnerabilities were reported by US-CERT.⁴³ Significantly, both of these reported vulnerabilities involved the ability of untrusted applets to compromise the security envelope.

Additionally, since unsigned code can take advantage of errors in underlying signed code, there is no guarantee that complex combinations of untrusted and trusted code will not lead to security compromises.

17.3.2 Virtual Machines. Isolation can also be realized using virtual machines. Theoretically, each virtual machine runs its own copy of an operating system on a separate virtual copy of the hardware. Properly implemented, mainframe systems have run multiple, independent machines on a single physical system for decades. Virtual machine technology is a foundational technology for so-called cloud computing. In such configurations, the Achilles' heel is the layer that separates the different virtual machines. If there is a deficiency in that layer, or in its interfaces, the compartmentalization may be subject to compromise.

In 2012, just such an incident occurred with multiple virtual machine implementations implementing virtual machines on Intel's x86 processors.⁴⁴ In these cases, a specially crafted stack frame could lead to a break of the virtual machine's containment mechanisms. Virtual machines from Xen, FreeBSD, Red Hat, and Microsoft Windows were subject to this attack.

17.3.3 Operating Environment Summary. The details of virtualization are significant in this context. Errors in virtualization software implementation can create security compromises in the separation between different virtual machines, or indeed it may compromise the underlying integrity of the host. Such hazards are not theoretical; they occur in application-level sandboxes and in systems-level virtual machines.

17.4 DISCUSSION. Mobile code security raises important issues about how to handle relationships in an increasingly interconnected computing environment.

17.4.1 Asymmetric, and Transitive or Derivative, Trust. It is common for cyberrelationships to be *asymmetric* with regard to the size or power of the parties. This fact increases the potential for catastrophic interactions. It also creates opportunities for mass infection across organization boundaries. Large or critical organizations often can unilaterally impose limitations on the ability of partner organizations to defend their information infrastructure against damage.

The combination of a powerful organization and insufficient controls on signing authority, or, alternatively, the obligatory execution of unsigned (or self-signed) ActiveX controls, is a recipe for serious problems. The powerful organization is able to obligate its partners to accept a low security level, such as would result, for example, from using unsigned ActiveX controls, while abdicating responsibility for the resulting repercussions.

All organizations should, for security and performance reasons, use the technology that requires the least degree of privilege to accomplish the desired result. JavaScript/ECMAScript can provide many functions, without the need for the

DISCUSSION 17 · 13

functionality provided by Java, much less ActiveX. It remains common for large organizations to force the download of ActiveX controls for purposes that do not require the power of ActiveX, merely using the justification that they perceive Internet Explorer to be the more prevalent browser. Often this requires running the installation script from an account with Administrator privileges, a second security violation. This is particularly surprising since these same organizations often offer parallel support for Firefox, Opera, Safari, and other non-ActiveX supporting browsers on Linux, Apple, and other platforms. This Trust Me concept forces the risk and burden of consequences on the end user, who is far less able to deal with the consequences.

As noted earlier in this chapter, Web servers represent an attractive vector for attacks. Signing (authentication) methods are a way to control damage potential, if the mechanisms used for admitting executable code are properly controlled. Failure to control these mechanisms leads to severe side effects.

In Chapter 30 in this *Handbook*, it is noted that protecting Web servers requires that the server contents be managed with care. It is appropriate and often necessary to isolate Web servers on separate network segments, separated from both the Internet and the organizational intranet by firewalls. These precautions are even more necessary when servers are responsible for supplying executable code to clients.

Security practitioners should carefully examine the different functions performed by each server. In some cases, such as OpenVMS hosts, where network servers commonly run as unprivileged processes in separate contexts and directory trees, it is feasible to run multiple services on a single server. In other systems, such as UNIX and Windows, where it is common for applications services to execute as privileged, with full access to all system files, a logic error in a network service can compromise the security of the entire server, including the collection of downloadable applets.

Far more serious and equally subtle is *transitive* (or derivative) trust: Alpha trusts Beta who trusts Gamma. A security compromise—for example, an unsigned Java applet or a malfunctioning or malevolent ActiveX control supplied by Gamma—compromises Beta. Beta then causes problems with Alpha’s systems. This cascade can continue repeatedly, leading to numerous compromised Web services and systems far removed geographically and organizationally from the original incident.

17.4.2 Misappropriation and Subversion. The threat space has mutated over the last several years. Where the main danger from mobile code was attacks on the target machine, today’s threat is far more diverse. In November 2007, John Schiefer of Los Angeles pled guilty to installing software designed to capture usernames and passwords. According to news reports, he was also involved in running a number of networks of compromised computers, often referred to as bots, which are often used to initiate distributed denial-of-service (DDoS) and other attacks.⁴⁵ In this particular case, the announcement by the U.S. Department of Justice⁴⁶ mentions two specific episodes: 250,000 machines infected with spybots to obtain user’s usernames and passwords for PayPal and other systems; and a separate scheme involving a Dutch Internet advertising company in which a network of 150,000 infected computers were used to sign up for one of the advertising company’s programs.

This was one of the cases stemming from Bot Roast II,⁴⁷ an FBI operation against several botnet networks.

17.4.3 Multidimensional Threat. Mobile code is a multidimensional threat, with several different aspects that must each be treated separately. Signing code, such as Java applets or ActiveX controls, addresses the problem of authenticity and authority

17 · 14 MOBILE CODE

to release the code. However, the integrity of the signature mechanism requires that the integrity of the PKI infrastructure be beyond reproach. In a very real sense, the PKI infrastructure is beyond the control of the organization itself. Any compromise or procedural slip on the part of the Certificate Authority or signer invalidates the presumptions of safety.

Signing, however much it contributes to resolving the question of authenticity, does not address safety or validity. As an example, the Windows Update ActiveX control, distributed by Microsoft as part of the various Windows operating systems, has as its underlying purpose the update of the operating system. A failure of that control would be catastrophic. Fortunately, Microsoft gives users the choice of using the automatic update facility or doing updates manually. Many Web applications are not so accommodating.

The problem is not solely a question of malfunctioning applets. It is possible that a collection of applets involved in a client's overall business activities may collide in some unanticipated fashion, from attempting to use the same Windows registry key in contradictory ways, to inadvertently using the same temporary file name. Similar problems often occur with applications that presume they have a monopoly on the use of the system, an all-too-common syndrome.

These issues are, for the most part, completely unrelated to each other. A solution in one area would neither improve nor worsen the situation with regard to the other issues.

17.4.4 Client Responsibilities. The expanding threat presents a challenge to those responsible for ensuring the integrity of desktop computing. Put simply, there is a complex, multidimensional threat, and it is not easily defended against using the techniques of portals, firewalls, and scanners.

The danger from browsing the World Wide Web is the danger that the browser will permit an attacker, directly or indirectly, to cause a modification to the persistent state of the system. The simplest step in the correct direction is not to browse the World Wide Web from within a protection context that has access to critical system files and settings. Limiting this access by using a nonprivileged user account for browsing significantly decreases the hazard, provided of course that the system files are protected from access by such an account.

The mass availability of *virtual machine* technology presents an additional alternative. Virtual machine technology, pioneered by IBM in the 1960s on mainframes, has emerged in a new guise on platforms down to the desktop level. The general availability of this capability in the desktop world opens up a whole new defensive strategy against mobile malware: the expendable Web browser.

An expendable WWW browser is an instantiated desktop within a virtual machine environment, from a known system image. If it is compromised, it is merely rewritten from a known, uncompromised system image. It allows one to create a low-security, at-risk, browsing enclave within an otherwise higher security environment. This is an approach that has been used, in a physical sense to be sure, by some organizations providing public access personal computers. Rather than attempting to fortify the machines against compromise or attack, they are reinitialized from a known image after each user. This allows the end user to indulge the foibles of trading partners' attempts to impose unsafe computing practices in an expendable environment that can be isolated. Using Windows as an example, while it is an unsafe practice to install software as Administrator, it is far less damaging to do so in a virtual machine, where the machine can be deleted at convenience with little side effect.

DISCUSSION 17 · 15

Similarly, *disposable virtual machines*^{48,49} can be used by normal users to separate the normal operating environment from one where there is a perceived increase in risk exposure (e.g., browsing suspect sites; forced installation of ActiveX controls of uncertain provenance or safety).

17.4.5 Server Responsibilities. As noted earlier in this chapter, Web servers represent an attractive vector for attacks. Signing (authentication) methods are a way to control damage potential, provided the mechanisms used for admitting executable code are properly controlled. Failure to control these mechanisms leads to severe side effects.

The concept of *minimum necessary privilege* applies to mobile code. There is little reason to impose the use of ActiveX for the purposes of changing the color of a banner advertisement. This effect can often be accomplished using Cascading Style Sheets (CSS). JavaScript/ECMAScript is capable of many powerful, display-related operations with a high degree of safety in situations where CSS is not adequate. Using Java to maintain a shopping cart (price, quantities, and contents) is reasonable and does not require the use of a signed applet, with its attendant greater capabilities and risks. At the other end of the scale, it is plausible that a system update function (e.g., the Windows Update function, which automatically downloads and installs changes to the Windows operating system) requires the unbridled power of a signed ActiveX control.

When the power of signed applets or controls is required, good software engineering practice provides excellent examples of how to limit the potential for damage and mischief, as discussed in Chapter 38 in this *Handbook*.

These problems extend beyond an individual organization's direct exposure. An increasing trend over the past several years has been to use third-party Web servers as an infection vector, a springboard for infecting Website visitors. If components are downloaded from third-party servers, the integrity of your WWW site may be vulnerable to compromise.

Good software implementation isolates functions and limits the scope of operations that require privileged access or operations. Payroll applications do not directly manipulate printer ports, video display cards, network adapters, or disk drives. Privileged operating system components, such as device drivers and file systems, are responsible for the actual operation. This separation, together with careful parameter checking by the operating system kernel and the privileged components, ensures safety.

The same techniques can be used with applets and controls. Because they require more access, they should be programmed carefully, using the same defensive measures as are used when implementing privileged additions to operating systems. As an example, there is little reason for a Simple Mail Transfer Protocol (SMTP) server to be privileged. An SMTP server requires privileges for a single function, the delivery of an individual email message to a recipient's mailbox. This can be accomplished in two ways:

1. Implement the application in a nonprivileged way, by marking users' email files and directories with the necessary access permissions for the mail delivery program to create and modify email files and directories. Such a mechanism is fully in conformance with the NCSC *Orange Book*'s C2-level of security.
2. Implement a separate subcomponent whose sole responsibility is the actual message delivery of the message. The subcomponent must be written defensively to check all of its parameters, and does not provide an interface for the execution of arbitrary code. This approach is used by HP's OpenVMS operating system.

17 · 16 MOBILE CODE

The UNIX *sendmail* program, by contrast, is a large, multifunctional program that executes with privileges. *sendmail* has been the subject of numerous security problems for over a decade and has spawned efforts to produce more secure replacements.⁵⁰

17.5 SUMMARY. Mobile code provides many flexible and useful capabilities. The different mechanisms for implementing mobile code range from the innocuous (HTML), to fairly safe (JavaScript/ECMAScript), and with increasing degrees of power and risk through Java and ActiveX.

Ensuring security and integrity with the use of mobile code requires cooperation on the part of both the provider and the client. Clients should not accept random signed code and controls. Providers have a positive responsibility to:

- Follow good software engineering practices
- Grant minimum necessary privileges and access
- Use defensive programming
- Limit privileged access, with no open-ended interfaces
- Ensure the integrity of the signing process and the associated private keys

With appropriate caution, mobile code can be a constructive, powerful part of intranet and Internet applications, both within an organization and in cooperation with its customers and other stakeholders.

17.6 FURTHER READING

- Carl, Jeremy. "ActiveX Security: Under the Microscope," *Web Week*, 2, No. 17, November 4, 1996; www.Webdeveloper.com/activex/activex_security.html
- CERT. "NIMDA Worm," September 11, 2001, www.cert.org/advisories/CA-2001-26.html
- CERT. "sadmind/IIS Worm," May 8, 2001, www.cert.org/advisories/CA-2001-11.html
- CERT. "Unauthenticated 'Microsoft Corporation' Certificates," March 22, 2001, www.cert.org/advisories/CA-2001-04.html
- Dormann, Will, and Jason Rafail. "Securing Your Web Browser," CERT, January 23, 2006, www.cert.org/tech_tips/securing_browser/index.html on March 19, 2013.
- Evers, J. FAQ: JavaScript Insecurities. CNET July 28, 2006 <http://news.cnet.com/2100-7349_3-6100019.html>
- Felten, Edward. "Security Tradeoffs: Java vs. ActiveX," last modified April 28, 1997, www.cs.princeton.edu/sip/faq/java-vs-activex.html
- Felten, E., and J. Halderman (2006, February 14). *Lessons from the Sony CD DRM Episode*, Center for Information Technology Policy, Department of Computer Science, Princeton University, USENIX Security 2006 http://static.usenix.org/events/sec06/tech/full_papers/halderman.pdf
- Felten, E., and G. McGraw. *Securing Java: Getting Down to Business with Mobile Code*. New York: John Wiley & Sons, 1999. Also free and unlimited Web access from www.securingjava.com
- Gehtland, J., B. Galbraith, and D. Almaer. *Pragmatic Ajax*. Raleigh, NC: Pragmatic Bookshelf, 2006.
- Grossman, J., and T. C. Niedzialkowski. "Hacking Intranet Websites from the Outside," *Black Hat (USA)*, Las Vegas, August 3, 2006. www.blackhat.com/presentations/bh-usa-06/BH-US-06-Grossman.pdf

NOTES 17 · 17

- Hensing, R. "W32/HLLP.Philis.bq, Chinese Gold Farmers and What You Can Do about It," December 2, 2007, http://blogs.technet.com/robert_hensing/archive/2006/12/04/w32-hllp-philis-bq-chinese-gold-farmers-and-what-you-can-do-about-it.aspx on
- Holzman, S. *Ajax Bible*. Hoboken, NJ: John Wiley & Sons, 2007.
- Keizer, G. "FBI Planted Spyware on Teen's PC to Trace Bomb Threats," *PCWorld*, July 19, 2007. www.pcworld.com/article/134823/article.html
- McGraw, G., and E. W. Felten. *Securing Java: Getting Down to Business with Mobile Code*, 2nd Edition. John Wiley & Sons, 1999.
- McGraw, Gary & Ed Felten. "Java Security Web Site." www.digital.com/javasecurity/index.html
- Microsoft. "Introduction to Code Signing" (with appendix), 2001, <http://msdn.microsoft.com/en-us/library/ms537361%28v=vs.85%29.aspx>
- Rhoads, C. "Web Scammer Targets Senior U.S. Executives," *Wall Street Journal*, November 9, 2007.
- Schwartz, J. "iPhone Flaw Lets Hackers Take Over Security Firm Says," *New York Times*, July 23, 2007.
- VeriSign, "Microsoft Security Bulletin MS01-017: Erroneous VeriSign-Issued Digital Certificates Pose Spoofing Hazard," March 22, 2001, www.microsoft.com/TechNet/security/bulletin/MS01-017.asp
- VeriSign, "VeriSign Security Alert Fraud Detected in Authenticode Code Signing Certificates," March 22, 2001.
- Zakas, N., J. McPeak, and J. Fawcett. *Practical Ajax*, 2nd ed. Hoboken, NJ: John Wiley & Sons, 2007.

17.7 NOTES

1. ECMA-262 retrieved from www.ecma-international.org/publications/files/ECMA-ST-ARCH/ECMA-262%205th%20edition%20December%202009.pdf
2. Carter Dotson, "Is Apple Losing Its Grip on App Store Security?" Mashable Web site, March 1, 2012, <http://mashable.com/2012/03/01/app-store-security-risks>
3. Scott Thurm and Yukari Iwatani Kane, "Your Apps Are Watching You," *The Wall Street Journal*, December 17, 2010, <http://online.wsj.com/article/SB10001424052748704694004576020083703574602.html>
4. J. Schwartz, "iPhone Flaw Lets Hackers Take over, Security Firm says," *The New York Times*, July 23, 2007, www.nytimes.com/2007/07/23/technology/23iphone.html
5. CBC News (2013, January 28) "WhatsApp breaches privacy laws," *CBC News*, January 28 2013, www.cbc.ca/news/technology/story/2013/01/28/technology-whatsapp-privacy.html
6. Office of Privacy Commissioner of Canada, "Report of Findings Investigation into the personal information handling practices of WhatsApp Inc.," January 15, 2013, www.priv.gc.ca/cf-dc/2013/2013_001_0115_e.asp
7. Memorandum by Arthur M. Money, Assistant Secretary of Defense for C3I, and CIO to Secretaries of the Military Departments, Chairman of the Joint Chiefs of Staff, Chief Information Officers of the Defense Agencies, et al., "Use of Mobile Code Technologies in DoD Information Systems," October 23, 2006, www.dtic.mil/whs/directives/corres/pdf/855201p.pdf

17 · 18 MOBILE CODE

8. J. Penny, “40 Years Too Long in Norwich Porn Case?” *Norwich Bulletin*, January 9, 2007; G. Smith, “Teacher Facing Porn Charges,” *Norwich Bulletin*, November 11, 2004.
9. Robert Cringley “The Julie Amero Case: A Dangerous Farce,” *PCWorld*, December 2, 2008, http://pcworld.com/article/154768/julie_amero.html
10. Robert Cringley, “Julie Amero Case, Part II: May the Farce Be With You,” *PC-World*, December 4, 2008, www.pcworld.com/article/154937/julie_amero.html
11. Robert McMillan, “How Spyware Nearly Sent a Teacher to Prison,” *PC-World*, November 26, 2008, www.pcworld.com/article/154611/spyware_teacher_prison.html
12. Jesse James Garrett, “AJAX: A New Approach to Web Applications,” *Adaptive Path*, February 18, 2005, <http://adaptivepath.com/ideas/essays/archives/000385.php>
13. Ian Fette and Alexey Melnikov, “The WebSocket Protocol,” RFC 6455, December 2011, www.rfc-editor.org/rfc/rfc6455.txt
14. D. Izenberg, “Trojan Horse Developers Indicted,” *Jerusalem Post*, March 5, 2006.
15. Sean Gallagher, “RSA Details March cyberattack, Blames ‘Nation State’ for SecurID Breach,” *ars technica*, October 12, 2011, <http://arstechnica.com/business/2011/10/rsa-details-march-cyber-attack-blames-nation-state-for-securid-breach>
16. J. Kirk, “Hackers Target C-level Execs and Their Families,” *Network World*, July 2, 2007, www.networkworld.com/news/2007/070207-hackers-target-c-level-execs-and.html
17. G. Keizer, “FAQ: What We Know (Now) about the FBI’s CIPAV Spyware,” *Computerworld*, July 29, 2007, www.computerworld.com/s/article/9028298/FAQ_What_we_know_now_about_the_FBI_s_CIPAV_spyware
18. Mark Russinovich, “Sony, Rootkits and Digital Rights Management Gone Too Far,” Mark Russinovich’s Blog, December 2, 2007, <http://blogs.technet.com/markrussinovich/archive/2005/10/31/sony-rootkits-and-digital-rights-management-gone-too-far.aspx>
19. *The State of Texas v. Sony BMG Music Entertainment, LLC*, Case GV-505065, District Court of Travis County, Texas, 126th Judicial District.
20. *James Michaelson and Ori Edelstein v. Sony BMG Music, Inc. and First 4 Internet*, Case 05 CV 9575, United States District Court, Southern District of New York.
21. *Alexander William Guevara v. Sony Music Entertainment, et al.*, Case BC342359, Superior Court of the State of California, County of Los Angeles.
22. R. McMillan, “Sony Pays \$ 1.5M to Settle Texas, California Root Kit Suits,” *Computerworld*, December 20, 2006.
23. Federal Trade Commission, “Sony BMG Settles FTC Charges,” January 30, 2007, www.ftc.gov/opa/2007/01/sony.htm
24. Sorin Dudea, “Backdoor.IRC.Snyd.A,” Bitdefender Toolbox Web site, December 2, 2007, www.bitdefender.com/VIRUS-1000058-en-Backdoor-IRC-Snyd-A.html
25. Roger Grimes, “Authenticode,” Microsoft TechNet Web site, n.d., <http://technet.microsoft.com/en-us/library/cc750035.aspx>
26. CERT/CC, *Results of the Security in ActiveX Workshop, Pittsburgh, Pennsylvania, August 22–23, 2000*, www.cert.org/archive/pdf/activeX_report.pdf
27. CERT/CC, pp. 6–9.

NOTES 17 · 19

28. F. McLain, "The Exploder Control Frequently Asked Questions (FAQ)," January 30, 1997, updated February 7, 1997, www.halcyon.com/mclain/ActiveX/Exploder/FAQ.htm
29. D. Weber-Wulff, "Electronic Funds Transfer without Stealing PIN/TAN," *RISKS* 18, No. 80 (1997), <http://catless.ncl.ac.uk/Risks/18.80.html>
30. D. Wallach, "RE: Electronic Funds Transfer without Stealing PIN/TAN," *RISKS* 18, No. 81 (1997), <http://catless.ncl.ac.uk/Risks/18.81.html>
31. B. Atkinson, "Comments and Corrections Regarding Authentication," *RISKS* 18, No. 85 (1997), <http://catless.ncl.ac.uk/Risks/18.85.html>
32. *RISKS* 18, No. 86 (1997), et seq. <http://catless.ncl.ac.uk/Risks/18.86.html>
33. Microsoft, "Erroneous VeriSign-Issued Digital Certificates Pose Spoofing Hazard," Microsoft Security Bulletin MS01-017, <http://support.microsoft.com/kb/293818>
34. R. Lemos, "Microsoft Warns of Hijacked Certificates," *ZDNet News*, March 22, 2001, http://news.cnet.com/Microsoft-warns-of-hijacked-certificates/2100-1001_3-254586.html
35. Cory Doctorow, "Stolen Government of Malaysia Certificate Used to Sign Malicious Fake Adobe Software Update," *BoingBoing.net* Web site, November 14, 2011, <http://boingboing.net/2011/11/14/stolen-government-of-malaysia.html>
36. Brad Arkin, "Inappropriate Use of Adobe Code Signing Certificate," Adobe Secure Software Engineering Team (ASSET) Blog, September 27, 2012, <http://blogs.adobe.com/asset/2012/09/inappropriate-use-of-adobe-code-signing-certificate.html>
37. Patrick Morley, "Bit9 and Our Customers' Security," Bit9 Website, February 8, 2013, <https://blog.bit9.com/2013/02/08/bit9-and-our-customers-security/>
38. Alexander Sotirov, "Creating a Rogue CA Certificate" *Security Research* Web site, December 30, 2008, www.phreedom.org/research/rogue-ca
39. Dan Goodin, "Fraudulent Google Credential Found in the Wild: Did Counterfeit SSL Cert Target Iranians?" *The Register*, August 29, 2011, www.theregister.co.uk/2011/08/29/fraudulent_google_ssl_certificate/
40. For full text, see <http://csrc.nist.gov/publications/secpubs/rainbow/std001.txt>
41. CERT, "Java Implementations Can Allow Connection to an Arbitrary Host," CERT Advisory, March 5, 1996, updated September 24, 1997, www.cert.org/advisories/CA-1996-05.html
42. CERT, "Netscape Allows Java Applets to Read Protected Resources," CERT Advisory, August 10, 2000, www.cert.org/advisories/CA-2000-15.html
43. CERT, "Sun Java JRE Vulnerable to Unauthorized Network Access," Vulnerability Note VU#336105, October 5, 2007, revised October 12, 2007, www.kb.cert.org/vuls/id/336105; and "Sun Java JRE Vulnerable to Privilege Escalation," Vulnerability Note VU#102289, January 9, 2007, revised May 16, 2007, www.kb.cert.org/vuls/id/102289
44. CERT, "SYSRET 64-Bit Operating System Privilege Escalation Vulnerability on Intel CPU Hardware," Vulnerability Note VU#649219, June 12, 2012, revised September 4, 2012, www.kb.cert.org/vuls/id/649219
45. J. Serjeant, "'Botmaster' Admits Infecting 250,000 Computers," *Reuters*, November 9, 2007.
46. United States Attorney's Office, Central District of California (2007, November 9) "Computer Security Consultant Charges with Infecting up to a Quarter

17 · 20 MOBILE CODE

Million Computers that Were Used to Wiretap, Engage in Identity Theft, Defraud Banks,” Press Release No. 07-143, November 9, 2007, www.justice.gov/criminal/cybercrime/press-releases/2007/schieferCharge.pdf

47. FBI, “‘Bot Roast II’: Cracking Down on CyberCrime,” November 29, 2007, www.fbi.gov/news/stories/2007/november/botnet_112907; more extensive details in “‘Bot Roast II’ Nets 8 Individuals,” November 29, 2007. www.fbi.gov/news/pressrel/press-releases/bot-roast-ii-nets-8-individuals
48. Robert Gezelter, “Disposable Virtual Machines: Deliberately Expendable” Ruminations—An IT Blog, August 23, 2010, www.rlgsc.com/blog/ruminations/disposable-virtual-machines.html
49. Robert Gezelter, “Disposable Virtual Machines,” 2010 Trenton Computer Festival, Ewing, New Jersey, April 24, 2010, www.rlgsc.com/trentoncomputerfestival/2010/disposable-virtual-machines.html
50. For example, the National Vulnerability Database (NVD, <http://nvd.nist.gov>) shows that the Common Vulnerabilities and Exposures (CVE) database includes a total of 29 unique vulnerabilities involving sendmail, of which 15 are dated 2000 and 2001. This trend continues, with the NVD Version 2.0 showing an additional 16 sendmail-related issues from 2002 through 2007.

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 18

DENIAL-OF-SERVICE ATTACKS

Gary C. Kessler

18.1	INTRODUCTION	18·2	18.7	DISTRIBUTED DENIAL-OF-SERVICE ATTACKS	18·14
18.2	HISTORY OF DENIAL-OF-SERVICE ATTACKS	18·2	18.7.1	Short History of Distributed Denial of Service	18·15
18.3	COSTS OF DENIAL-OF-SERVICE ATTACKS	18·4	18.7.2	Distributed Denial-of-Service Terminology and Overview	18·16
18.4	TYPES OF DENIAL-OF-SERVICE ATTACKS	18·6	18.7.3	Distributed Denial-of-Service Tool Descriptions	18·19
18.5	SPECIFIC DENIAL-OF-SERVICE ATTACKS	18·6	18.8	DENIAL-OF-SERVICE USING EXPLOITABLE SOFTWARE	18·23
18.5.1	Destructive Devices	18·6	18.8.1	Code Red	18·23
18.5.2	Email (and Email Subscription)		18.8.2	NIMDA	18·24
18.5.3	Bombing	18·6	18.9	DEFENSES AGAINST DISTRIBUTED DENIAL-OF-SERVICE ATTACKS	18·25
18.5.4	Buffer Overflow	18·8	18.9.1	User and System Administrator Actions	18·25
18.5.5	Bandwidth Consumption	18·8	18.9.2	Local Network Actions	18·26
18.5.6	Routing and Domain Name System Attacks	18·10	18.9.3	Internet Service Provider Actions	18·27
18.5.7	SYN Flooding	18·11	18.9.4	Exploited Software Defensive Actions	18·28
18.5.8	Resource Starvation	18·12	18.9.5	Other Tools under Development or Consideration	18·28
18.5.9	Java, PHP, and ASP Attacks	18·13	18.10	MANAGEMENT ISSUES	18·30
18.5.10	Router Attacks	18·13	18.11	FURTHER READING	18·31
	Other Denial-of-Service Attacks	18·13	18.12	NOTES	18·32
18.6	PREVENTING AND RESPONDING TO DENIAL-OF-SERVICE ATTACKS	18·14			

18 · 2 DENIAL-OF-SERVICE ATTACKS

18.1 INTRODUCTION. This chapter discusses denial-of-service (DoS) and distributed denial-of-service (DDoS) attacks. These attacks seek to render target systems and networks unusable or inaccessible by saturating resources or causing catastrophic errors that halt processes or entire systems. Furthermore, they are increasingly easy for even *script kiddies* (persons who follow explicit attack instructions or execute attack programs) to launch. Successful prevention of and defense against these attacks will come only when there is widespread cooperation among all Internet service providers (ISPs) and other Internet-connected systems worldwide and when antimalware software is universally installed and kept up to date on user systems.

Working in a variety of ways, the DoS attacker selects an intended target system and launches a concentrated attack against it. Although initially deemed to be primarily a nuisance, DoS attacks can incapacitate an entire network, especially those with hosts that rely on Transmission Control Protocol/Internet Protocol (TCP/IP). DoS attacks on corporate networks and ISPs have resulted in significant damage to productivity and revenues. DoS attacks can be launched against any hardware or operating system platform because they generally aim at the heart of Internet Protocol (IP) implementations. Because IP is the typical target, the DoS attack tools that run under one operating system (Linux is a common choice) can be aimed at any operating system running IP. Additionally, because IP implementations are similar for different platforms, one DoS attack may target several operating systems and work on each. Once written for one platform and released, new DoS attacks appear to evolve (via the examination and participation of hackers and crackers) so that in a short period of time (approximately two weeks) mutations of the DoS attack appear that work on virtually all platforms.

Because of the critical impact that DoS attacks can have, they cannot be taken lightly. Targeted DoS attacks have been around in one form or another since the 1980s; in 1999, they evolved into DDoS attacks, primarily due to the heavy use of internal networks and the Internet. DDoS tools launch coordinated DoS attacks from many sources against one or more targets simultaneously.

This chapter describes first DoS and then DDoS attacks because DoS attacks historically predate DDoS attacks. In addition, some DDoS attacks make use of DoS techniques. However, the attacks themselves, the terminology, and the defenses are sufficiently different that they warrant separate discussion.

18.2 HISTORY OF DENIAL-OF-SERVICE ATTACKS. Historically, any act that prevented use of a system could be called a denial of service. For example, in some mainframe and minicomputer systems, typing a malformed command could cause a system failure; for example, on the HP3000 minicomputer around 1982, using a faulty TERMTYPE parameter in the logon string could instantly crash the system. On the HP3000, typing one or more characters on the system console without pressing the RETURN key would block all further system messages on the console, causing system buffers to fill up with undisplayed messages; when the system buffers (16 by default) were full, no system action requiring notification to the console could take place (e.g., logons, logoffs, special-form requests, or tape-mount requests). A classic DoS without software could be the result of having a cleaner unplug a power cord from a wall socket without realizing that the cord was powering one or more workstations. Configuring a permanent lockout of users after a certain number of bad-password attempts (instead of a timed lockout) led to DoS; malicious users could even try to lockout root users by deliberately entering bad passwords on their accounts. However, such events were usually the result of bugs in the operating system, inadequate numbers of system resources, poor system management, or accident.

HISTORY OF DENIAL-OF-SERVICE ATTACKS 18 · 3

One of the first major DoS incidents that made headlines and ripples was probably accidental. It took place in December 1987, when an employee of IBM in Europe sent out a holiday greeting email message. The email message, however, contained a program that drew a nice Christmas tree on the recipient's terminal—and read the recipient's NAMES and NETLOG files listing the user's address book as well as the email addresses of individuals that this user had recently sent mail to or received mail from. The Christmas Tree program then triggered the automatic transmission of copies of itself to all email addresses found in NAMES and NETLOG—and failed to prevent remailing to systems that had already received or sent the greeting. The resulting mailstorm overloaded the IBM corporate network worldwide and brought it crashing down both in Europe and the United States. In addition, messages escaped from IBM's corporate network and wreaked havoc on the BITNET/EARN education and research networks in North America and Europe. Although the cause of the outage was originally thought to be a computer virus, and the incident has been referred to as a virus, a worm, and a prank, the result was a true denial of service.

Perhaps the most famous Internet DoS resulted from the Morris Worm, also referred to as the Internet Worm, which was released in November 1988. Cornell graduate student Robert T. Morris had cowritten an article about vulnerabilities in *Sendmail* and *fingerd* on UNIX systems during the previous year. Most of the TCP/IP community dismissed the vulnerabilities as theoretical; apparently Morris wanted to demonstrate that the vulnerabilities actually could be practically exploited. To demonstrate the problem, his program had to invade another system, guess some passwords, and then replicate itself. Morris said after the incident that he wanted the program to replicate just a few times to demonstrate that it was real, and he included code to prevent rapid spread; unfortunately, a programming error led the worm to replicate often and quickly, in addition to superinfecting already-infected systems. The worm clogged the Internet with hundreds of thousands of messages and effectively brought the entire network down; most sites that did not actually crash were disconnected from the network by their system administrators to avoid being infected and to allow disinfection. Regardless of intent, the Morris Worm inadvertently caused a DoS. Because sites administrators did not have a way to communicate with other admins, and because solutions that Morris himself posted could not be trusted, there was a lot of time spent in eradicating the worm. One result was the creation of Computer Emergency Response Team Coordination Center, now CERT/CC, at Carnegie Mellon University.

On Friday, September 6, 1996, PANIX, a public access Internet provider in Manhattan, was struck by a DoS attack that consisted of many messages flooding the server with massive amounts of data. The attacks were made against mail, news, name, and Web servers as well as user shell account machines. The attackers successfully eluded tracking and the attacks went on for several days. About a week after the attacks began, PANIX went public with the story, and dozens of other online service providers acknowledged that they too were the victims of similar attacks.

In 1997, a disgruntled former employee of Forbes, Inc., used a former colleague's password to remotely access the firm's computer systems. Deliberately tampering with the system, he deleted budgets and salary information and caused five of the eight network servers to crash. When Federal Bureau of Investigation (FBI) agents caught the perpetrator, his home was filled with hacking tools, proprietary information from Forbes, Inc., and other incriminating material.

In February 1998, hackers from Israel and Northern California attacked the U.S. Department of Defense (DoD). In a carefully organized attack that exploited buffer

18 · 4 DENIAL-OF-SERVICE ATTACKS

overflow, these hackers systematically perpetrated a DoS attack that lasted over a week on 11 DoD sites.

In March 1998, all across the United States, system administrators found their Windows NT servers under apparently automated attack. Systems crashed repeatedly until they were updated to the latest patches from Microsoft. There appeared to be no file damage from the attacks, which lasted more than a day. Sites affected included several NASA and other military sites, as well as several University of California and other college campuses.

Yet another mail storm erupted in May 1998 when an Australian official set autoreply on his email package while he was away. Unfortunately, he inadvertently set his destination for these largely useless messages to be all 2,000 users on his network—and requested auto confirmation of delivery of each autoreply, which generated yet another autoreply, and so on *ad infinitum*. Within four hours, his endless positive-feedback loop had generated 150,000 messages before his autoreply was shut down. The ripples lasted for days, with the perpetrator saddled with 48,000 messages in his IN basket and a stream of 1,500 a day pouring in. This was another case of an inadvertent DoS attack.

In January 1999, someone launched a sustained DoS attack on Ozemail, an important Australian Internet service provider. Email service was disrupted for users in Sydney.

In March 1999, the Melissa email-enabled virus/worm swept the world in a few days as it sent copies of itself to the first 50 addresses in victims' email address books. Because of this high replication rate, the virus spread faster than any previous virus in history. On many corporate systems, the rapid rate of internal replication saturated email servers with outbound automated junk email. Initial estimates were in the range of 100,000 downed systems. Antivirus companies rallied immediately, and updates for all the standard products were available within hours of the first notices from the CERT/CC. The Melissa macro virus was quickly followed by the PAPA MS-Excel macro virus with similar properties but that, in addition, launched DoS attacks on two specific IP addresses.

And after all of the publicity dealing with the Y2 K problems prior to the turn of the millennium, the ILOVEYOU, or Love Letter, worm was released in May 2000. The Love Letter worm was in a Visual BASIC attachment to emails with the subject line "ILOVEYOU." This worm used a so-called double extension exploit; the attachment was named LOVE-LETTER-FOR-YOU.txt.vbs, but because the default Windows setting is to not display the extension, many (or most) recipients thought that this was a harmless text file. The worm modified the infected system's registry to autostart the program on system startup, replaced files with certain extensions with copies of itself, and sent copies of itself to the first 50 entries in the Outlook address book.

More recent attacks are discussed in the sections describing modern DoS tools.

18.3 COSTS OF DENIAL-OF-SERVICE ATTACKS. What are the effects of these DoS attacks in terms of productivity and actual financial costs? It is difficult to place an exact monetary figure on DoS attacks. DoS attacks can interrupt critical processes in an organization, and such interruption can be costly. When a company's computer network is inaccessible to legitimate users and they cannot conduct their normal business, productivity is lowered. The negative effect is bound to carry over to the financial aspects of the business. However, putting exact figures to these effects is uncertain at best, and the estimates are widely disputed, even among security and business experts. In addition, many companies do not comment on the exact losses they suffer because they fear that the negative publicity will decrease their market share. This latter point is significant: In an early 1990s study of Wall Street firms, some of

COSTS OF DENIAL-OF-SERVICE ATTACKS 18 · 5

the companies suggested that if they were to be without their network for two to three days, they might never reopen their doors.

In the case of the Christmas Tree worm, it took IBM several days to clean up its network and resulted in the loss of millions of dollars, both for cleansing the system and for lost business because of lost connectivity and related productivity. Additionally, there was the embarrassment suffered by IBM, a noted technology company. The individual who launched the worm was identified and denied access to any computer account, while IBM had to write a letter of apology to the European Academic and Research Network (EARN) administrators.

In 1988, at the time of the Morris Worm, the Internet consisted of 5,000 to 10,000 hosts, primarily at research and academic institutions. As a result, although the Morris Worm succeeded in bringing many sites to a halt and gained worldwide notoriety, the financial and productivity impact on the commercial world was minimal. A similar incident today would wreak havoc and cost tens or hundreds of millions of dollars in losses.

By 1996, however, commercial reliance on the Internet was already becoming a matter of course. Working around the clock, the management at PANIX and Cisco, the ISP's router vendor, kept the service provider up and running, but the network received 210 fraudulent requests per minute. Although the systems did not crash, thousands of subscribers were unable to receive their email messages. Other sites were attacked in the same time frame as PANIX, including Voters Telecommunication Watch. No one took responsibility for these attacks and it has been widely assumed that they were triggered by articles on SYN DoS attacks (see the description below) that had recently appeared in *2600 Magazine* and *Phrack*, journals that cater to hackers.

According to Forbes, Inc., the losses suffered by the firm because of the DoS attack perpetrated by the disgruntled former employee exceeded \$100,000. Could it have been prevented? According to the firm, it is highly unlikely, since Forbes had no reason to suspect the individual was maintaining either the firm's confidential and sensitive material at home or that he was thinking of hacking into the computer system and deliberately doing damage. Although the firm had security on its systems, the perpetrator used the password of a legitimate, authorized user.

The DoS attack launched against the Department of Defense computers in 1998 proved that attackers could deny access to vital military information. In this particular instance, the attack was directed at unclassified machines that had only administrative and accounting records, but it was a blow to the confidence of the Department of Defense. Tying up the computers for over a week presumably reduced productivity, but the government would not comment on the actual cost from loss of machine time and personnel productivity.

These cases show that a DoS attack on a computer or network can be devastating to an organization. Important equipment and networks and even an entire organization can be disabled by such attacks.

Early DoS incidents often were described as annoying, frustrating, or a nuisance. However, with increasing sophistication and dependency on networking, it has become difficult to keep a sense of humor about such incidents. Especially in corporations, where the mission is to make a profit for the shareholder, company managers find it increasingly difficult to excuse being incapacitated because of a DoS or DDoS attack.

As these forms of attack become more sophisticated, so must the tools and methods for detecting and fighting them. Current products scan equipment and networks for vulnerabilities, trigger alerts when an abnormality is found, and frequently assist in eliminating the discovered problem.

18 · 6 DENIAL-OF-SERVICE ATTACKS

18.4 TYPES OF DENIAL-OF-SERVICE ATTACKS. DoS attacks, whether accidental or deliberate, result in loss of service; either a host or a server system is rendered inoperable or a network is rendered inaccessible. DoS attacks are launched deliberately by an *intruder* (the preferred term for attacker in this context). Systems and networks that are compromised are referred to as the *victims*. And while DoS attacks can be launched from the intruder's system, they often are launched by an automated process that allows the intruder to start the attack remotely with a few keystrokes. These programs are known as *daemons*, and they are often placed on another system that the hacker has already compromised.

There are four basic types or categories of DoS attack:

1. **Saturation.** This type of attack seeks to deprive computers and networks of scarce, limited, or nonrenewable resources that are essential in order for the computers or networks to operate. Resources of this type include CPU time, disk space, memory, data structures, network bandwidth, access to other networks and computers, and environmental resources such as cool air and power.
2. **Misconfiguration.** This type of attack destroys or alters configuration information in host systems, servers, or routers. Because poor or improperly configured computers may fail to operate or operate inadequately, this type of attack can be very severe.
3. **Destruction.** This type of attack results in network components being physically destroyed or altered. To guard against this type of attack, it is necessary to have good physical security to safeguard the computers and other network components. This chapter does not deal with physical damage, but Chapters 22 and 23 of this *Handbook* do so in some detail.
4. **Disruption.** This attack interrupts the communications between two devices by altering state information—such as the state of a TCP virtual connection—such that effective data transfer is impossible.

18.5 SPECIFIC DENIAL-OF-SERVICE ATTACKS. The discussion below describes some specific DoS attacks as examples of the general methods that can be employed in a DoS or DDoS. The list is not exhaustive.

18.5.1 Destructive Devices. Destructive devices are programs that accomplish either harassment or destruction of data. There are mixed opinions regarding how severe destructive devices are, but if they threaten a computer's or network's ability to function properly and efficiently, then they may be the instruments of DoS attacks. Viruses, email bombs, and DoS tools all can be considered destructive devices. In fact, viruses and email bombs are known to cause DoS attacks. Some viruses or other malware are known to attack control systems for hardware, such as those employed the Stuxnet and Flame attacks. Viruses and other malicious software and their actions are covered in Chapter 16 of this *Handbook*; nonreplicating or *traditional* DoS and DDoS tools are discussed here.

18.5.2 Email (and Email Subscription) Bombing. Email and email subscription *bombings* were among the first documented DoS attacks. An *email bomb* consists of large numbers of email messages that fill up a victim's electronic mailbox. A huge number of messages can tie up an online connection, slow down mail delivery, and even overload the email server system until the system crashes. Most email bombings

SPECIFIC DENIAL-OF-SERVICE ATTACKS 18 · 7

are thought to be deliberate attacks by disgruntled people; specific targets may be victims of someone with a particular grudge. For example, a San Francisco stockbroker received 25,000 messages on September 23, 1996, consisting of the word “Idiot” from a consultant with whom he had had a disagreement. The flood of messages prevented him from using his computer, so in December the victim sued the perpetrator’s employer for \$25,000 of damages. On occasions in the past, such as with the Christmas Tree worm and the Internet Worm, the DoS is thought to have been accidental.

Email bomb packages automate the process of launching and carrying out an *email bombing* DoS attack. With names like Up Yours, Kaboom, Avalanche, Gatemail, and the Unabomber, these packages can be placed on a network server during a DoS attack and used to attack other systems. Administrators who are aware of these names and others should regularly scan their drives for associated filenames and eliminate them.

To safeguard computers and/or servers, mail filters and exclusionary schemes can automatically filter and reject mail sent from a source address using email bomb packages. Mail filters are available for UNIX, Windows, Macintosh, and Linux systems. Most computer operating systems and most ISPs now offer filtering tools for eliminating unsolicited commercial email and other email. Although perpetrators often disguise their identity and location by using a false address, most filters can be set to screen and eliminate these addresses.

With *email subscription bombing*, also known as *list linking*, a user is subscribed to dozens of mailing lists by the attacker without the user’s knowledge. For example, someone calling himself “johnny xchaotic” perpetrated one of the earliest subscription-bombing incidents. In August 1996, he claimed the blame for a massive mail-bombing run based on fraudulently subscribing dozens of victims to hundreds of mailing lists. In a rambling and incoherent letter posted on the Internet, he made rude remarks about famous and not-so-famous people whose capacity to receive meaningful email was obliterated by up to thousands of unwanted messages a day. Today, filtering packages have point-and-click mechanisms that provide automatic list linking. A user could, conceivably, start to receive hundreds or thousands of mail messages per day if linked to just 50 to 100 lists. Once linked to the various lists, the victim has to manually unsubscribe from each individual mailing list. If an attack takes place while the victim is away and without access to email, a user could have a backlog of thousands of messages by the time he or she returns.

List server software should never accept subscriptions without sending a request for confirmation to the supposed subscriber; in the current environment, most list servers require confirmation from the supposed subscriber before initiating regular mailings. The default for the confirmation messages is to ignore the announcement if it is fraudulent. However, even this safety mechanism can generate a wave of many single email messages if a mail bomber abuses the list servers. As a result, many sites have incorporated Completely Automated Public Turing test to tell Computers and Humans Apart (CAPTCHA) images to interfere with bots that attempt to subscribe victims to their mailings.

Speaking of being away, vacation messages and return receipts are another way in which individuals can inadvertently start an *email storm*. Many users set their email clients to automatically request a return receipt of all messages sent. Then the users go on vacation and set up an auto-reply vacation message. When they get a message, the client sends back the vacation message and also requests a receipt. The returned receipts, in turn, generate more auto-reply vacation messages.

Another variant on this feedback loop occurs when an employee goes on vacation and forwards all email to an external ISP that has a local access number in the locale

18 · 8 DENIAL-OF-SERVICE ATTACKS

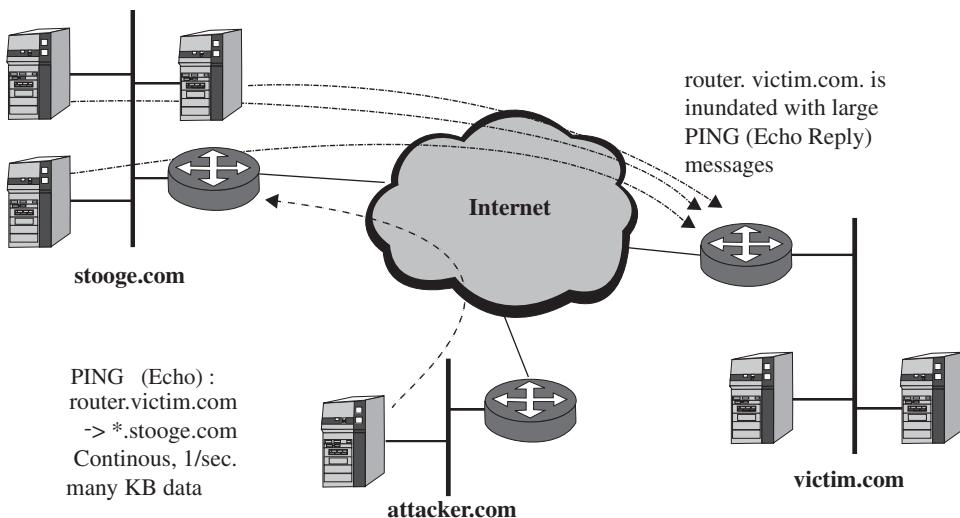
of the vacation. If the employee decides not to check the mail while away, either due to a high local access fee or so as not to interfere with the vacation, the ISP mailbox will fill up with forwarded messages. If the mailbox fills up, the ISP will send a bounce message back to the corporate server—which then forwards the bounce message back to the ISP, which generates yet another bounce message. Eventually, even the corporate mail server will fill up with a single individual’s messages, causing an email DoS.

18.5.3 Buffer Overflow. *Buffer overflow* attacks can be insidious and damaging. It is possible to send an input string to a target program that contains actual code and is long enough to overflow the memory space or input buffer. Sometimes this surreptitious code is placed on the process stack (the area in a computer’s memory where the operating system keeps track of the program’s input and related code used for processing the inputs), and the code then is processed. An overflow can occur when the input data overflows its buffer space and flows into the stack, where it overwrites the previous data and return address. If the program is written so that the stack address points to the malicious code located in the return buffer, the code executes with the original program’s privileges. Buffer overflow is the result of poor programming, where the programmer does not check the size of the input compared to the input buffer. Although the bad-programming basis of buffer overflows should have been eradicated by now, new buffer overflow attacks pop up monthly. As of January 2013, the National Vulnerability Database¹ included 6,273 buffer overflows out of a total of 64,398 vulnerabilities—about 10 percent. This percentage has been dropping over the years; in the period from January 2010 through December 2012, 1,281 of the 14,510 records were for buffer overflows—only 9 percent; in contrast, of the 39,822 records for vulnerabilities entered before 2010, 5,047 (13 percent) involved buffer overflows.

Not all buffer overflows allow the user to insert executable code. DoS attacks such as the *ping of Death* merely attach a data block that is larger than allowed by the IP protocol (i.e., greater than 65,536 bytes). Because the packets are broken into fragments for transmission, they manage to get through the network and, probably, the router and firewall. Once reassembled at the target, however, the packets cause the IP kernel’s buffer to overflow and, if not managed properly, the system crashes.

Another example is an old flaw in Microsoft’s Internet Information Services² (IIS) that could be exploited to allow the Web service to be halted. To do this, an attacker would request a document with a very long URL from an IIS-based Website (and how does one identify an IIS site? If a Website uses pages with the .htm or .asp extensions, it is a good guess that the site is running IIS). Upon receipt of the request, an access violation occurred and the server would halt. Although Microsoft issued a patch for this vulnerability, successful attacks continued to take place for years.

18.5.4 Bandwidth Consumption. *Bandwidth consumption* involves generating a large number of packets directed to the network under attack. Such attacks can take place on a local network or can be perpetrated remotely. If the attacker has, or can access, greater bandwidth than the victim has available, the attacker can flood the victim’s network connection. Such saturation can happen to both high-speed and low-speed network connections. Although any type of packet may be used, the most common are Internet Control Message Protocol (ICMP) *echo* messages (generated by *pinging*). By engaging multiple sites to flood the victim’s network connection, attackers can amplify the DoS attack. To do this successfully, the attackers convince the amplifying system to send traffic to the victim’s network. Tracing the intruder who

SPECIFIC DENIAL-OF-SERVICE ATTACKS 18 · 9**EXHIBIT 18.1** smurf DoS Attack

perpetrates a bandwidth consumption attack can be difficult since attackers can spoof their source addresses.

One very common bandwidth consumption attack is the *smurf attack*. This attack is particularly interesting (and clever) since it uses tools that reside on every IP system and employs a third-party site without actually having to take control of any system anywhere. As Exhibit 18.1 shows, the smurf attack starts when the intruder (at *attacker.com*) sends a continuous string of very large *ping* messages to the IP broadcast address (* in the exhibit, representing the all-ones, or broadcast, IP address) of the unsuspecting third party, *stooge.com*. The intruder spoofs the source IP address of the *ping* message so that it appears that these messages come from, say, the router at the target network, *router.victim.com*. If the intruder sends a single 10,000-byte *ping* message to the broadcast address of an intermediate site with 50 hosts, for example, the responses will consume 4 megabits (Mb).³ Even if *victim.com* has a T1 line (with a bandwidth of 1.544 Mb per second), the attacker can swamp a victim's line merely by sending a single large *ping* per second to the right stooge site. The intermediate site is called, for obvious reasons, the *amplifying network*; note that the attacker does not need to compromise any systems there. The ratio of originally transmitted packets to the number of systems that respond is known as the *amplification ratio*.

A variant of the smurf attack is called *fraggle*. In this variant, the attackers send spoofed User Datagram Protocol (UDP) packets instead of *echo* messages to the broadcast address of the amplifying network. Each system on the amplifying network that has the specific broadcast address port enabled will create large amount of traffic by responding to the victim's host; if the port is not enabled, the system on the amplifying network will generate ICMP *host unreachable* messages to the victim's host. In either case, the victim's bandwidth is consumed.

Kernel panic attacks are not due to programming error, per se, but to a hole in a program's logic. As an example, the Intel Pentium chip that could not correctly divide two particular legal inputs had a programming flaw. An IP kernel that fails when receiving a packet that should never occur in the first place is, indeed, a gap in the program's logic but not the same as failing to handle legal input. In this case, the program did not know how to fail gracefully when it received unexpected input.

18 · 10 DENIAL-OF-SERVICE ATTACKS

A specific example of kernel panic occurs in Linux kernel v.2.2.0 when a program usually used for printing shared library dependencies is used instead to print some core files. Under certain circumstances, *munmap()*, a function call used to map and unmap devices into memory, overwrites critical areas of kernel memory and causes the system to panic and reboot.

And there are other examples of these kinds of attacks. A *land attack* occurs when a spoofed packet is sent to a target host where the TCP source and destination ports are set to the same value, and the IP source and destination addresses are set to the same value. Since this is confusing to the host operating system, it results in 100 percent CPU utilization and then a halt. Land attacks have been directed at just about all operating systems.

Teardrop attacks are also the result of behavior when receiving impossible packets. If an IP packet is too large for a particular network to handle, the packet will be fragmented into smaller pieces. Information in the IP packet header tells the destination host how to reassemble the packet. In a Teardrop attack, the attacker deliberately crafts IP packets that appear to overlap when reassembled. This, too, can cause the host to crash. Teardrop attacks have been targeted against Microsoft operating systems and all variants of UNIX.

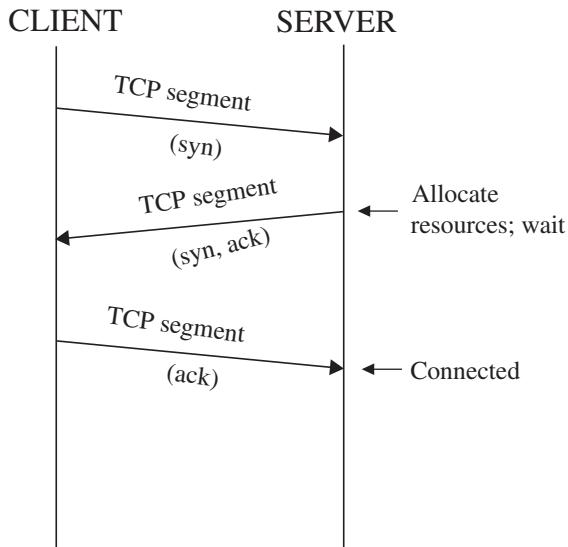
18.5.5 Routing and Domain Name System Attacks. *Routing and Domain Name System (DNS) attacks* are clever attacks that are achieved repeatedly. By tampering with the DNS, a site's domain name resolves to the IP address of any other site that the attacker wishes. In August 1999, for example, a disgruntled engineer for PairGain Technologies provided a direct link to a bogus but authentic-looking Bloomberg News Service page he had created. The engineer, who owned PairGain shares, posted false information about PairGain's supposed acquisition by an Israeli company, pumping up the price of PairGain stock worldwide and creating havoc in the market through his page-jacking exploits. Although a nontraditional DoS attack, these actions did deny users access to the site that they desired, with serious consequences.

Another example occurred in July 1997, when Eugene Kashpureff filed fraudulent information with InterNIC for its DNS updates for that month, forcing domain name servers around the globe to recognize temporary and unauthorized Internet addresses ending in .xxx, .mall, .nic, and .per. A few weeks later, Kashpureff inserted false information that forced people trying to access the Website of Network Solutions Inc. to end up at the AlterNIC site, which he owned. He was arrested later that year on charges of wire fraud.

In another example—this from 2000—of a routing and DNS attack, RSA Security, Inc., after announcing that it had developed a method to combat Website hackers, found that users were unwittingly being rerouted to a counterfeit RSA Website. The fraudulent site looked exactly like the original RSA Website but made fun of the fact that the hacker had managed to achieve his DoS goal.

As far back as 1997, weaknesses were found in and documented about BIND implementations in versions preceding 4.9.5+P1. The earlier versions would cache bogus DNS information when DNS recursion was enabled. The vulnerability enabled attackers to exploit the process of mapping IP addresses to hostnames in what is known as *pointer (PTR) record spoofing*. This type of spoofing provides the potential for a DNS DoS attack.

DNS attacks are still seen widely today. *Phishing* is a form of social engineering attack whereby users are directed to bogus, but authentic-looking, bank or credit card company Websites and enticed to enter personal information used for identity theft.

SPECIFIC DENIAL-OF-SERVICE ATTACKS 18 · 11**Exhibit 18.2** Normal TCP 3-Way Handshake

If users were to look closely at the URL of the site, however, they would observe a suspicious URL. *Pharming* is a variant of phishing that relies on some form of DNS poisoning so that a user going to the actual URL of a bank or credit card company will be redirected to the bogus site. Other variants of phishing have also emerged, such as *spearphishing* (directing a phishing attack at a specific person, job function, or group), *vishing* (a phishing attack via the telecommunications network, such as Voice-over-IP [VoIP]), and *smishing* (phishing via Short Message Service [SMS], or text, messages). For more details of such social engineering attacks, see Chapter 19 in this *Handbook*.

18.5.6 SYN Flooding. *SYN flooding* is a DoS attack that fits into the consumption-of-scarce-resource category. This DoS attack exploits the three-way handshake used by TCP hosts to synchronize the logical connection prior to data exchange. In a normal TCP host-to-host connection, the two hosts exchange three TCP segments prior to exchanging data, as shown in Exhibit 18.2:

1. The client sends a segment to the server with its initial sequence number (ISN). The SYN (synchronization) flag is set in this segment.
2. The server responds by sending a segment containing its ISN and acknowledges the client's ISN. This segment will have both the SYN and ACK (acknowledgment) flags set. At this point, the server allocates resources for the soon-to-be-established connection and waits for the third segment.
3. The client sends a segment acknowledging the server's ISN. This and all subsequent segments until the end of the session will have only the ACK flag set.

A SYN flood takes advantage of the three-way handshake and the fact that a server can have only a finite number of open TCP connections. The attack is launched when an attacker initiates connections for which there will never be a third segment (see Exhibit 18.3). After the server sends the segment in step 2, it waits for a response. Under normal circumstances, the client will respond within a few seconds. The server

18 · 12 DENIAL-OF-SERVICE ATTACKS

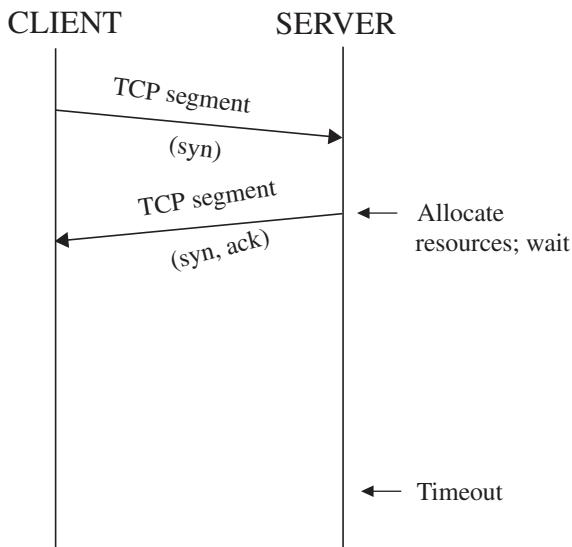


EXHIBIT 18.3 TCP SYN DoS Attack

might wait, say, 10 seconds before timing out and releasing the resources. But suppose the attacker sends hundreds of connection messages per second on a sustained basis. When these bogus connection attempts flood the target faster than they can time out, there will not be any resource left in which to establish a legitimate connection. This is the type of attack that was launched against PANIX in 1996.

18.5.7 Resource Starvation. *Resource starvation* is a catchall category for many other types of DoS attacks that are the result of some scarce resource—be it bandwidth, processing power, disk space—being consumed and exhausted. One example uses a novel UDP DoS attack, where an intruder forges UDP packets and uses them to connect the *echo* service of one machine (UDP port 7) to the character generator (*chargen*) service on another machine (UDP port 19). When this occurs, the two machines consume all available network bandwidth, and all machines on the same networks as the target machines can be affected and starved for resources.

But such attacks can happen locally, as well. Unfortunately, many authorized users carry out *local DoS attacks* that either consume system resources or deny users access. Numerous resource starvation attacks and programming flaws attack specific systems. For instance, by exceeding imposed quotas on disk space, multiuser systems can suffer from a resource starvation attack. Most operating systems limit the amount by which a user can exceed disk quota. Windows 2000 has a twist on this; although individual files can exceed their quota by only a small amount, if *every* file exceeds the quota, a user can consume a lot of extra disk space.

Another variant of this type of attack is called *Slowloris*, which targets Web servers and uses a relatively low amount of bandwidth by the attacker. Using the Slowloris scheme, an attacker sends a large number of connection requests to a target Web server and holds them open as long as possible. Under some circumstances, the affected Web server will exhaust its pool of connection resources on unfulfilled requests, thus denying legitimate requests. Many types of Web servers are susceptible to Slowloris, including several versions of Apache.

SPECIFIC DENIAL-OF-SERVICE ATTACKS 18 · 13

A Slowloris-variant is called *R-U-Dead-Yet* (RUDY). RUDY sucks connections resources from a Web server by sending large numbers of POST messages and randomly large Content-Length values in the HTTP message header.

18.5.8 Java, PHP, and ASP Attacks. There are many attacks that target program vulnerabilities in Websites that use Java, PHP (PHP: Hypertext Preprocessor), Active Server Pages (ASP), and other server-side scripting languages. One example is a 2011 DOS attack that resulted in CPU resource exhaustion. In this particular attack, an attacker would send variables to a Web server using the Hypertext Transfer Protocol (HTTP) POST method. By sending tens or hundreds of thousands of specially selected variable names, collisions in the hash table names could be forced, eventually consuming all CPU cycles on handling a single HTTP POST request.

18.5.9 Router Attacks. *Router attacks* have been developed for a wide variety of routers. Because routers are the backbone of the Internet and the gateway through which organizations connect to the Internet, killing a router denies network service for hundreds of machines. Productivity and financial repercussions from a hardware attack can be very severe. Popular targets for these kinds of attacks are routers from Ascend, Cisco, Juniper, Lucent, and 3Com. Unfortunately, many network managers make the attacks easy by employing Telnet or HTTP for remote access and not properly securing the network against remote access by anyone over the Internet.

18.5.10 Other Denial-of-Service Attacks. Even this list of specific DoS attacks is not exhaustive. For example, *Bonk* and *boink* (aka *bonk.c*) are DoS attacks that cause a target Windows 9x/NT system to crash. *Arnudp* (*arnudp100.c*) forges UDP packets to implement a DoS against UDP ports 7 (*echo*), 13 (*daytime*), 19 (*chargen*), and 37 (*time*), services that frequently run on UNIX and Windows NT systems. And *cacb.c* is a *cancelbot* that destroys existing Usenet news postings that fit certain criteria. (Some argue that *cacb.c* does not actually carry out a DoS attack; however, once activated, the program denies access to postings for targeted Usenet users.)

An example of network-based attack on a non-network resource would be *billion laughs* (an XML bomb), a DoS attack targeting Extensible Markup Language (XML) parsers. In XML, an *entity* is a piece of code that allows the programmer to define something that can be easily modified or reused in other documents, not unlike a programming object. In *billion laughs*, the attacker creates an XML document containing 10 entities, each defined as being composed of 10 instances of the previous entity ($10^{10} \approx 1$ billion, hence the name). Billion laughs was first introduced as early as 2003 but remained a problem late into at least the first decade of the 2000s.

On a final note, DoS attacks are not limited to networks. A *black fax* attack, for example, is a simple yet effective resource depletion attack; the attacker merely faxes a page containing a large black box to the victim, which causes the victim's fax to use a large amount of black toner. Sending a large number of these faxes will ultimately deplete the toner at the victim's fax and/or prevent legitimate faxes from being delivered. A security expert once glued the ends of a long strip of paper together to create an endless loop in a fax machine to attack a fax-spammer who refused to stop sending him unwanted travel offers. Another attack on faxes uses automated scripts for a computer to dial and redial the target machine with thousands of copies of (usually offensive) faxes; irritated nerds have been known to use this attack on fax-spammers unwise enough to use real telephone numbers for their own fax machines.

18 · 14 DENIAL-OF-SERVICE ATTACKS

18.6 PREVENTING AND RESPONDING TO DENIAL-OF-SERVICE ATTACKS.

DoS attacks are best prevented; handling them in real time is very difficult. And the most important way to protect a system is to harden the operating systems:

- Install them with security in mind.
- Monitor sites to be aware of security vulnerabilities.
- Maintain the latest versions of software where possible.
- Install all relevant security patches.

But a large measure of the prevention consists of packet filtering at network routers. Because attackers frequently hide the identity of the machines used to carry out the attacks by falsifying the source address of the network connection, techniques known as *egress filtering* and *ingress filtering* are commonly used as protective measures. As discussed later in this chapter, egress and ingress filtering are methods of preventing packets from leaving or entering the network, respectively, with an invalid source address. Blocking addresses that do not fit the criteria for legitimate source addresses and making certain that all packets leaving an organization's site contain legitimate addresses can thwart many DoS attacks.

Other packet-filtering methods that will help prevent DoS are to block all broadcast messages and most ICMP messages. There is no reason that a site should accept messages being broadcast to all hosts on the site. Furthermore, there is probably no good reason to allow all hosts to respond to *ping* or traceroute messages; in fact, most ICMP messages probably can be blocked.

In some instances, victims have set up response letters triggered to send and resend in large quantities so that they flood the attacker's address. Doing this is generally not a good idea. If these messages are sent to a legitimate address, the attacker may get the message and stop. But the attackers generally spoof the source IP address, so responding back to a possibly bogus IP network is not a good defensive posture because it may harm innocent victims. The best defense will involve the ISP.

In instances where the attacker's service provider can be identified and contacted, the victim can request that the service provider intervene. In these instances, it is usual for the ISPs to take appropriate action to stop the attack and find the perpetrator. However, in instances where a DoS appears to emulate or mimic another form of attack or when it continues for an unusually long period of time, the victim may want to take more aggressive action by contacting CERT/CC, the FBI, and other authorities that have experience with DoS attacks and some jurisdiction if the perpetrators are caught.

Real-time defenses are difficult but possible. Many routers and external intrusion detection systems (IDSs) can detect an attack in real time, such as too many connection requests per unit time from a given IP host or network address. A router might block the connection requests or an IDS might send a pager message to a security administrator.

However, attacks such as smurfs can suck up all of the bandwidth even before the packets get to the target site. Cooperation by ISPs and end-user sites is required to fully combat DoS attacks. This will be addressed further as part of the discussion of responding to DDoS.

18.7 DISTRIBUTED DENIAL-OF-SERVICE ATTACKS.

DDoS tools use amplification to augment the power of the attacker. By subverting poorly secured systems

DISTRIBUTED DENIAL-OF-SERVICE ATTACKS 18 · 15

into sending coordinated waves of fraudulent traffic aimed at specific targets, intruders can overwhelm the bandwidth of any given victim.

In a DDoS attack, the attacking packets come from tens or hundreds of addresses rather than just one, as in a standard DoS attack. Any DoS defense that is based on monitoring the volume of packets coming from a single address or single network will fail since the attacks come from all over. Rather than receiving, for example, 1,000 gigantic *pings* per second from an attacking site, the victim might receive one *ping* per second from each of 1,000 attacking sites.

One of the other disconcerting things about DDoS attacks is that the handler can choose the location of the agents. So, for example, a handler could target several North Atlantic Treaty Organization (NATO) sites as victims and employ agents that are all in countries known to be hostile to NATO. The human attacker might be sitting in, say, Canada.

Like DoS attacks, all of the DDoS attacks employ standard TCP/IP messages—but employ them in some nonstandard ways. Common DDoS attacks include *Tribe Flood Network* (TFN), *Trinoo*, *Stacheldraht*, and *Trinity*. The sections that follow present some details about these attacks.

18.7.1 Short History of Distributed Denial of Service. Denial-of-service attacks under a number of guises have been around for decades. Distributed DoS attacks are much newer. In late June and early July 1999, groups of hackers installed and tested a DDoS tool called *Trinoo* (see Section 18.7.3.1) to launch medium to large DDoS attacks. Their tests involved over 2,000 compromised systems and targets around the world.

Most of the literature suggests that the first documented large-scale DDoS attack occurred in August 1999, when *Trinoo* was deployed in at least 227 systems (114 of which were on Internet2) to flood a single University of Minnesota computer; this system was down for more than two days.

On December 28, 1999, CERT/CC issued its Advisory CA-1999-17 (www.cert.org/advisories/CA-1999-17.html) reviewing DDoS.

On February 7, 2000, Yahoo! was the victim of a DDoS during which its Internet portal was inaccessible for three hours. On February 8, Amazon, Buy.com, CNN, and eBay were all hit by DDoS attacks that caused them either to stop functioning completely or to slow down significantly. And on February 9, 2000, E*Trade and ZDNet both suffered DDoS attacks. Analysts estimated that during the three hours Yahoo! was down, it suffered a loss of e-commerce and advertising revenue that amounted to about \$500,000. According to bookseller Amazon.com, its widely publicized attack resulted in a loss of \$600,000 during the 10 hours it was down. During their DDoS attacks, Buy.com went from 100 percent availability to 9.4 percent, while CNN.com's users went down to below 5 percent of normal volume and Zdnet.com and E*Trade.com were virtually unreachable. Schwab.com, the online venue of the discount broker Charles Schwab, was also hit but refused to give out exact figures for losses. One can only assume that to a company that does \$2 billion weekly in online trades, the downtime loss was huge. Another type of damage caused indirectly by the DDoS was the decline in stock values in the 10 days following the attacks: eBay suffered a 24 percent decline, Yahoo! dropped 15 percent, and Buy.com dropped 44 percent.

These attacks were orchestrated by 15-year-old Michael Calce, a youngster living in the west end of Montreal island in Quebec whose pseudonym was *Mafiaboy*. He eventually pled guilty to 56 charges of computer crime and served eight months of juvenile detention.⁴

18 · 16 DENIAL-OF-SERVICE ATTACKS

These types of DDoS attacks have continued since the summer of 1999. One of the best-known incidents was a series of DDoS attacks again Steve Gibson's GRC.com Website in May 2001. The attacker was a 13-year-old using an Internet Relay Chat (IRC) bot, automated programs that exploit systems using IRC clients to become DDoS zombies.

The DDoS attack that had the highest potentially devastating impact, however, occurred on October 21, 2002, when all of the top-level DNS root servers were subjected to a sustained attack by thousands of zombies. Nine of the 13 DNS root servers were knocked off the Internet; the remaining four were able to keep operating during the attack. All of the major ISPs and many large private networks maintain their own DNS systems, although most servers ultimately rely on the root servers to find noncached DNS entries. The attack lasted for just an hour or two; had it continued for much longer, the remaining servers probably would have been overwhelmed, effectively blocking DNS host name/address translation.

A disturbing and growing trend is the use of DDoS as an extortion tool. An increasing number of criminals are using DDoS tools as a way to threaten an attack rather than actually disrupting a target organization's network. Although several providers of network, security, and consulting services claim that they and many of their customers have received such extortion demands, few are public about naming the targets—many of which are acceding to the threats and paying the blackmail.⁵ Most experts agree that the extortion demands should not be met; doing so only encourages the criminal behavior. If such a threat is received, the organization's ISP and law enforcement authorities should be contacted immediately.

18.7.2 Distributed Denial-of-Service Terminology and Overview. To describe and understand DDoS attacks, it is important to understand the terminology that is used to describe the attacks and the tools. Although the industry has more or less settled on some common terms, that consensus did not come about until well after many DoS/DDoS attacks had already appeared in the hacker and mainstream literature. Early descriptions of DDoS tools used a jumble of terms to describe the various roles of the systems involved in the attack. At the CERT/CC Distributed System Intruder Tools Workshop held in November 1999, some standard terminology was introduced and those terms are used in the paragraphs below. To align those terms and the terms used by the hacker literature as well as early descriptions, here are some synonyms:

Intruder—also called the *attacker* or *client*.

Master—also called the *handler*.

Daemon—also called an agent, bcast (broadcast) program, or zombie.

Victim—also called the *target*.

DoS/DDoS attacks actually have two victims: the ultimate target and the intermediate system(s) that were exploited and loaded with daemon software. In this chapter, the focus is on the end-of-the-line DoS/DDoS victim.

DDoS attacks always involve a number of systems. A typical DDoS attack scenario might follow roughly these three steps:

1. The *intruder* finds one or more systems on the Internet that can be compromised and exploited (see Exhibit 18.4). This is generally accomplished using a stolen

DISTRIBUTED DENIAL-OF-SERVICE ATTACKS 18 · 17

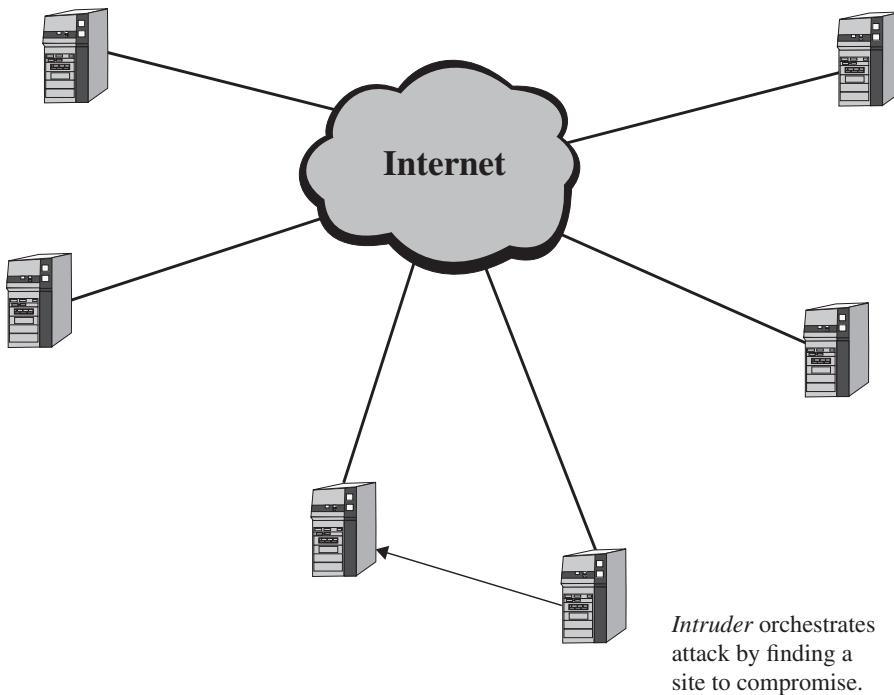


EXHIBIT 18.4 DDoS Phase 1

account on a system with a large number of users or inattentive administrators, preferably with a high-bandwidth connection to the Internet. (Many such systems can be found on college and university campuses.)

2. The compromised system is loaded with any number of hacking and cracking tools, such as scanners, exploit tools, operating system detectors, rootkits, and DoS/DDoS programs. This system becomes the DDoS *master*. The master software allows it to find a number of other systems that can themselves be compromised and exploited. The attacker scans large ranges of IP network address blocks to find systems running services known to have security vulnerabilities. This *initial mass-intrusion phase* employs automated tools to remotely compromise several hundred to several thousand hosts, and installs DDoS agents on those systems. The automated tools to perform this compromise are not part of the DDoS toolkit but are exchanged within groups of criminal hackers. These compromised systems are the initial victims of the DDoS attack. These subsequently exploited systems will be loaded with the DDoS *daemons* that carry out the actual attack (see Exhibit 18.5).
3. The intruder maintains a list of *owned systems* (sometimes spelled *pwned* by hackers), the compromised systems with the DDoS daemon. The actual *denial-of-service attack phase* occurs when the attacker runs a program at the master system that communicates with the DDoS daemons to launch the attack. Here is where the intended DDoS victim comes into the scenario (see Exhibit 18.6).

Communication between the master and daemons can be obscured so that it becomes difficult to locate the master computer. Although some evidence may exist on one or more machines in the DDoS network regarding the location of the master, the daemons

18 · 18 DENIAL-OF-SERVICE ATTACKS

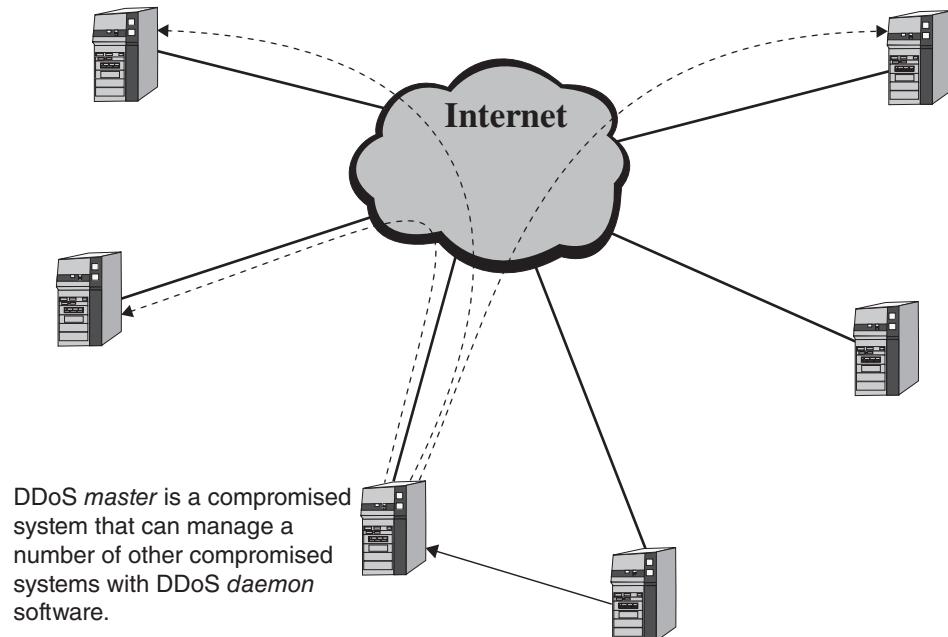


EXHIBIT 18.5 DDoS Phase 2

normally are automated so that it is not necessary for an ongoing dialog to take place between the master and the rest of the DDoS network. In fact, typically techniques are employed to deliberately camouflage the identity and location of the master within the DDoS network. These techniques make it difficult to analyze an attack in progress and difficult to block attacking traffic and trace it back to its source.

In most cases, the system administrators of the infected systems do not even know that the daemons have been put in place. Even if they do find and eradicate the DDoS

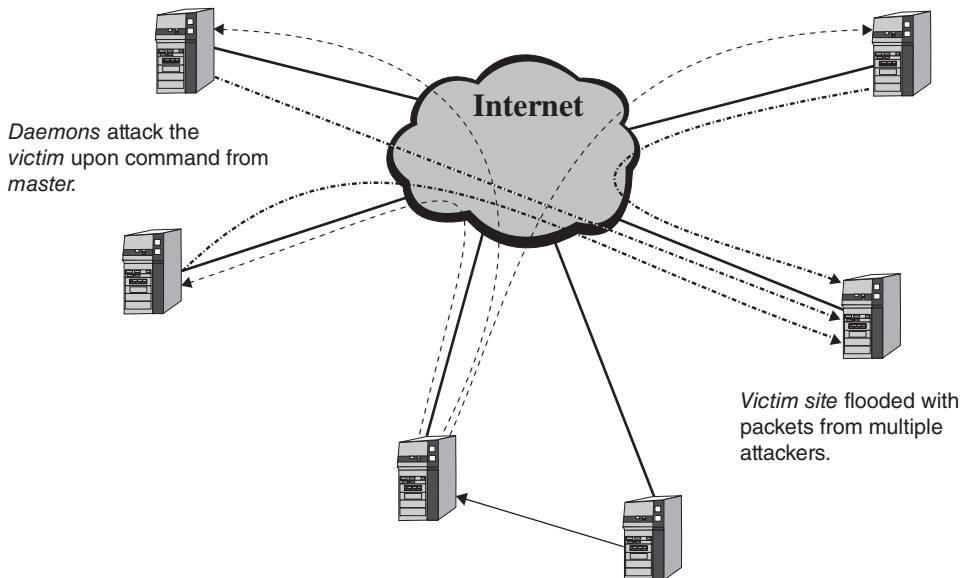


EXHIBIT 18.6 DDoS Phase 3

DISTRIBUTED DENIALOFSERVICE ATTACKS 18 · 19

software, they cannot help anyone determine where else the software may have been placed. Popular systems to exploit are a site's Web, email, name, or other servers, since these systems are likely to have a large number of open ports, a large amount of traffic, and are unlikely to be quickly pulled off-line even if an attack can be traced to them.

18.7.3 Distributed Denial-of-Service Tool Descriptions. This section provides some details on how some of the major DDoS tools work.

18.7.3.1 Trinoo (TrinOO). *Trinoo* or *Trin00* was the first known DDoS tool, appearing in the summer of 1999. The typical installation of Trinoo is similar to the scenario painted above where an attacker plants handler software on a system and the handler, in turn, loads the attack software on the agents. Trinoo is a distributed SYN DoS attack.

Trinoo uses a number of TCP and UDP ports:

- Masters listen on TCP port 27665 for attacker-to-master communication.
- Daemons listen on UDP port 27444 for master-to-daemon communication.
- Masters listen on UDP port 31335 for daemon-to-master communication.

These are default port numbers, and variants used other ports. The human attacker can control a Trinoo master (the handler) remotely via a connection to TCP port 27665. After connecting, the attacker gives the expected password, *betalmostdone*.

The Trinoo master program is typically named *master.c* and the daemon is *ns.c*. Communication between the Trinoo master (handler) and daemons (agents) is via UDP. Master-to-daemon communications employ UDP datagrams on port 27444. All commands contain a password, the default being *l44adsl*. All valid commands contain the substring *l44*.

Communication from the Trinoo daemons to the master use UDP datagrams on port 31335. When the daemon starts, it sends a message to the master containing the string **HELLO**. The Trinoo *keep alive* function is accomplished by an exchange between the master and daemon: The master sends a Trinoo *mping* command, which sends the string *png* to a daemon; the daemon responds by sending the string *PONG* to the master.

The passwords are there to prevent system administrators from being able to take control of the masters and daemons that form the Trinoo network. Other default passwords in the initial attacks were *gOrave* to start the Trinoo master server and *killme* to control the master's *mdie* command to kill the Trinoo processes. Like the port numbers, the passwords can be changed easily by the attackers.

Intrusion detection software or system management routine analysis can look for a number of things that might indicate the presence of Trinoo:

- A system listening on UDP port 27444 could be a Trinoo daemon.
- Trinoo daemon communication will contain the string *l44*.
- The SYN flood mechanism picks the destination port using a random number generator function.
- A Trinoo daemon will send the string *PONG* if it receives a *png* command.
- A system listening on TCP port 27665 could be a Trinoo master.

18 · 20 DENIAL-OF-SERVICE ATTACKS

- A system listening on UDP port 27444 could be a Trinoo master.
- UDP packets will contain the string *l44adsl*.

18.7.3.2 Tribe Flood Network. The *Tribe Flood Network* (TFN) appeared after Trinoo. TFN runs primarily on compromised UNIX systems exploited using buffer overrun bugs in the remote procedure call (RPC) service. TFN client and daemon programs implement a DDoS network capable of employing a number of attacks, such as ICMP flood, SYN flood, UDP flood, and smurf-style attacks.

TFN is noticeably different from Trinoo in that all communication between the client (attacker), handlers, and agents uses *ICMP echo* and *echo reply* packets. Communication from the TFN client to daemons is accomplished via *ICMP echo reply* packets. The absence of TCP and UDP traffic sometimes makes these packets difficult to detect because many protocol monitoring tools are not configured to capture and display the ICMP traffic.

Remote control of the TFN network is accomplished by executing a program on the client system. The program also can be executed at the handler system by the client via some host-to-host connection method, such as connecting to an exploited TCP port or using a UDP- or ICMP-based remote shell. The program must be supplied:

- The IP address list of hosts that are ready to carry out the flood attack
- The type of attack to be launched
- The IP address list of target hosts
- The port number for SYN attack

No password protection is associated with TFN. Each command to the daemons is sent in the Identifier field of the ICMP packet; values *345*, *890*, and *901* start the *SYN*, *UDP*, and *ICMP* flood attacks, respectively. The Sequence Number field in the *echo reply* message is always set to *0x0000*, which make it look like the response to the initial *echo* packet sent out by the *ping* command.

The TFN client program typically is named *tribe.c* and the daemon is *td.c*.

18.7.3.3 Stacheldraht. Stacheldraht (German for *barbed wire*) is a DDoS tool that appeared in August 1999 and combines features of Trinoo and TFN. It also contains some advanced features, such as encrypted attacker–master communication and automated agent updates.

Stacheldraht uses a Trinoo-like client/server architecture. The handler listens on TCP port 16660 for client (intruder) commands, and the agents listen on TCP port 65000 for commands from the handler. Agent responses to the handler employ *ICMP echo reply* messages. The possible attacks are similar to those of TFN; namely, ICMP flood, SYN flood, UDP flood, and smurf attacks.

Trinoo and TFN exchange commands in plaintext. Trinoo, being TCP-based, is also subject to common TCP attacks, such as session hijacking. Stacheldraht addresses these deficiencies by employing an encrypting *telnet alike* client. (*Telnet alike* is a Stacheldraht term.) The client uses secret-key cryptography.

The Stacheldraht network comprises a number of programs. The attacker uses an encrypting client called *telnetc/client.c* to control one or more handlers. The handler program is called *mserv.c*, and each handler can control up to 1,000 agents. The agent software, *leaf/td.c*, coordinates the attack against one or more victims upon command from the handler.

DISTRIBUTED DENIAL-OF-SERVICE ATTACKS 18 · 21

18.7.3.4 TFN2 K. *Tribe Flood Network 2 K* (TFN2 K) was released in December 1999 and targets UNIX and Windows NT servers. TFN2 K is a complex variant of the original TFN with features designed specifically to:

- Make TFN2 K traffic difficult to recognize and filter
- Remotely execute commands
- Hide the true source of the attack using IP address spoofing
- Transport TFN2 K traffic over multiple transport protocols including UDP, TCP, and ICMP
- Confuse attempts to locate other nodes in a TFN2 K network by sending “decoy” packets

TFN2 K, like TFN, can consume all of a system’s bandwidth by flooding the victim machine with data. But TFN2 K, unlike TFN, also includes attacks designed to crash or introduce instabilities in systems by sending malformed or invalid packets, such as those found in the Teardrop and Land attacks.

TFN2 K uses a client server architecture in which a single client issues commands simultaneously to a set of TFN2 K agents. The agents then conduct the DoS attacks against the victim(s). The agent software is installed in a machine that already has been compromised by the attacker.

18.7.3.5 Shaft. Trinoo, TFN/TFN2 K, and Stacheldraht were the first, and possibly best studied, DDoS tools. Other tools followed, of course, and have become increasingly more complex, but these early tools remain available and some systems remain vulnerable to these forms of attack.

In November 1999, for example, the *Shaft* DDoS tool became available. A Shaft network looks conceptually similar to a Trinoo network with client-managing handler programs (*shaftmaster*) that, in turn, manage agent programs (*shaftnode*). Like Trinoo, handler-agent communication uses UDP, with the handler(s) listening on port 20433 and the agent(s) listening on port 18753. The client communicates with the handler by telnetting to TCP port 20432. The attack itself is a packet flooding attack, and the client controls the size of the flooding packets and duration of the attack. One signature of Shaft is that the sequence number for all TCP packets is always 0x28374839.

18.7.3.6 HTTP Apache Attack. In August 2000, a DDoS attack against Apache Web servers was first detected. The attack took advantage of a vulnerability whereby a URL sent to an Apache Web server containing thousands of forward slashes (/) would put the server into a state that would consume enormous CPU time. This particular attack was launched by over 500 compromised Windows computers and would, presumably, succeed against Apache Web servers prior to version 1.2.5.

18.7.3.7 Trinity. In September 2000, a new DDoS tool called *Trinity* was reported. Trinity is capable of launching several types of flooding attacks on a victim site, including ACK, fragment, RST, SYN, UDP, and other floods. Trinity agent software must be placed on Linux systems compromised by a buffer overflow vulnerability. The agent binary code typically is found in */usr/lib/idle.so*. Communication from the handler or intruder to the agent, however, is accomplished via Internet Relay Chat (IRC) or America Online’s ICQ instant messaging software. Whereas the attacker has

18 · 22 DENIAL-OF-SERVICE ATTACKS

to keep track of the IP addresses of compromised systems with Trinoo and TFN, all of the Trinity agents report back to the attacker by appearing in the same chat room. The original reports were that the Trinity agent communicated over an IRC channel called `#b3eblebr0x`; other IRC channels presumably are being used for DDoS, as well. IRC uses TCP ports 6665 to –6669, and Trinity appears to use port 6667. In addition, a binary called `/var/spool/uucp/uucico` is a backdoor program that listens on TCP port 33270 for connections; an attacker connecting on that port and providing the password `!@#` will achieve rootshell on the affected system.

18.7.3.8 SubSeven. Zombie software is not always distributed by an attacker exploiting a vulnerability of an exposed system. Indeed, very often the user is the culprit. Trojan horses are often the mechanism for distributing the zombie code. The *SubSeven* software, for example, is a backdoor virus. SubSeven often gets on a user’s system because it is distributed within programs available via Usenet and other Internet sites, such as some game or pornography programs (e.g., `SexxyMovie.mpeg.exe`). Potential attackers frequently scan computer systems today, particularly residential systems connected to the Internet via DSL or cable modem, for the presence of SubSeven, which provides a potential backdoor into users’ systems; system administrators also are learning to scan for this dangerous program on their own systems.

18.7.3.9 Mydoom. The description above of some of the first DDoS tools is meant to show the early baseline from which later tools evolved, many of which are simply a variant of the early tools. *Mydoom* (aka W32.MyDoom@mm and other aliases) is an example of a DoS attack via an email worm, first launched in 2004. Mydoom appears as an email with some sort of error-type message in the subject line (e.g., “Error” or “Mail System Failure”). Interestingly, the email could be prepared in several different languages, including English, French, and German. The email contained an attachment that, if executed, forwarded the worm to addresses found in a variety of files on the local system, including the email address book. The worm could also be copied to shared folders of a user’s peer-to-peer networking sites. Early versions of Mydoom contained a payload that would launch a DoS attack against the SCO Group on February 1, 2004; the SCO Group was unpopular because they asserted, in a \$1 billion lawsuit against IBM, that several Linux distributions were violating SCO’s UNIX intellectual assets. Other variants of Mydoom were subsequently released, including a version employed in attacks on South Korea and the United States in July 2009 that included the launching of a botnet.

18.7.3.10 LOIC, HOIC, and HULK. Another variant of these types of attacks is the *Low Orbit Ion Cannon* (LOIC). Designed to be used as a network stress test, it can also be deployed as a DoS or DDoS tool. LOIC works, basically, by flooding a target with TCP and UDP packets in order to learn whether the routers can handle the load and how they will respond—or to usurp all of the bandwidth. Originally written in C#, a JavaScript variant (LS LOIC) and Web version (Low Orbit Web Cannon) have also appeared. LOIC has been used extensively by the hacktivist collective *Anonymous*.

The *High Orbit Ion Cannon* (HOIC) is a modification of LOIC. It has an easy-to-use graphical user interface that allows attackers to specify targets and *booster* scripts to target multiple pages on the targeted server, much like a shotgun instead of a rifle. The *intensity* setting defines the number of requests per second (2/sec for *low* and 8/sec for *high*).

DENIAL-OF-SERVICE USING EXPLOITABLE SOFTWARE 18 · 23

The *HTTP Unbearable Load King* (HULK) is a variant of other tools that bombard a Web server with a overwhelming number of packets. Most prior tools, however, send TCP SYN requests or other predictable packets, which allows the firewall and/or server to detect the attack and mount a defense. HULK generates a unique, nonpredictable set of requests designed to frustrate defenses based on pattern recognition and packet filtering.

18.8 DENIAL-OF-SERVICE USING EXPLOITABLE SOFTWARE. The tools just discussed employ a common DoS approach: An attacker exploits a vulnerability of a potential victim and uses that system to launch attacks on the intended victim. Later rounds of DDoS attacks, however, use code that is commonly available and that has known vulnerabilities. And all too often, the vulnerabilities are exploited after a patch to fix them has been released but ignored by some system administrators. A case in point is the flurry of such releases in mid-2001.

In May 2001, a buffer overflow exploit was discovered in the Microsoft IIS Indexing Service. In mid-June, Microsoft released a security bulletin warning that administrative scripts (.ida files) and Internet data queries (.idq files) did not do proper bounds checking. As it happens, what seems like the vast majority of IIS servers did not get the patch, and, in essence, every unpatched IIS server became a DDoS zombie.

18.8.1 Code Red. In July 2001, eEye Digital Security and several other security organizations around the Internet saw an alarming number of TCP port 80 scans on the Internet. What they eventually discovered was what became known as the *Code Red Worm*.

Code Red had three distinct phases. The *propagation phase* occurred during the first 19 days of the month. During this phase, the attacking system scanned target systems on TCP port 80 and sent a specially crafted HTTP GET request that exploited the IIS buffer overflow (even if the Index Service is not running). A common log entry might appear as:

If the exploit was successful, the worm ran in RAM of the infected server and spawned 99 new threads to attack a quasi-random set of IP addresses. If the exploited server's native language was English, the server's Web page was defaced with a message that said "Welcome to <http://www.worm.com>! Hacked by Chinese!" This message would stay up for 10 hours and then disappear.

The *flood phase* occurred on days 20 to 27 of the month. This is when the attack really happened; every day between 8:00 and 11:59 P.M. UTC, the compromised servers sent 100 KB packets to the IP address 198.137.240.91, which formerly was assigned to www.whitehouse.gov. (Once the actions of Code Red were discovered, the IP address of www.whitehouse.gov was changed from the targeted address.)

Days 28 to 31 of the month were the *termination phase*, when the worm became dormant. Code Red was relatively innocuous, compared to what it could have been;

18 · 24 DENIAL-OF-SERVICE ATTACKS

Propagation by Network Shares

- Worm code written to writeable share.

Propagation by E-mail

- E-mail client using IE HTML reader.
- Autoexecutes README.EXE.



Propagation via Web Browser

- Browser downloads README.EML via JS.
- Vulnerable browser autoexecutes.

Propagation by IIS

- Found by port 80 scan.
- Vulnerability exploited.
- ADMIN.DLL obtained via tftp.

Infected System

- GUEST added to Admin.
- Sharing turned on.
- README.EXE file forwarded via e-mail.
- If Web server, all HTML content pages compromised.

EXHIBIT 18.7 NIMDA Propagation Vectors

once asleep, the worm stayed asleep, although it could be reawakened. Removing the worm program from RAM required only a reboot, and the patch from Microsoft would prevent further infection.

As an aside, although only IIS servers could be exploited, many other devices that listen on port 80 were also affected. Cisco 600 DSL routers and HP JetDirect devices, for example, listen on port 80 and would crash when they received the buffer overflow packet.

Three different variants of Code Red existed on the Internet, all acting as described. In August 2001, a couple of new variants appeared that were called *Code Red II*. Unlike Code Red, Code Red II did not deface Web pages nor did it launch a DDoS attack on any given site. Instead, this worm was destructive, installing backdoors on infected servers, changing many registry settings, installing a Trojan horse version of *explorer.exe* (Windows Explorer), and disabling the System File Checker (SFC) utility. The worm also spread quickly, employing up to 300 threads at a time looking for other systems to infect.

18.8.2 NIMDA. The next evolution appeared in September 2001 and was called *NIMDA*. NIMDA (*admin* backward) was unique because it exploited multiple vulnerabilities in Microsoft code, namely IIS, Internet Explorer (IE), and the Message Application Program Interface (MAPI). As a result, NIMDA had four distinct propagation vectors (see Exhibit 18.7):

- 1. IIS.** When a Web server was found, the attacker attempted to exploit various IIS vulnerabilities, including IIS sadmind, a Code Red II root.exe or other backdoor program, or IIS Directory Traversal. If successful, the attacker used the Trivial File Transfer Protocol (tftp) from cmd.exe to send the worm code (admin.dll) to the victim.
- 2. Web browser.** The worm on an infected server created a copy of itself in a file called *readme.eml*. The worm also altered every Web-content file at the infected site with a small JavaScript code that pointed to this file. When a user browsed to

DEFENSES AGAINST DISTRIBUTED DENIAL-OF-SERVICE ATTACKS 18 · 25

the infected Web server, the infected page's JS code was activated, and *readme.eml* was downloaded. Vulnerable versions of Internet Explorer would auto-execute the file while most other browsers would not.

3. **Email.** NIMDA sent itself to all of the email addresses found in the InBox and Address Book of an infected server in a MIME-encoded, 56KB attached file named *readme.exe*. The file had an “audio/x-wav” section that contained the worm. Email clients using IE 5.1 or earlier to display HTML would automatically execute the attachment if the message was opened or previewed.
4. **Network shares.** When on an infected system, the worm copied itself to all local directories on the victim host and to all open, writeable network shares. The worm also set up shares on the victim host.

In addition, the GUEST account on the infected system was activated and made a member of Administrator group. (Interestingly, after a summer of increasing DDoS worms, almost all such activity ceased after 9/11 for many months.)

Over the years, DDoS attacks have continued, but there is a dual purpose. Some DDoS attacks are intended to knock a server or network off the Net the old-fashioned way, namely by sucking up all of the bandwidth or other resource. Other attacks, however, use the target site's defenses against itself; by performing reconnaissance to understand how a site will respond to various stimuli, the attacker can actually send a crafted set of packets to the target site that will cause the target to self-limit bandwidth and/or acceptable source addresses.

18.9 DEFENSES AGAINST DISTRIBUTED DENIAL-OF-SERVICE ATTACKS.

As with DoS attacks, a site cannot in isolation defend itself from DDoS attacks. Members of the Internet community must work together to protect every site against becoming the source of attacks or forwarding the attacks. This section discusses some ways to help prevent the spread of DDoS attacks by limiting the distribution of the tools and by limiting the propagation of the offending attack packets.

Although not discussed in detail here, another point needs to be made about DDoS attack responses. As discussed in Chapter 53 in this *Handbook*, victims of such an attack should maintain detailed logs of all actions they take and events they detect. These logs may prove invaluable in understanding the attack, in preventing other attacks at the initial target and others, and in aiding law enforcement efforts to track down the perpetrators.

18.9.1 User and System Administrator Actions. These steps should be taken to minimize the potential that an individual system will be compromised and attacked or used as a stepping-stone to attack others:

1. Keep abreast of the security vulnerabilities for all of the site's hardware, operating systems, and application and other software. This sounds like a Herculean task, but it is essential to safeguarding the network. Apply patches and updates as soon as possible. Standardize on certain hardware, operating systems, and software where feasible to help manage the problem.
2. Use host-based firewall software on workstations to detect an attack.
3. Monitor systems frequently to test for known operating system vulnerabilities. Regularly check to see what TCP/UDP ports are in use using the *netstat -a*

18 · 26 DENIAL-OF-SERVICE ATTACKS

command; every open port should be associated with a known application. Turn off all unused applications.

4. Regularly monitor system logs and look for suspicious activity.
5. Use available tools to periodically audit systems, particularly servers, to ensure that there have been no unauthorized/unknown changes to the file system, registry, user account database, and so on.
6. Every client system must use anti-malware software that updates itself continuously to keep abreast of the constantly changing threat landscape. Effective anti-malware tools include integration with browsers for automatically blocking access to Websites that are known to harbor dangerous code. Such tools also integrate into email clients to block automatic opening of attachments. Users should not be permitted to disable the anti-malware software.
7. Do not download software from unknown, untrusted sites. If possible, know the author of the code. Even better, download source code, review it, and compile it on a trustworthy system rather than downloading binaries or executables.
8. Keep up with and follow recommendations from NIST, US-CERT, anti-virus research labs, hardware and software vendors, and other sources of best practices.

18.9.2 Local Network Actions. Even if users lock down their systems so that no vulnerability has gone unpatched and no exposure unprotected, the local network itself still can be at risk. Local network managers and network administrators can take several steps to protect all of their own users as well as the rest of the Internet community:

1. Every network connected to the Internet should perform egress address filtering at the router. *Egress filtering* means that the router should examine the Source Address field of every outgoing IP packet sent to the Internet to be sure that the NET_ID matches the NET_ID of the network. Historically, firewalls have been used to protect networks from attacks from the outside world. But those attacks come from somewhere, so sites should also use the firewall to *protect* the outside world.
2. Networks should block incoming packets addressed to the broadcast address (the all-ones HOST_ID). There is no legitimate reason that an external network device should be sending a broadcast message to every host on a network.
3. To prevent a site from being used as a broadcast amplification point, turn off the Directed Broadcast capability at the router unless it is absolutely essential. If it is essential, reexamine the network to see if there is not a better way or if the broadcast scope can be minimized. Even where Directed Broadcasts are useful, typically they are needed only *within* the enterprise and are not required for hosts on the outside.
4. RFC 1918 defines three blocks within the IP address space that are reserved for private IP networks; these addresses are not to be routed on the Internet.

IP Address Range	Network ID/ Subnet Mask	Number of Equivalent Classful IP Networks
10.0.0.0–10.255.255.255	10/8 prefix	1 Class A network
172.16.0.0–172.31.255.255	172.16/12 prefix	16 Class B networks
192.168.0.0–192.168.255.255	192.168/16 prefix	256 Class C networks

DEFENSES AGAINST DISTRIBUTED DENIAL-OF-SERVICE ATTACKS 18 · 27

In addition, there are a number of reserved IP addresses per RFC 3330 that are never assigned to public networks or hosts, including:

0.0.0.0/32	Historical broadcast address
127.0.0.0/8	Loopback network identifier
169.254.0.0/16	Link-local networks
192.0.2.0/24	TEST-NET
224.0.0.0/4	Multicast address range
240.0.0.0/5	Reserved for future use
255.255.255.255/32	Broadcast

Attackers commonly use IP address spoofing, generally by using one of the RFC 1918 private addresses or one of the other reserved addresses. Firewalls should immediately discard *any* packet that contains any RFC 1918 or reserved IP address in the Source Address or Destination Address field; such packets should never be sent to the Internet.

5. Block all unused application ports at the firewall, particularly such ports as IRC (6665–6669/tcp) and those known to be associated with DDoS software.
6. Use stateful firewalls that can better examine a packet within the context of the entire packet exchange. Stateless firewalls only look at packets according to a simple set of rules but are not aware of the overall packet traffic on the network (e.g., a stateless packet filter might forward an ICMP *echo* Reply packet to be sent out because the rules allow it; a stateful packet filter will forward the packet only if the rule allows it *and* there was a corresponding ICMP *echo* Request).
7. Use intrusion detection and intrusion prevention to protect the network. For example, personal firewall software can be installed on every workstation to help detect an attack on individual systems; this strategy is particularly useful at sites that have a large number of systems *in front* of a firewall (e.g., colleges). It is no coincidence that so many daemons reside on college and university computers that have been *Owned* (i.e., taken over by hackers).
8. Regularly monitor network activity so that aberrations in traffic flow can be detected quickly.
9. Educate users about events to watch for on their systems and how to report any irregularity that might indicate that someone or something has tampered with their system. Educate the help desk and technical support to assist those users who make such reports. Have an intelligence-gathering system within the organization so that such reports can be coordinated centrally to spot trends and to devise responses.
10. Follow NIST, US-CERT, and other best practices procedures.

18.9.3 Internet Service Provider Actions. ISPs offer the last hope in defeating the spread of a DDoS attack. Although the ISP cannot take responsibility for locking down every customer's host systems, ISPs have—and should accept—the responsibility to ensure that their network does not carry packets that contain obviously “bad” packets. Some of the steps that ISPs can take include:

1. As mentioned, attackers commonly employ IP address spoofing using an RFC 1918 private address or other reserved address. Amazingly, many ISPs will route

18 · 28 DENIAL-OF-SERVICE ATTACKS

these packets. Indeed, there is no entry in their routing table telling them where to send the packets; they merely forward them to a default upstream ISP. Any packet that contains any RFC 1918 or reserved IP address in the IP Source Address or Destination Address field should be discarded immediately.

2. Perform ingress (and egress) address filtering. *Ingress filtering* means that ISPs should examine every incoming packet to their network from a customer's site and examine the IP Source Address field to be sure that the NET_ID matches the NET_ID assigned to that customer. Doing this will require additional configuration at the router and may even result in slight performance degradation, but the trade-off is certainly well worth the effort. The ISPs also should perform egress filtering to check their outbound packets to upstream and peer ISPs.
3. Disable IP directed broadcasts.
4. Pay careful attention to high-profile systems (servers) and customers.
5. Educate customers about security and work with them to help protect themselves.

Most of the ISP community takes at least some of these steps. Users should insist that their ISPs provide at least these protections and should not do business with those that do not. RFC 3013 and the North American Network Operators' Group (NANOG) are good sources of information for ISPs.

18.9.4 Exploited Software Defensive Actions. There are a number of defensive steps that can be taken to avoid or mitigate problems due to Code Red/NIMDA-like attacks that exploit software. Several of these recommendations are controversial due to their implicit boycott of products from a specific vendor.

- If using IIS, consider using alternate Web server software. If use of IIS is essential, keep IIS and the operating system up to the latest patch revision. Microsoft's IIS Cumulative Patch, for example, does *not* clean a system of many of the backdoors used for exploits.
- If using Internet Explorer, consider using alternate browser software. If you must use IE, secure it against MIME auto-execution. Note that as late as September 2012, unpatched vulnerabilities were being reported by US-CERT in IE6, IE7, IE8, and IE9, and some authorities (e.g., the government of Germany) were recommending against the use of any version of IE.
- Disable all unused accounts on servers and other systems. In particular, enable the Guest account or anonymous access only if absolutely necessary.
- Disable JavaScript, Java, and ActiveX on browsers unless absolutely necessary.
- Do not execute or open *any* email attachment unless expected, known, and verified.
- Use the most up-to-date antivirus signature files.
- Unbind file and print sharing from TCP/IP. In some cases, this will require installing NetBEUI for file and print sharing.

18.9.5 Other Tools under Development or Consideration. Responses to DDoS attacks are not limited to the defensive steps just listed. Indeed, proactive responses to the prevention and detection of DDoS attacks are an active area of research. These proposals are merely samples of some ways for dealing with DDoS attacks; the first adds new hardware to the Internet, the second requires changing Web server and

DEFENSES AGAINST DISTRIBUTED DENIAL-OF-SERVICE ATTACKS 18 · 29

client software, and the latter two require incrementally changing software in all of the Internet's routers and hosts, respectively. Upgrading Web browsers is probably the most practical strategy even though there are millions of copies in distribution; the vast majority come from just a handful of vendors, and users tend to upgrade frequently anyway.

18.9.5.1 Distributed Traffic Monitor. One method that has been proposed is to examine the network at the ISP level and build a type of intelligent, distributed network traffic monitor; in some sense, this would be like an IDS for the Internet. ISPs, peering points, and/or major host servers would have traffic monitor hardware using IP and the Internet for communications, much like today's routing protocols. Each node would examine packets and their contents, doing a statistical analysis of traffic to learn the normal patterns. These devices would have enough intelligence to be able to detect changes in traffic level and determine whether those changes reflected a normal condition or not. As an example, suppose that such hardware at Amazon.com were to identify a DoS attack launched from an ISP in Gondwanaland; the traffic-monitoring network would shut off traffic to Amazon coming from that ISP as close to the ISP as possible. In this way, the distributed network of monitors could shut traffic off at the source.

The hardware would need to be informed about traffic-level changes due to normal events, such as a new Super Bowl commercial being posted on YouTube or a new fashion show at Victoria Secret's Website. The hardware also would need to prevent the attacker community from operating under the cover of these normal events.

18.9.5.2 Client Puzzle Protocol. RSA Laboratories has proposed cryptographic methods as a potential defense to DDoS attacks against Web servers. This approach would use a *client puzzle* protocol designed to allow servers to accept connection requests from legitimate clients and block those from attackers. A client puzzle is a cryptographic problem that is generated in such a way as to be dependent on time and information unique to the server and client request.

Under normal conditions, a server accepts any connection request from any client. If an attack is detected, the server selectively accepts connection requests by responding to each request with a puzzle. The server allocates the resources necessary to support a connection only to those clients that respond correctly to the puzzle within some regular TCP time-out period. A bona fide client will experience only a modest delay getting a connection during an attack, while the attacker will expend an incredible amount of processing power to continue sending the number of requests necessary for a noticeable interruption in service at the target site, quickly rendering the attack ineffective (in effect, a reverse DoS). This scheme might be effective against a DDoS attack from a relatively small number of hosts each sending a high volume of packets but might have limited effectiveness against a low-volume attack from a large number of systems or from a botnet.

18.9.5.3 IP Traceback. A third mechanism that was the subject of considerable research during the early to mid-2000s was *IP Traceback*. The problem with DoS/DDoS attacks is that packets come from a large number of sources, and IP address spoofing masks those sources. Traceback marking, in concept, is a relatively straightforward idea. Every packet on the Internet goes through some number of ISP routers. Processing power, memory, and storage are available for routers to mark packets with partial path information as they arrive. Since DoS/DDoS attacks generally comprise a large number

18 · 30 DENIAL-OF-SERVICE ATTACKS

of packets, the traceback mechanism does not need to mark every packet, but only a sample size that is statistically likely to include attack packets (e.g., 1 packet out of every 20,000, or 0.005% of the IP traffic). This feature would allow the victim to locate the approximate source of the attack without the aid of outside agencies and even after the attack had ended. Another traceback proposal would define an ICMP Traceback message that would be sent to the victim site containing partial route information about the sampled packet. There are many issues related to traceback that need to be resolved, such as the minimum number of marked packets required to reconstruct the path back to the attacker, the actual processing overhead, and the ability to perform traceback while an attack is under way. In addition, any traceback solution will require a change to tens of thousands of routers in the Internet; how effective can traceback be during a period of gradual implementation? The upside, of course, is that the solution is backward compatible and results in no negative effects on users.

18.9.5.4 Host Identity Payload. A fourth proposal was to modify IP to be less prone to address spoofing by making the protocol less dependent on the address field for anything more than routing. The Host Identity Payload (HIP), for example, defines a protocol for the exchange of a cryptographic *Host Identity* between two communicating systems. This feature relegates the IP address for use solely as a mechanism for packet forwarding rather than as an identifier of the sender. Sender identification, instead, is accomplished by the Host Identity value, and all of the higher layer protocols are bound to the Host Identity. HIP is not yet widely deployed but is available in some TCP/IP implementations.

18.10 MANAGEMENT ISSUES. One of the greatest shortcomings in many organizations is that the highest levels of management do not truly understand the critical role that computers, networks, information, and the Internet play in the life of the organization. It is difficult to explain that there is an intruder community that is actively working on new tools all the time; and history has shown that as the tools mature and become more sophisticated, the technical knowledge required of the potential attacker goes down and the number of attacks overall goes up. Too many companies insist that “no one would bother us” without realizing that *any site* can become a target just by being there.

DoS attacks come in a variety of forms and aim at a variety of services, causing increased complexity and difficulty for system defense. DoS attacks should be taken seriously because of the potential threat they present, and attempts should be made to educate operational staff before such attacks occur, to document DoS attacks if they do occur, and to review the documentation and actions taken after the incident is over. Discussion of what steps were taken, what actions went into effect, and what the overall result was will help in determining whether the procedures carried out and techniques utilized were those best suited to the situation. A frank review and discussion will help achieve the best, most rapid, and most effective deployment of resources.

If anything proves the intertwined nature of the Internet, it is the defense against DDoS attacks. DDoS attacks require the subversion and coordination of hundreds or thousands of computers to attack a few victims. Defense against DDoS attacks requires the cooperation of thousands of ISPs and customer networks. Fighting DDoS requires continued diligence in locking down all of the hosts connected to the Internet as well as fundamental changes in the nature of TCP/IP connection protocols. In addition, many of the same techniques used to insert viruses and worms that lead to DoS and DDoS attacks are being used to insert Advanced Persistent Threat (APT)-class malware.

FURTHER READING 18 · 31

As with most security concerns on the Internet, a large part of the solution lies in user education. Users must be educated as to what current threats are currently on the Internet and affecting relevant applications, what to do if something suspicious occurs, and how to respond to attackers that are increasingly creative. The world of everyone having their own computers, mobile devices with an incredible amount of intelligence, and a bring-your-own-device (BYOD) corporate atmosphere requires that every user in the enterprise is a potential vulnerable point and needs to be properly educated.

18.11 FURTHER READING

- Anonymous. *Maximum Security*, 4th ed. Indianapolis: SAMS, 2003.
- CERT Coordination Center. *Results of the Distributed-Systems Intruder Tools Workshop*. CERT Website, December 7, 1999. www.cert.org/reports/dsit-workshop.pdf
- Denning, D. E. *Information Warfare and Security*. Reading, MA: Addison-Wesley, 1999.
- Ferguson, P., and D. Senie. “Network Ingress Filtering: Defeating Denial of Service Attacks which Employ IP Source Address Spoofing” (RFC 2827/BCP 38). The Internet Society, May 2000. www.ietf.org/rfc/rfc2827.txt
- Gao, Z., and N. Ansari. “Tracing Cyber Attacks from the Practical Perspective.” *IEEE Communications Magazine* 43, no. 5 (May 2005): 123–131.
- Gibson, S. “The Strange Tale of the Denial of Service Attacks Against GRC.COM.” Computer Crime Research Center Website, June 2, 2001. www.crime-research.org/library/grcdos.pdf
- McClure, S., J. Scambray, and G. Kurtz. *Hacking Exposed, Network Security Secrets & Solutions*, 6th ed. Berkeley, CA: Osborne/McGraw-Hill, 2009.
- Mirkovic, Jelena, Sven Dietrich, David Dittrich, and Peter Reiher. *Internet Denial of Service: Attack and Defense Mechanisms*. Upper Saddle River, NJ: Prentice-Hall, 2005.
- Raghavan, S. V., and E. Dawson, eds. *An Investigation into the Detection and Mitigation of Denial of Service (DoS) Attacks: Critical Infrastructure Protection*. Springer, 2011.
- Rekhter, Y., B. Moskowitz, D. Karrenberg, G. J. de Groot, and E. Lear. “Address Allocation for Private Internets” (RFC 1918/BCP 5). The Internet Society, February 1996. www.ietf.org/rfc/rfc1918.txt
- Skoudis, E. *Counter Hack Reloaded: A Step-by-Step Guide to Computer Attacks and Effective Defenses*, 2nd ed. Upper Saddle River, NJ: Prentice-Hall, 2006.
- Spafford, E. H. “The Internet Worm: An Analysis.” *Computer Communication Review* 19, no. 1 (January 1989): 17–57.
- Most of the best current references are on the Web. Some of the sites that readers should monitor for DoS/DDoS information are:
- ATLAS Summary Report, Global Denial of Service: <http://atlas.arbor.net/summary/dos>
CERT/CC: www.cert.org
- Dave Dittrich’s “Distributed Denial of Service DDoS Attacks/tools” Web page: <http://staff.washington.edu/dittrich/misc/ddos/>
- Denialinfo.com’s “Denial of Service (DoS) Attack Resources” Web page: www.denialinfo.com/dos.html
- DShield: www.dshield.org
- SANS Internet Storm Center: <http://isc.sans.org>
- US-CERT: www.us-cert.gov

18 · 32 DENIAL-OF-SERVICE ATTACKS

18.12 NOTES

1. DHS National Cyber Security Division, US-CERT, “National Vulnerability Database,” Website, 2013, <http://web.nvd.nist.gov/view/vuln/search?execution=e2s1>
2. Formerly Internet Information Server.
3. $10\text{ KB} \times 8\text{ bits/byte} \times 50\text{ hosts} = 4\text{ Mb}$.
4. Michael Calce and Craig Silverman, *Mafiaboy: A Portrait of the Hacker as a Young Man*, Lyons Press, 2011.
5. Brian Bloom, “The New DDoS: Silent, Organized, and Profitable,” *PCWorld*, May 29, 2012, www.pcworld.com/article/256431/the_new_ddos_silent_organized_and_profitable.html

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 19

SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

**Karthik Raman, Susan Baumes,
Kevin Beets, and Carl Ness**

19.1 INTRODUCTION	19·1	19.5.1 Consequences	19·14
19.2 BACKGROUND AND HISTORY	19·2	19.5.2 Case Studies—Examples from Business	19·15
19.3 SOCIAL-ENGINEERING METHODS	19·4	19.5.3 Success Rate	19·16
19.3.1 Impersonation	19·4	19.5.4 Small Business versus Large Organizations	19·16
19.3.2 Seduction	19·5	19.5.5 Trends	19·17
19.3.3 Intimidation	19·6		
19.3.4 Low-Tech Attacks	19·6	19.6 DETECTION	19·17
19.3.5 Network and Voice Methods	19·8	19.6.1 People	19·17
19.3.6 Reverse Social Engineering	19·11	19.6.2 Audit Controls and Event Logging	19·18
		19.6.3 Technology for Detection	19·18
19.4 THE PSYCHOLOGY OF SOCIAL ENGINEERING	19·11	19.7 RESPONSE	19·19
19.4.1 Introduction	19·11	19.8 DEFENSE AND MITIGATION	19·19
19.4.2 Psychology	19·12	19.8.1 Training and Awareness	19·20
19.4.3 Social Psychology	19·13	19.8.2 Technology for Prevention	19·20
19.4.4 The Social Engineer Profile	19·14	19.8.3 Physical Security	19·21
19.5 DANGERS OF SOCIAL ENGINEERING AND ITS IMPACT ON BUSINESSES	19·14	19.9 CONCLUDING REMARKS	19·21
		19.10 NOTES	19·22

19.1 INTRODUCTION. According to Greek mythology, the Greeks defeated the Trojans in the Trojan War with the help of a wooden statue. After fighting a decade-long war in vain, the Greeks withdrew from their stronghold on the beach. Outside the gates of Troy, they left a giant wooden horse. The statue confused the Trojan soldiers, but it was brought within the fortified walls of Troy. Inside the statue hid several Greek soldiers. When darkness fell, these soldiers emerged from the statue and opened the

19 · 2 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

gates of Troy. The Greek army entered Troy and took the soldiers and citizens of Troy by surprise. After the deception, Greece won the war quickly.

The Trojan horse built by the Greeks was effective because it used deception to achieve the desired result: penetrating the enemy's established defenses. The Trojan horse accomplished what those in information security label *social engineering*.

Social engineering may be defined as obtaining information or resources from victims using coercion or deceit.

Social engineering refers to lying, cheating, tricking, seducing, extorting, intimidating, and even threatening employees into revealing confidential information that can then be used to break into systems. Social engineering is based on deception and on violation of social norms of fairness and honesty.¹

During a social-engineering attack, attackers do not depend on subverting technology, for example, scanning networks, cracking passwords using brute force, or exploiting software vulnerabilities. Rather, social engineers operate in the social world by manipulating the trust or gullibility of human beings. Thus, social engineering relies on manipulating human nature to extract information or access by tricking a victim. Related to social engineering are low-tech attacks, and the two often go hand-in-hand. Low-tech attacks are similar to social-engineering attacks in that they do not utilize technology. They are physical attacks performed against companies or individuals' property.

Not all social-engineering and low-tech attacks will give attackers all the information they are seeking at once. Social engineers will collect small pieces of information that seem innocuous to the individuals that divulge them. Social engineers may gather these snippets of information in a seemingly random order but then assemble them into intelligence that is used to launch larger attacks that can be devastating to an organization's information security, resources, finances, reputation, or competitive advantage. Indeed, the purpose of a social-engineering attack can be as varied as the attack method employed. The result, however, is generally the same: a loss of intellectual property, money, business advantage, credibility, or all of the above.

Social-engineering attacks have been and will remain successful due to the weakest security link in an organization: the people. They are important because they exploit human nature, which is immutable and therefore perpetually vulnerable.

This chapter presents the history of social-engineering and low-tech attacks, its methods, the social science behind it, and its business impact. In addition, it covers detection and mitigation policies for managers and information security officers to defend against and mitigate social-engineering and low-tech attacks.

The purpose of a social-engineering attack can be as diverse as the attack method. Nevertheless, the result to the victims is generally the same, loss of intellectual property, money or business, credibility, or all of the previous. This chapter presents the history of social-engineering and low-tech attacks, its methods, the social science behind it, and its business impact. In addition, it covers detection and mitigation policies for managers and information security officers to defend against social-engineering and low-tech attacks.

19.2 BACKGROUND AND HISTORY. Social engineering is not a new tactic, nor is it an invention of modern-day hackers. The term has its foundations in political history, where a person or group manipulates a group of people, large or small, in an attempt to persuade or manipulate social attitudes or beliefs. Often, governments or political parties engage in this practice. The modern term's origins date back before World War II began, and the Nazis are thought to have actually coined the term.² To this day,

BACKGROUND AND HISTORY 19 · 3

the term carries a negative connotation because of its roots in history, especially during Nazi-controlled Germany. Some researchers also consider the realm of social engineering to cover everything from advertising and modern media to political action groups.

Social engineering today is better known for its use as a security-penetration technique.

Deception has been integral to espionage through the centuries. Sun Tzu wrote in the fifth century BCE,

- “All warfare is deception.”
- “Having doomed spies, doing certain things openly for purposes of deception, and allowing our spies to know of them and report them to the enemy.”³

Deception—specifically, feeding false information to the enemy’s spies—was an integral part of the Second Punic War between Rome and Carthage in the third century BCE.⁴ During the Second World War, the Allies used “Operation Bodyguard” to deceive the Axis powers into believing that the D-Day invasion would occur at a different time and place from its true schedule and target.⁵ Deception was integral to the spy vs. counter-spy competition between the Communist block and the West during the Cold War.⁶

A well-known example of social engineering is the escapades of Frank W. Abagnale, subject of the fictionalized movie *Catch Me If You Can*.⁷ Abagnale was able to successfully impersonate authority figures, including a physician, pilot, attorney, and teacher, even though he was a teenager. He also used social-engineering techniques to persuade and manipulate innocent and good-natured individuals to help him carry out many of his frauds. Many of Abagnale’s techniques were highly successful and well engineered. Today, he helps organizations recognize and defend against such attacks through his speaking engagements and consulting business.⁸

One of the best-known social engineers is Kevin Mitnick. Though Mitnick is now a computer security consultant, lecturer, and author, his present career was preceded by a number of years spent as a computer hacker and social engineer. Long reformed, Mitnick has written several books discussing his observations and techniques as a computer hacker. Mitnick maintains that social engineering is the most powerful tool in the hacker’s toolbox.

Ever since social engineering has become a popular, and more importantly, successful, technique, its frequency of use has increased. One of the most visible reminders of this is the large number of phishing and pharming attacks, discussed in detail in Chapter 20 in this *Handbook*.

Social engineering has even become a spectator sport at some criminal-hacker conventions. At the 2012 DEF CON in Las Vegas, Shane MacDougall won the Social-Engineering Capture-the-Flag Contest by tricking a Walmart executive in a Canadian store into revealing every single piece of confidential information listed in the contest objectives.

Darnell asked the manager about all of his store’s physical logistics: its janitorial contractor, cafeteria food-services provider, employee pay cycle, and staff shift schedules. He learned what time the managers take their breaks and where they usually go for lunch.

Keeping up a steady patter about the new project and life in Bentonville, Darnell got the manager to give up some key details about the type of PC he used. Darnell quickly found out the make and version numbers of the computer’s operating system, Web browser, and antivirus software.⁹

19 · 4 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

Social engineering, when used by itself as a standalone attack tool, is very effective. But this tool can be used as part of a larger attack, or even as part of a technical attack. Social engineering is commonly used at the beginning of a larger, more substantial attack. It is important to remember that some attacks have taken place over weeks, months, or years. This is not always the case, and just because a specific social-engineering attempt was averted, the larger attack may not be completely mitigated. An example of this may be that an initial attack was unsuccessful because the attacker did not have sufficient or accurate information to carry out the attack. The attacker may revert once again to social engineering and low-tech attack strategies to gather information for a different, possibly successful, technical attack.

Even though an organization may have the best firewall, intrusion detection system, and risk management tools in the world, an attacker may use social engineering to circumvent these technical defenses. The attacker will not be stopped by the best technical defenses if the attacker is able to extract a valid user name and password from an unsuspecting employee. Once the attacker has this information, it may be enough information to carry out a massive attack on an organization's information systems—possibly undetected.

Computer criminals have been able to use pure social engineering and low-tech attack techniques, without relying on any serious technology, to cause large-scale damage to an organization as well. An action as simple as taking discarded, potentially damaging, information out of a dumpster and sending it to local media can cause substantial damage to an organization's image. This act might be part of a larger campaign or operation by a group of individuals seeking political or social action against an organization. An example of this was recently chronicled by Home Box Office (HBO) involving a group of citizens who routinely gathered the trash from county election centers and city halls in an effort to collect evidence that they could present to the media to prove voter fraud. Documents that were improperly disposed of have had extremely negative consequences for election officials.¹⁰

19.3 SOCIAL-ENGINEERING METHODS. Social-engineering attacks can take many different forms and expert social engineers are capable of changing their methods of attack very quickly in order to succeed. The underlying principle of most attacks is *pretexting*, which is defined as “the collection of information ... under false pretenses.”¹¹ Two distinct underlying methodologies are used during social-engineering attacks, impersonation and seduction. The basis for most attacks is either of these two methods. Certain low-tech attacks, which do not involve any human contact, are exceptions.

Targets of social-engineering attacks also vary widely. Depending on how complex the attack is or how much knowledge the attacker has will help determine the target for the attack. However, in many instances, a social-engineering attack is an attack of opportunity and the victim will be randomly chosen. In certain well-planned attacks, a particular target may be identified as the victim.

19.3.1 Impersonation. Impersonation is defined as pretending to be someone else, and it is one of the most popular methods that social engineers employ. Social engineers may use cross-functional impersonation attacks to target an employee at any level within the targeted organization. They may pretend to be a real, named individual, or they may pretend to have a particular role or authority.

Corporate executives are common targets of social engineering.¹² Helpdesk employees and systems administrators are also common targets of impersonation. Most

SOCIAL-ENGINEERING METHODS 19 · 5

organizations have helpdesks for assisting employees with issues related to information technology (IT). Employees generally will follow the instructions from helpdesk personnel simply because they are perceived as technologically well-informed and trustworthy. Social engineers understand this trust and will exploit it to steal information. The attacker will impersonate helpdesk personnel and abuse this blind trust to gather whatever information the attacker needs.

Helpdesk personnel can also be the victims of social-engineering attacks with the social engineer impersonating a user in need of technical assistance. An attack against the AOL helpdesk is a historical example of a successful social-engineering attack against helpdesk personnel. An attacker posing as a customer was able to infiltrate the AOL environment by conning an AOL helpdesk staff member.¹³

There have also been instances where social engineers impersonate corporate officers or managers. In one case, the payroll giant ADP released personnel and brokerage account information on hundreds of thousands of customers. The attacker impersonated a corporate board member, requested and received the information.¹⁴ This case highlights the need for controls and education for all levels of employees. The victim in this case might have been concerned about not pleasing a powerful member of the organization and may have feared retribution.

It is important to note that attackers will also impersonate regular employees of an organization by dressing, speaking, and blending into the organization's environment. In doing so, the attacker may gain physical access to a facility by piggybacking or tailgating. In one case reported by a student in a 1993 security course, a social engineer collaborating with a corrupt manager was able to commandeer an unused desk in a large company for several *months* before a security guard asked for his company identification. The criminal disappeared immediately, but not before ingratiating himself to the employees as a new employee and even being invited out to social events by the unsuspecting victims.¹⁵

For a social engineer, there are no boundaries for impersonation; they may try to gain information via physical access to an organization using any number of ruses including impersonation of:

- Temporary employees (e.g., contractors or auditors)
- Utility or telecommunications company employees
- Emergency personnel
- Janitorial or maintenance employees
- New employees
- Delivery personnel

Physical access can greatly increase the success rate of an attack depending on whether the organization has proper controls in place. Highly motivated social engineers may go so far as to gain employment at the company, or at a company that is a client of the company, in order to have easier access to the victim.

19.3.2 Seduction. The word seduce is defined as “to lead away from duty, accepted principles, or proper conduct”.¹⁶ In general, a social-engineering attack using seduction will take longer to complete than an impersonation attack. The attacker, using seduction, will identify a target and will form a bond with that individual, through social settings, online, or through another mechanism. It is during this relationship that the victim’s information is divulged to the attacker.

19 · 6 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

For example, a social engineer who wishes to gain access to a building may befriend a security guard of that organization. After some time has passed and the relationship has progressed, the attacker may request a tour of the facility. The security guard, wanting to please the new friend, may allow a tour. The social engineer, once inside, can plant clandestine listening devices, look for usernames or passwords, and read documents left in the open.

19.3.3 Intimidation. Social engineers can use fear to achieve their ends. For example, in 2012, scammers spread alarming emails about FB I prosecutions to innocent victims.

A new extortion technique is being deployed by cybercriminals using the Citadel malware platform to deliver Reveton ransomware. The latest version of the ransomware uses the name of the Internet Crime Complaint Center to frighten victims into sending money to the perpetrators. In addition to instilling a fear of prosecution, this version of the malware also claims that the user's computer activity is being recorded using audio, video, and other devices.

As described in prior alerts on this malware, it lures the victim to a drive-by download Website, at which time the ransomware is installed on the user's computer. Once installed, the computer freezes and a screen is displayed warning the user they have violated United States Federal Law. The message further declares that a law enforcement agency has determined that a computer using the victim's IP address has accessed child pornography and other illegal content.

To unlock the computer, the user is instructed to pay a fine using prepaid money card services. The geographic location of the user's PC determines what payment services are offered. In addition to the ransomware, the Citadel malware continues to operate on the compromised computer and can be used to commit online banking and credit card fraud.¹⁷

Another social-engineering attack using intimidation is the phony office-supply scam.¹⁸ A low-level employee accepts an offer of free printer toner, glass cleaner, or paper and receives a box of materials—followed by an invoice. Attempts to return the unwanted, no-longer-free materials are refused and the criminals use increasingly strong threats to intimidate the employee into passing the invoice along for payment. Employees may be told that their employer will be sued or that a complaint about their behavior will be sent to upper echelons of the organization. A variation involves sending a gift to the victim in the hope that accepting it can be used to embarrass the recipient into complying with the scammers' demands. Such scams succeed often enough to generate millions of dollars of revenue for the criminals.¹⁹

19.3.4 Low-Tech Attacks. Low-tech attacks are invaluable to an attacker as part of reconnaissance, information theft, and surveying for easy targets. On occasion, these methods may even reward the attacker with a very large amount of useful information in a very short amount of time. Low-tech attack methods may seem simple or improbable, but the methods described can easily be overlooked by security managers. They are not urban legends—they have been used in the past and continue to be utilized today.

19.3.4.1 Dumpster® Diving. In the context of social engineering, *Dumpster® diving* is the social engineer's act of searching through an organization's garbage in an attempt to find documents, hardware, software, or anything that could be of value to meet the goals of the attacker. Even with the widening use of sensitive document destruction, *Dumpster®* diving is a popular social-engineering technique because it is easy and often successful. Oracle hired detectives to purchase Microsoft's trash during

SOCIAL-ENGINEERING METHODS 19 · 7

Microsoft's antitrust trial. (The detectives were unsuccessful.)²⁰ Social engineers do not need to deceive anyone to perform the attack. In many cases, the materials disposed may sit in open containers for weeks. Dumpster® diving is most often carried out at night when no one is around, as there is less risk of being caught. So as to not draw attention to themselves, Dumpster® divers have been known to dress in dark clothing or even use janitorial uniforms.

All organizations must understand the legal ramifications of Dumpster® diving. Local and state laws may vary widely and be murky; however, no organization or individual should have *any* realistic expectation of privacy relating to materials in refuse containers. It may even be legal for an attacker to remove and take ownership of anything left in a garbage receptacle if it is placed outside the limits of private property. Confidential information must be destroyed, not simply discarded.

An unfortunate example of making confidential data accessible in garbage came to light in November 2012, when strips of confetti thrown into the crowd at the Macy's Thanksgiving parade were discovered to have intact, readable confidential data:

Among the easily identifiable records from the Nassau County Police Department were what appears to be details of Mitt Romney's motorcade route to and from the final presidential debate at Hofstra University.

Confetti collected by spectators near 65th Street and Central Park West also contained arrest records, incident reports, and personal information, and identified undercover officers, WPIX-TV says.

"There are phone numbers, addresses, more Social Security numbers, license plate numbers," said Ethan Finkelstein, 18, of Manhattan, who gathered up some of the confetti with friends. "And then we find all these incident reports from police."²¹

The papers had been shredded the wrong way in a single-cut shredder, leaving horizontal strips of readable data for anyone to see.

19.3.4.2 Theft. The age-old crime of theft is another popular social-engineering technique. Social engineers may pick up anything they can get their hands on, literally, and leverage the information obtained to carry out other attacks. Targets of theft include, but are not limited to, printed materials, CD-ROMs, USB flash drives, backup media, tablets, smart devices, and laptops. Thieves may obtain objects on or off company premises. While on company premises, they may look for objects they can grab and quickly conceal. An attacker may bring in empty laptop cases, backpacks, or even large purses to aid in their efforts. Most employees are more likely aware of theft techniques outside of the organization, such as taxicabs, airports, and other public places. Employees may not realize that a social engineer or other criminal may also be an insider.

19.3.4.3 Leveraging Social Settings. Social engineers may use social settings to gain information because people relax in social settings and may believe that information security practices are for the workplace only. A social engineer may use a social setting, such as a bar, to take advantage of drinking employees to gain information. The attacker may actively engage a target or passively eavesdrop on a conversation. While this type of attack may seem far-fetched, there are many situations where the attacker may be in the right place at the right time to gain knowledge that can be later used as part of a larger attack. People in social settings are less likely to have their defenses up and security on their minds.

19 · 8 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

Restaurants, corporate functions, impromptu meetings outside of a company building, or loud phone conversations are all areas or situations where eavesdropping can occur. Many times, employees will work on commuter trains or conduct business in other public areas. A social engineer can exploit an organization's mobile workforce to gain information. Employees should be cognizant of who is around them while performing any work-related task.

19.3.4.4 Exploiting Curiosity or Naïveté. Social engineers may trick a victim into unknowingly aiding in an attack by piquing the user's curiosity. For example, an attacker may leave an intriguingly labeled CD-ROM in a break room, hoping that a victim is curious about the contents. When the victim places the CD-ROM in his or her computer, a malicious program, such as a virus, could automatically execute and spread. This technique has also been executed using USB drives, iPods, and corrupted music CDs.²²

19.3.4.5 Data Mining and Data Grinding. Search engines can catalogue a surprising amount of information that is sensitive and confidential. Social engineers can create special searches, use search engine application programming interfaces (APIs), and use advanced search capabilities of many search engines to mine information about a company. Another attack vector is the caching feature of many search engines. A search engine may cache a Web page with sensitive information. If an organization requests the cache be removed, still there is a delay before the search engine cache is updated, during which time the organization's information is exposed. Social engineers could also exploit data about individuals or organizations gathered from social media. See Section 19.3.5.6, "Social Media," in this *Handbook* for details about the use of this method. Attackers could mine this data for victims' behavioral trends. They could garner more information about the victims than was readily presented by social media data.

Documents published by a company are another source of unintended information disclosure. In a technique known as data grinding, social engineers can use metadatareader software to extract information such as the author's name, organization, computer name, network name, email address, user identification (ID), and comments from Microsoft Office documents.²³ This potentially damaging information is included in most document types.

19.3.4.6 Piggybacking or Tailgating. A common and very successful method of social engineering is *piggybacking* or *tailgating*. The method allows an attacker access to a facility without showing proper credentials. The victim in these cases is being polite and holding the door for the attacker who enters the facility using the credentials of the victim. Once inside a facility, an attacker is free to roam the building in search for information. If questioned upon entering, attackers may use the excuse of forgetting their credentials or being new employees. In general, if a piggybacking is going to be attempted, social engineers will dress and act according to other members of the organization so they can blend into the facility. If a facility is secured by multiple layers of physical security, attackers may try to use social engineering to breach the first few layers then proceed to use other attack methods.

19.3.5 Network and Voice Methods. The Internet has taken off, and so have social-engineering attacks. These attacks differ from traditional social-engineering attacks due to minimal or no human interaction. Still, the basis of the attacks still relies

SOCIAL-ENGINEERING METHODS 19 · 9

on the trusting nature of humans. The annoying Nigerian 419 attack is a combination attack using both email and human interaction.²⁴

The methods used for these attacks include *phishing*, *pharming*, *spim*, *malware*, and *vishing*. Phishing and Trojans are discussed in detail in Chapter 20 in this *Handbook*.

19.3.5.1 Phishing. Phishing is one of the most widely used and successful social-engineering attacks. It is defined as the “act of sending an email to a user falsely claiming to be an established legitimate enterprise in an attempt to scam the user into surrendering private information.”²⁵ Phishing, as with all social-engineering attacks, relies on the trusting nature of people. Naïveté about using the Internet also plays a role in the success of this attack.

An interesting new wrinkle is that being *misrepresented as the origin* of phishing emails was the second most common type of attack reported by respondents (39 percent) to the Computer Security Institute’s “2010/2011 Computer Crime and Security Survey.”²⁶

Newer, more targeted phishing attacks called *spear phishing* are also gaining in popularity. Spear phishing is used in a dangerous type of threats known as Advanced Persistent Threats (APTs). Spear phishing was used in a portion of a breach of the security firm RSA in 2011.²⁷

19.3.5.2 Spim. Many people and businesses rely on the synchronous communication that instant messaging (IM) offers. Social engineers have noted the increase in the use of IM software and developed *spim*. Spim is “instant spam or IM spam.”²⁸ A spim attack is very similar to a phishing attack except the vector is IM software instead of email. For example, an attacker will develop a fraudulent Website that resembles a legitimate one and send its link to many IM accounts. Victims will visit the Website and login, revealing their credentials to the attacker. The fraudulent site could contain malicious software that infects the victim’s computer. Surprisingly, despite the increasing use of IM software, the growth of spim has been slow.²⁹

19.3.5.3 Pharming. In *pharming* attacks, attackers attempt to make victims visit spoofed Websites to reveal sensitive personal information. Attackers achieve this by manipulating the victim’s local or global DNS directory. Pharming attacks may fool users more easily than other low-tech attacks because the user may be given no indication that an attack is under way. Users may type in the URLs of their banking or credit card Websites into their browsers as usual and not normally notice that they are in fact visiting fraudulent sites. In 1997, Dan Parisi registered the domain whitehouse.com and snared people mistakenly looking for whitehouse.gov into viewing pornography—an early form of pharming.³⁰ In a 2005 pharming attack, users of the Internet mail service Hushmail were redirected to a fraudulent site where their information was harvested. Hushmail suffered negative publicity from this attack and was forced to update its users daily about its investigation into the attack.³¹

19.3.5.4 Malware. Malware are programs or files that are harmful or dangerous to the end user. Social engineers frequently use Trojans and viruses. Although the programs themselves do not make up a pure social-engineering attack, the basis of the attack still relies on manipulating a victim’s trust. For example, a social engineer may send a victim a phishing email with a link to a malicious Website. If the user visits the malicious Website, a Trojan is installed on the victim’s computer and the malicious program will begin gathering the victim’s information. Another example involves a

19 · 10 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

USB drive loaded with auto-executing malware. A victim curious about the contents will plug the drive into a machine and execute the malicious code.³² Any number of vectors can be used, including documents, email messages, Websites, CDs, and USB drives. The premise is the same for all types of these attacks: The unsuspecting users, curious about the content, will inadvertently install these dangerous programs. Malware infections were the most common type of attack reported by respondents (67 percent) to the Computer Security Institute's "2010/2011 Computer Crime and Security Survey," and bots or zombies were reported by 29 percent.³³

Malware has recently begun to use a type of social-engineering attack whereby the trust placed in a valid certificate is exploited in conjunction with a long-known vulnerability called Data Link Library (DLL) preloading. In this attack, a package is prepared that includes a valid, trusted application that has a valid signature. In this package, however, a malicious DLL is included. Most users, who may know to look for a valid certificate, will indeed find one in the application and therefore trust the entire package. The victim launches the valid, trusted application, and due to the preloading vulnerability, the malicious DLL is loaded from the local directory before the valid DLL is loaded from its regular path.³⁴

See Chapters 16, 17, and 18 in this *Handbook* for details of malware.

19.3.5.5 Vishing. Attackers who lure victims with the use of email and the telephone or just the telephone are performing a *vishing* attack.³⁵ Social engineers may send an email or call victims in an organization about an issue and request a call back. The number given is either staffed by accomplices or answered by a legitimate-sounding automated system. Victims are prompted to release potentially damaging information, about themselves in the case of identity theft, or about their company.

In another example, an attacker can leverage automated phone answering systems which allow the caller to input the first few letters of the last name of a contact to reach their extension. The attacker can try multiple extensions and stumble on some that reveal details regarding a person's position, title, or office status (e.g., "I'm out on vacation until April 4"). A social engineer can leverage that information and call other employees to gain additional information or access.

It should also be noted that with the proliferation of Voice over Internet Protocol (IP) (VoIP) there is now Voice Phishing, which is propagating a vishing-phishing attack. Since the call is routed over IP, the call number can be changed easily. This is analogous to how phishing Websites operate. Since this is primarily a voice-based social-engineering attack, it is considered a voice attack instead of a Web attack.

VoIP has also internationalized the occurrence of such frauds. In an incident in 2013, criminals thought to be in Africa who were advertising nonexistent free kittens in thousands of pages on the Web used VoIP to contact victims in an attempt to fool well-meaning cat-lovers into paying for imaginary air-shipments of imaginary cats.³⁶

19.3.5.6 Social Media. Social media have exploded in popularity in recent years, and so have the opportunities for the social engineer to victimize that frequent these outlets.³⁷

Useful details may be easily harvested just from searching the victim's social media page, or the pages of relatives who include the victim on their pages. Data harvested could include things such as:

- Event dates
- Friends/relatives/colleagues of interest

THE PSYCHOLOGY OF SOCIAL ENGINEERING 19 · 11

- Pictures
- Contact information
- Locations frequented

In addition to strictly harvesting data, a social engineer could use impersonation to befriend the individual, so that details hidden from public view are viewable. This may allow the attacker to target his attack to his needs more effectively.

Of particular concern is the LinkedIn network, which reached 187 million members worldwide by the end of 2012; 57 percent were located in the United States and 63 percent elsewhere.³⁸ LinkedIn makes it possible to track job titles, job history, technical background, education, and linkages to other people in any given organization—ideal information for social engineers to commandeer. Although details are supposed to be kept away from people who have not linked to each other, some LinkedIn users routinely accept link requests from anyone who sends them an invitation—even though LinkedIn’s User Agreement explicitly states that members must not “Invite people you do not know to join your network.”³⁹ Criminals have demonstrated their interest in LinkedIn contacts by selling catalogs of information gleaned on this network.⁴⁰

19.3.6 Reverse Social Engineering. Reverse social engineering is an effective attack usually executed by an experienced social engineer. A reverse social-engineering attack has three distinct parts. First, a social engineer will create a problem, for example, a user ID issue. Second, the social engineer will publicize that they are the only person capable of fixing the issue. In the final part of the attack, the social engineer will assist the victim and “fix” the issue. It is during the third segment of the attack that a social engineer will gather information. The success rate of reverse social-engineering attacks tends to be high simply because the victim is satisfied that the fabricated problem is fixed.⁴¹

For example, the attacker may change the name of a file (the problem). The victim searches for and cannot find the file. The attacker will announce that they have been able to retrieve lost information but will require a user ID and password to gain access to the system (the publicity). The victim, flustered by the thought of losing an important document, will divulge the information. Finally, the attacker will “find” the missing file (the fix). The victim, being pleased that the file is returned, will forget having let out access credentials to the system.⁴²

19.4 THE PSYCHOLOGY OF SOCIAL ENGINEERING

19.4.1 Introduction. Well-known security expert Bruce Schneier identified that behavioral economics, psychology of decision making, psychology of risk, and neuroscience can help explain why our feeling of security deviates from reality.⁴³ This section focuses on one of these aspects—psychology—and studies the science underlying the success of social engineering. Section 19.4.2 uses some well-established principles of psychology and social psychology to analyze social engineering from two angles, the psychological perspective of the victim and social-psychological perspective of the social engineer and victim. Section 19.4.3 explains that there is no single social-engineer stereotype. The terms used in this section can be found in undergraduate psychology and social psychology textbooks and therefore academic references have been minimized.

19 · 12 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

Social engineering succeeds because of human nature: In this section, examples of social-engineering attacks illustrate scientific terms used to characterize social engineering.

19.4.2 Psychology. A *cognitive bias* is defined as a mental error caused by humans' simplified information-processing strategies.⁴⁴ People become victims of social-engineering attacks due to cognitive tendencies, also called "heuristics" or rules of thumb, which are inherent in all humans. These cognitive tendencies are useful, mostly. Psychologist Robert Cialdini writes:

We can't be expected to recognize and analyze all the aspects in each person, event, and situation we encounter in even one day. We haven't the time, energy, or capacity for it. Instead, we must very often use our stereotypes, our rules of thumb, to classify things according to a few key features and then to respond mindlessly when one or another of these trigger features is present.⁴⁵

Although cognitive biases are found in everyone, the universal presence of cognitive biases does not imply that they are impossible to counter.

Following are some cognitive biases that can explain why people fall prey to social-engineering attacks:

- **Choice-supportive bias.** People tend to remember an option they chose that had more positive aspects than negative aspects.⁴⁶ Information technology (IT) helpdesk operators may provide employee names, extensions, or both without verifying the identity of callers. Helpdesk operators remember this practice as being good because most callers are genuine and the callers thank them for providing the information. Social engineers can masquerade as genuine callers to exploit helpdesk operators' choice-support bias. According to an article on security analysis site SecurityFocus, a demonstration of this attack by the Computer Security Institute actually succeeded.⁴⁷
- **Confirmation bias.** People tend to collect and interpret evidence in a way that confirms their conceptions.⁴⁸ If an organization has a contract with a custodial service and employees see custodians all wearing the same uniform, then a social engineer wearing the uniform may not be challenged to identify himself because of the employees' confirmation bias. The movie *Ocean's Eleven* contains many scenes in which the owner and employees of three Las Vegas casinos are conned by picaresque protagonists dressed to blend into their roles and environment.⁴⁹
- **Exposure effect.** People tend to like things that are familiar.⁵⁰ A social engineer may call victims and explain they are performing a survey for a popular local restaurant and ask about the organization by which the victim is employed. The victims are comfortable providing that information because they are familiar with the restaurant. According to malware researcher Elodie Grandjean, some of the themes of baited emails, URLs, or instant messages used in social-engineering attacks include⁵¹:
 1. Pornographic links and images
 2. Using a female name in the sender field
 3. Political agendas, including solicitations for contributions in the name of a popular candidate

THE PSYCHOLOGY OF SOCIAL ENGINEERING 19 · 13

4. Fake emails for banks, online payment services, and other financial services. These request a confirmation or an update of login credentials or credit card information.
 5. Threatening emails, mentioning jail sanctions or jury-duty procedures
 6. Free games and screensavers containing a Trojan, or free antispyware tools, which are often rogue programs themselves
 7. Big events, such as sports, extreme weather disaster, or urgent news
 8. Celebrity names and reports on their adventures and misbehavior
 9. Potentially trusted or secret relationships such as affiliation with a 12-point program
 10. Social networking Websites, fake friends, school classmates or relatives, and secret lovers
- **Anchoring.** People tend to focus on one trait when making decisions.⁵² If a social engineer has a soothing voice, the victim may focus on that attribute versus the questions being asked. Security expert Winn Schwartau describes a successful social-engineering attack where some employees of a New York financial services company fell for a bogus letter written on company letterhead claiming to be from the information security department; 28 percent of the 1,200 employees wrote back with their personal details.⁵³ The employees had focused, to their detriment, too much on the letterhead and had ignored other aspects of the request, such as its suspiciousness, graveness, and suddenness.

19.4.3 Social Psychology. Social psychologists define the *schema* as the inherent picture of reality used by humans to make judgments and decisions. From the perspective of social psychology, social engineers exploit the fact that most peoples' schema includes rules to be trustful of other people and their intentions. People are taught from the very beginning of their socialization that being nice to others is a good thing. In the context of information security, peoples' tendency to blindly trust others can spell disaster.

Below is a list of common errors that people make, and examples of how social engineers will exploit those mistakes to attack an organization.

- **Fundamental attribution error:** In this common error, people assume that the behaviors of others reflect stable, internal characteristics. Someone committing the fundamental attribution error might see a colleague in a bad mood and think, "She is *always* moody." In reality, the colleague might be pleasant in general but be suffering a headache at the time. Social engineers will act pleasant and charming to lead victims to commit the fundamental attribution error, to be impressed that the attackers are nice people *in general* and so to help them.
- **Salience effect:** Given a group of individuals, people tend to guess that the most or least influential person is the one who stands out the most. For example, from a group of ten people, nine of whom are six feet tall and one who is five feet tall, if asked to guess who the most intelligent person in the group is, an observer might say that it is the five-foot-tall person. Social engineers attempt to blend into their victim's environment to take advantage of the salience effect. They are acutely aware of company lingo, events, and regional accents.

19 · 14 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

- **Conformity, compliance, and obedience:** People respond to the social pressures of conformity, compliance, and obedience by adjusting their behaviors. A social engineer impersonating a high-powered executive demanding admittance into the company premises may persuade a new security guard with the weight of his assumed authority. The authority figure's promise of reward or threat of punishment may further influence the security guard's decision to carry out the request of the attacker.

For more details of the implications of social psychology for information assurance, see Chapter 50 in this *Handbook*.

19.4.4 The Social Engineer Profile. The profile of a social engineer is not that of the stereotypical computer hacker often portrayed in the movies or on television.⁵⁴ Social engineers are most likely not going to be loner teenagers who spend all their time with their computers in a dark basement. A social engineer is often outgoing, confident, an excellent communicator, and well-educated. Social engineers may use their own personality, or adopt a persona that greatly differs from their normal personality—oftentimes a personality they have developed over months or even years. To achieve their goals, they will blend into the environment. Social engineers seek to be unnoticeable, unremarkable. They will dress according to the dress code of the environment in which they operate. Interestingly, social engineers may be excellent actors, able to think on their feet and to adapt quickly to changing conditions. The attacker's confidence will often mask any nervousness or tension during the social-engineering attempt, which may lead to unwarranted credibility.

Social engineers may also exhibit a dark side. Attackers may have very little regard for the consequences of their actions on the victim. Even though the attackers may appear very polite or congenial towards their victims, they actually care very little about the victim or the people utilized to achieve a goal. The victim and others are simply a means to an end; they are only part of the social-engineering attack tool. The social engineer's motivations may vary widely and range from personal financial gain to revenge. There may also be significant external pressure on the attacker from acquaintances or organized crime syndicates.⁵⁵

For more details of the psychology of computer criminals, see Chapters 12 and 13 in this *Handbook*.

19.5 DANGERS OF SOCIAL ENGINEERING AND ITS IMPACT ON BUSINESSES. The ultimate goals of the attacker may be limited to theft but may also include disruption of business or even destruction of a business. An organization must evaluate the potential impact of even a seemingly minor social-engineering attempt.

19.5.1 Consequences. The consequences of a successful social-engineering attack are almost innumerable. Anything one can dream up as an attack vector has probably been attempted by an attacker somewhere. Much like disaster recovery planning, when trying to quantify and understand the impact of social-engineering attacks, all possibilities must be put on the table. Something as seemingly minor as an internal memo that has not been properly destroyed could have the potential to bankrupt an organization. As a result, a simple social-engineering attack can lead to a more complex attack which morphs into a major information security breakdown simply based on the information contained in one misplaced memo.

DANGERS OF SOCIAL ENGINEERING AND ITS IMPACT ON BUSINESSES 19 · 15

The danger is especially high for publicly traded companies that can lose value because of a loss of confidence from investors.⁵⁶ Many companies, within the last few years, have come under financial duress because of a security breach or lapse that drew considerable attention from the press. Social engineering was likely part of these security incidents. Many organizations are required by law to have safeguards in place when it comes to data security; more and more of these requirements require security controls that will aid the organization in defending against social-engineering and low-tech attacks.

Another serious consequence of a successful social-engineering attack is the uncertainty that follows and the difficulty of investigating such an attack. Organizations may never fully discover to what extent a social engineer was able to infiltrate the organization. There are so many different vectors and possibilities of intrusion that it may be impossible to fully understand exactly what or who was compromised during the attack. An especially dangerous and frustrating situation is one where an attacker has an insider or accomplice within the organization. Some organizations have never identified these individuals within the organization. This uncertainty is difficult to recover from and defend against in the future. It is also another situation where a company may have substantial credibility problems and a loss of confidence from its shareholders after an attack if such details are disclosed. This exact situation has also led to a tug-of-war between regulators and organizations about disclosure laws. These laws vary widely state to state and between countries. Businesses fight to protect themselves from having to do extensive disclosure for the very reason of protecting the company's image and financial value.

19.5.2 Case Studies—Examples from Business. There are many well-known examples from real-life scenarios from organizations that can demonstrate the effectiveness of social engineering. Here are some well-known examples without any specifics that would embarrass the involved parties. Social-engineering attacks are real—they are not simply computer-security theory.⁵⁷

Case 1: One very well-known social-engineering attack in the business world is piggybacking. In this type of social engineering, an attacker will depend on an individual's sense of courtesy. Most people remember from their early school days that they are taught to hold a door open for someone who is behind him or her. This courtesy is often extended in the workplace, including secure areas such as a datacenter. A potential attacker may attempt to enter a datacenter without proper credentials by following closely behind someone else who is entering a datacenter with proper identification and authorization. The authorized individual will probably exhibit courtesy and hold the door open for the piggybacker, even if that person's identity is unknown. Holding something innocuous often creates confidence.

Result: In this example, the result is that an attacker who is not authorized has gained access to a secure facility by relying on another individual's sense of courtesy. It can be very difficult to demand that employees refrain from allowing piggybacking, but a policy must be in place to demand exactly that. All persons entering a secure facility should be required to fully use identification and authorization mechanisms every time.

Case 2: Another example that has been carried out in different variations involves an attacker using several social-engineering techniques to take advantage of

19 · 16 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

several coinciding events to exploit data, information, or equipment from an organization. The attacker would often make several telephone calls or send several emails to find a specific date where a company official, perhaps the CFO or director of technology, is out of the office. The attacker would then show up at the organization making claims that the company official authorized the attacker to take a certain computer from the company's premises. Usually the attacker will try to show a sense of urgency and extend a very confident display of authority. Many times the employee will cave into the attacker's demands without checking the story, and now a very important computer has been taken from the company. (A very commonly used twist on this case, and one that is becoming more and more common, is to install malware, remote backdoor software, or a keylogger on the computer instead of physically removing it.)

Result: In this case, the organization has lost control and ownership of a computer and/or its data. If the computer does not have any solid safeguards, such as data encryption mechanisms, the data and the computer could be used for any number of destructive activities. If information about the incident is made public, great damage to the organization's reputation can be done. The data harvested from the computer could also be sold to competitors or used as part of a blackmail scheme.

19.5.3 Success Rate. Although there are few accurate statistics of the success rate of social engineering, as is the case in many areas of information security, most experts believe the rate to be high. If history has anything to teach the security community through example, social engineering will continue to be a powerful and successful tool for criminals. Few organizations are immune to social engineering, no matter their industry or product. If even well-trained military personnel are vulnerable, every organization must take the success rate seriously.

The high success rate must also stress the importance of continued education of employees and reevaluation of the organization's efforts to combat social engineering and protect all data assets. This is an area where many organizations, of all sizes, continue to allocate too few resources. The frequency of social-engineering attempts also dictates the need for proper, efficient, and swift reporting of suspicious activity targeting individuals or the organization. It is very difficult to defend against social-engineering attempts if the organization does not know it is under attack or know where or how to report an incident. The probability of a successful attack can be substantially reduced with properly trained, supported, and motivated employees.

19.5.4 Small Business verses Large Organizations. The impact and dangers of social-engineering and low-tech attacks vary widely between small businesses and large corporations. As discussed above, the consequences can be potentially serious, all the way up to the collapse of the organization. Small businesses are often much less prepared operationally and financially and much less equipped to survive a serious breach of security. Conversely, and perhaps shocking to many, small businesses may have the upper hand over large organizations because of a substantially smaller workforce and collapsed management structure. It is much easier to communicate and engage everyone within a small company when an attack is attempted or carried out. Small businesses also have an advantage of a much smaller workforce to train; this results in much better prepared employees. Small business employees are probably

DETECTION 19 · 17

more likely to identify people who do not belong in the organization or should not be asking for access to sensitive areas or data. They may also be more likely to deny access or question someone whose story does not seem likely or is suspicious.

Large organizations can be mired in bureaucracy, ineffective management, or overly complicated reporting procedures. An attacker could carry out an entire plan before the security team in a large company would be alerted to a social-engineering attempt. Many times an individual is less willing to challenge the credentials of a stranger in a large organization. This is especially true where there are employees that feel they may be punished for questioning someone of a higher rank or preventing someone else from doing his or her job. The employee might be more likely to just let the attacker pass unquestioned rather than risk possible negative ramifications or scrutiny.

Bottom line, all organizations must be on the lookout for this type of attack. Criminals do not always choose easy or obvious targets. Any business, small or large, family-owned or corporate conglomerate, may be a target of an attack that utilizes social engineering.

19.5.5 Trends. While it is important for any information security manager to always keep a skeptical eye toward statistics, it is equally important to keep abreast of security threat trends. Social engineering is no exception. It may be difficult for any survey or poll to gather facts about how many attempts were made in any given year or how many were successful. Many social-engineering attempts are probably never detected, and even less are reported or admitted to on surveys. When it comes to such a powerful and successful attack mechanism, assume the worst: It is increasingly being used by criminals even when organizations know about this type of attack and even train employees to defend against it. Criminals create new forms and tactics every day, and people will probably continue to fall for these tactics.

19.6 DETECTION. Detection of social-engineering and low-tech attacks can be difficult. The nature of most types of social-engineering attacks is to take advantage of people's willingness to trust and help, which allows the perpetrator to circumvent technical controls. In many cases, the detection relies on people's ability to recognize a potential attack, including suspicious activity and respond appropriately. Further complicating detection is the potential that a social-engineering attack may not be a single occurrence, but many smaller events culminating in the attackers accessing restricted resources, hide nefarious activities, or stealing proprietary information. Using slow or multistep social-engineering attacks is even more difficult to detect and potentially more damaging. More recently, successful social-engineering attacks through the targeted email campaigns to senior level personnel have been widely reported. One penetration-testing expert estimates that social-engineering tactics account for less than 20 percent of the time spent in an attack; the rest of the time is spent in technical exploitation of the information originally gathered through deceit.⁵⁸

There are three main avenues of detection for social-engineering attacks: people, audit controls, and technology.

19.6.1 People. Since people are the vector for the attack, they are, in general, the first line of defense. Organizations need to provide employees with the resources and continuous education to help discern a potential attack from a legitimate request. All areas of an organization need the training. Social engineers are using all available resources; there is no department or person that is immune from a potential attack. In addition, organizations need to provide information on what to do during and

19 · 18 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

immediately after an attack. For example, during a phone attack, employees should be trained to remember as many details as possible. Items that the employee should try to remember include:

- Was the attacker male or female?
- Was a caller ID displayed?
- Was there noise in the background?
- Did he or she have an accent?
- What questions were asked?
- What answers were provided?

Employees should also be aware of people asking many questions, some of which may not make sense. In addition, callers or emails requesting names of managers or IT personnel should prompt a phone call to the IT security or investigation department. Organizational culture is also another pivotal piece to help defend against social-engineering attacks. The culture must stress and *reward* proper verification before providing information. Upper management can support a strong culture by praising employees who ask *them* to demonstrate their authenticity, rather than expressing irritation that anyone would have the nerve to do so.

All businesses, including security firms, must adhere to that mantra.⁵⁹ Unfortunately, too often VIPs or senior leaders do not have to follow the same verification process to change passwords or receive information. Finally, employees should be aware of their own policies and business practices. For example, all employees should understand why it is inappropriate that a caller would request a password change via phone or text message.⁶⁰

Organizations also need to provide clear information to their associates as to what to do when confronted with possible social engineering. For example, who should be notified during the actual event, immediately afterward, or both (depending on when the employee recognizes the attack). To help in the notification process, an organization can create an incident-notification information system and disseminate the information across the organization so an employee can immediately and easily determine whom to contact for help.

19.6.2 Audit Controls and Event Logging. Auditing and logging email, Internet content, systems logins, and systems changes of an organization can be used to detect social-engineering attacks. If the review of events does not happen in real time, there will be a delay from the time of the attack and identification of incident. In the instances where notification is not real time, forensic examiners can use the audit information to help piece together the attacks and determine root cause. Awareness teams can also use the log results to help devise additional training for the targeted groups. During real-time auditing, organizations can immediately enact their incident management plans to help limit the potential damage. Modern anti-malware and anti-phishing tools can provide not only defenses against such attacks but can also keep detailed records for analysis.

19.6.3 Technology for Detection. Organizations can implement technology to help limit social-engineering attacks. Content-filtering software can limit email and Website traffic; any identified malware should be blocked or stripped from email.

DEFENSE AND MITIGATION 19 · 19

Email monitoring tools can be used in bidirectional mode to inspect content in both directions. In addition, email-monitoring and content-filtering mechanisms can be used to scan for keywords or phrases that may trigger early warning signals. Any suspicious traffic, attachment, or email should be quarantined and reviewed prior to delivery. By scanning and blocking content, organizations may reduce the number of suspicious emails entering their networks and reduce the number of suspicious Websites that employees visit. Email services and anti-malware products have been improving their spam-blocking capabilities so that even phishing scams are being trapped successfully.

Advances in technology research will provide additional protection from social-engineering attacks, including the development of the Social Engineering Defense Architecture (SEDA). This architecture attempts to detect social-engineering attacks over the phone by identifying a legitimate employee versus a social engineer. The system is referred to as a text-independent voice signature authentication system. The system uses voice recognition technology, which would reduce the risk of a successful helpdesk attack. In addition to detecting an unauthorized caller, it can detect an insider masquerading as an employee with a higher security classification. The logging included in the architecture would aid a forensic examiner during an investigation.⁶¹

Risk-based authentication methods can also help limit potential social-engineering attacks; for example, if a customer or user always calls from a certain area code or phone number and a potential attacker calls from a different area code, systems can identify this and prompt for additional information prior to authenticating the individual.⁶²

Attackers are becoming more sophisticated. Recent phishing attacks are targeting systems, network, and security professionals. In these types of attacks, a phishing email is sent to an organization from a “customer” informing it of a phishing site attempting to steal customer information. Security personnel respond to the email and investigate the site. The site installs malware that allows the attacker to remotely control the machine. This change in attack methodology requires security personnel to be even more suspicious of any seemingly innocuous email telling the organization that there are phishing sites targeting that company.

19.7 RESPONSE. Responding to social-engineering and other low-tech attacks should fit into an organization’s incident-management process and response. As with all incident management plans, the responses should be well-defined, communicated, and tested. It behooves an organization to plan for the inevitable, especially as network and physical attacks are increasingly using multiple vectors in order to be successful.

See Chapter 53 in this *Handbook* for additional information on monitoring and controlling systems and see Chapter 56 for information on computer-security incident-response teams.

19.8 DEFENSE AND MITIGATION. The prevention of social-engineering attacks should be multifaceted, repeatable, and part of an organization’s defense-in-depth strategy. Since the very nature of the attack is to bypass or circumvent technical defenses by using people’s good nature and willingness to trust other human beings, the steps in preventing such attacks should focus on distinct areas, such as policy, training and awareness, technology, and physical defenses.

Focusing on all areas will help mitigate the threat of a successful social-engineering attack. Each area should have a regular process for review and auditing. Well-written policies provide the baseline of behavior that is acceptable and the potential consequences if they are not followed, but they must be accessible to employees. In the case of training, it should be integrated into the organization’s overall security awareness

19 · 20 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

program, included in all employee's evaluations, and tied to bonus pay and compensation. Organizations must train their employees in acceptable behavior and provide policies and the tools to identify and report potential social-engineering attacks.

Physical defenses can potentially block an intruder from entering a locked building, but people's sense of camaraderie does make the practices such as piggybacking or tailgating relatively easy. A recent study indicated that smokers returning to a building sometimes allow nonemployees into a secure location. Technological advances are also important but cannot be relied on as the only method of defense.

19.8.1 Training and Awareness. Organizations need to provide tools and knowledge to employees to help identify potential attacks and react to suspected attacks. Providing awareness is a continual process and should not be only for new hires. Training and awareness, when possible, should include real-life examples so employees can relate to the issue and understand the level of trust that is implied through their access to the facility and the data used for their positions. Employees should understand the responsibilities the company bestows on them. Essentially, employees need to be retrained that it is acceptable to ask why certain information is being requested or to see the badge of a person behind them.

A basic awareness program could include posters, email communications, and laminated instruction cards (hard cards) containing emergency contact numbers or other information. More mature awareness programs can include videos or brown-bag informational lunches. Whenever possible, teaching employees to defend against social-engineering attacks should be a live presentation with real-world examples. Ideally, the presentation should be tailored specially for each audience. For example, if the audience is helpdesk personnel, examples of potential social-engineering attacks should be described or demonstrated and discussed.

Employee training can also include instruction to keep file cabinets locked when not in use, lock workstations, and use cable locks, and contain instructions on how to create and remember good passwords.

The backbone of awareness training is a well-defined and executed information-security policy program. Please see Chapter 44 in this *Handbook* for additional information regarding information-security policy development, and Chapter 49 for additional information regarding awareness programs.

19.8.2 Technology for Prevention. Technology is emerging as a defense against certain types of social-engineering attacks. The technology enables organizations to identify some social-engineering attacks without relying on employees. This proactive identification enables an organization to mitigate the risk. Technology should comprise only one layer of defense and not be relied on as the sole defense.

Technologies such as content-monitoring systems for both email and Web content can help identify phishing attacks. In addition, organizations can install timely security patches and use up-to-date antivirus and antispyware software to help mitigate the risk of viruses, Trojans, and worms. Most versions of browsers and browser plug-ins are allowing users to evaluate the trustworthiness of Websites.⁶³

Ideally, employees should be prevented from downloading and installing unapproved software. However, organizations can employ inventory systems or other methodologies to detect illegal programs on a network. Certain types of systems can prevent malware-infected machines from entering the network.

Desktop and laptop configuration changes can be made to reduce the risk of a successful attack; changes include disabling pop-up windows in browsers, disallowing

CONCLUDING REMARKS 19 · 21

the automatic installation of Active-X controls, multiple versions of Java, limiting the types of cookies that Websites can place on local machines, using automatic password-protected screen savers, and finally, using email certificates for authentication.

In the same regard, organizations should review technology processes and verify that they are not inadvertently supplying information to potential social engineers. For example, metadata in documents should be removed before being accessible to outsiders, and regular Web searches should be conducted to ensure former or current employees are not posting information on the Internet.

19.8.3 Physical Security. Physical security mechanisms can reduce the risk of a successful social-engineering attack. All employees should have identification cards that they are required to display at all times. Secured areas within an organization should be locked, have limited access, and be monitored for noncompliance. Door alarms that can detect tailgating or piggybacking can be installed in areas. Cameras or other closed circuit monitoring technology can thwart potential intruders. Security personnel should watch all facility access points. All office doors, desks, file cabinets, and other storage devices should have keys and remain locked when not being accessed. Dumpsters® or recycling bins should also have locks that would prevent the removal of documents meant for shredding or incineration.

- Desktops, laptops, and other computer hardware should be physically locked in place. Users should be required to have strong passwords, and in the case of laptops or desktops, have automatic screen savers that require a password to unlock. All magnetic media need to have secure storage.
- Most organizations have a certain percentage of mobile workforces. Special training should be provided to them to prevent the loss of equipment or information. Training should include information such as:
 - Laptops should remain with the traveler and not checked in luggage.
 - Laptops should be locked in a safe or to a secure surface at all times.
 - Peripheral devices such as USB drives and handheld devices should have strong passwords.
 - Conversations involving confidential information should be prohibited in public.
 - Travelers should be aware of their surroundings, and special considerations regarding electronic communication and equipment usage should be considered while not in the United States and Europe.

Please see Chapters 22 and 23 of this *Handbook* for more information regarding physical security.

19.9 CONCLUDING REMARKS. Social-engineering attacks are unlike technological computer attacks. The vector of the attack is human, the nature of the attack is to circumvent controls, and the success of the attack depends on people's willingness to trust others. Simply being polite and holding a door open for a stranger can be bad for an organization, if the stranger is a social engineer.

Social-engineering attacks are not new in society; they have been used for millennia. Social engineers use many different methods to execute social-engineering and low-tech attacks. These methods could involve human contact, no human contact, or a combination of technology and social-engineering tactics.

19 · 22 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

Psychologists and social psychologists have offered a number of reasons for the success of social engineering. They theorize that human nature allows attackers to trick or con reasonable people into divulging information. A social engineer's profile does not fit into a single model, and their attacks are difficult to detect. Social-engineering attacks have grown more prolific, effective, and dangerous to both organizations and individuals.

Even though social-engineering attacks are difficult to defend against, there are technical, process, and personal defenses that an organization can adopt to prevent them and minimize their effects. As with all information security issues, a defense-in-depth strategy can help mitigate the risks associated with social engineering.

Between July and August 2011, Dimensional Research surveyed 835 IT professionals in the Western economies. Some of their findings reinforce the details in this chapter⁶⁴:

- The threat of social engineering is real.
- Financial gains are the primary motivation of social engineering.
- Social-engineering attacks are costly, especially in large organizations.
- [There is a] Lack of proactive training to prevent social-engineering attacks.

Although the focus of information assurance practice seems to be on technical attacks, social-engineering and low-tech attacks will continue to remain relevant and dangerous threats, worthy of attention and mitigation.

19.10 NOTES

1. M. E. Kabay, "Social Engineering Simulations," Network World Security Strategies, December 18, 2000, www.mekabay.com/nwss/059_social_engineering_v03.pdf
2. Ira Winkler, *Spies Among Us* (Indianapolis, IN: Wiley Publishing, Inc., 2005): 111.
3. L. Giles, trans., *Sun Tzu on the Art of War: The Oldest Military Treatise in the World*, www.chinapage.com/sunzi-e.html
4. Charles R. King, *Learning Secrets: Intelligence in the Second Punic War*, Kindle ed., Amazon Digital Services, Inc., 2012, www.amazon.com/Learning-Secrets-Intelligence-Second-ebook/dp/B0082PZIYI
5. Ben Macintyre, *Double Cross: The True Story of the D-Day Spies* (New York: Crown, 2012).
6. Edward Lucas, *Deception: The Untold Story of East-West Espionage Today* (New York: Walker & Company, 2012).
7. Steven Spielberg (dir.), *Catch Me If You Can*, 2002, www.imdb.com/title/tt0264464
8. Frank W. Abagnale, Website, www.abagnale.com/index2.asp
9. Stacy Cowley, "How a lying 'social engineer' hacked Wal-Mart," *CNN-Money*, August 8, 2012, <http://money.cnn.com/2012/08/07/technology/walmart-hack-defcon/index.htm>
10. For more information on this documentary, see www.hbo.com/docs/programs/hackingdemocracy/index.html
11. Andy O'Donnell, "Caller ID Spoofing," *About.com: Internet/Network Security | Cybercrime*, May 2, 2012, <http://netsecurity.about.com/od/securityadvisories/a/Caller-Id-Spoofing.htm>

NOTES 19 · 23

12. Joan Goodchild, “4 Reasons Why Executives are the Easiest Social Engineering Targets,” CSOonline.com, July 14, 2012, www.csoonline.com/article/599456/4-reasons-why-executives-are-the-easiest-social-engineering-targets
13. “Social Engineering,” AuditMyPC.com, www.auditmypc.com/social-engineering.asp
14. Dan Arnell, “Payroll Giant Gives Scammer Personal Data of Hundreds of Thousands of Investors,” ABC News, June 26, 2006, <http://abcnews.go.com/Technology/story?id=2160425&page=1>
15. M. E. Kabay, personal communication, 2013.
16. “Seduce,” *The Free Dictionary*, www.thefreedictionary.com/seduce
17. Federal Bureau of Investigation, “Citadel Malware Continues to Deliver Reveton Ransomware in Attempts to Extort Money,” FBI New E-Scams & Warnings Website, November 30, 2013, www.fbi.gov/scams-safety/e-scams
18. M. E. Kabay, “Office Supply Scams,” M. E. Kabay Website, 2011, www.mekabay.com/nwss/891_office_supply_scams.pdf
19. Bureau of Consumer Protection, “Avoiding Office Supply Scams,” Federal Trade Commission, March 2000, <http://business.ftc.gov/documents/bus24-avoiding-office-supply-scams>
20. Ajay Gupta, “The Art of Social Engineering,” Addison-Wesley, August 23, 2002, www.informit.com/articles/article.aspx?p=28802
21. Michael Winter, “Confetti from police files tossed at Macy’s parade,” USA TODAY, November 26, 2012, www.usatoday.com/story/news/nation/2012/11/26/thanksgiving-macys-parade-confidential-police-confetti/1728375/
22. M. Horowitz, “Defending against Malicious CDs and USB Flash Drives,” Computerworld | Defensive Computing Blog, July 3, 2011, http://blogs.computerworld.com/18560/defending_against_malicious_cds_and_usb_flash_drives
23. Microsoft Support. “How To Minimize Metadata in Office Documents,” Microsoft Corporation, January 24, 2007, <http://support.microsoft.com/default.aspx?scid=kb;EN-US;Q223396>
24. The 419 Coalition, “The Nigerian Scam (419 Advance Fee Fraud) Defined,” Nigeria—The 419 Coalition Website, accessed March 30, 2007, <http://home.rica.net/alphae/419coal>
25. *Webopedia*, s.v. “Phishing,” accessed March 30, 2007, www.webopedia.com/TERM/P/phishing.html
26. Robert Richardson, *2010/2011 Computer Crime and Security Survey*, Information Week, Computer Security Institute, June 6, 2011, <http://reports.informationweek.com/abstract/21/7377/security/research-2010-2011-csi-survey.html>
27. RSA FraudAction Research Lab, “Anatomy of an Attack,” RSA Speaking of Security Blog, April 1, 2011, <http://blogs.rsa.com/rivner/anatomy-of-an-attack>
28. *SearchExchange*, s.v. “Spim (Instant Messaging Spam),” accessed March 30, 2007, http://searchexchange.techtarget.com/sDefinition/0,290660,sid43_gci952820,00.html
29. Will Sturgeon, “U.S. Makes First Arrest for Spim,” CNet.com, February 21, 2005, http://news.cnet.com/U.S.-makes-first-arrest-for-spim/2100-7355_3-5584574.html
30. Linda Rosencrance, “Porn Site WhiteHouse.com Domain Name Up for Sale,” Computerworld, February 10, 2004, www.computerworld.com/s/article/90035/Porn_site_WhiteHouse.com_domain_name_up_for_sale

19 · 24 SOCIAL-ENGINEERING AND LOW-TECH ATTACKS

31. Ryan Naraine, "Hushmail DNS Attack Blamed on Network Solutions," eWeek.com, April 29, 2005, <http://www.eweek.com/c/a/Security/Hushmail-DNS-Attack-Blamed-on-Network-Solutions/>
32. Steve Stasiukonis, "Social Engineering, the USB Way," Dark Reading, June 7, 2006, <http://www.darkreading.com/perimeter/social-engineering-the-usb-way/208803634>
33. Richardson, *2010/2011 Computer Crime and Security Survey*
34. Jonathan Ness, "More Information about the DLL Preloading Remote Attack Vector," Microsoft Security Research & Defense, August 23, 2010, <http://blogs.technet.com/b/srd/archive/2010/08/23/more-information-about-dll-preloading-remote-attack-vector.aspx>
35. Brian Koerner, "Vishing," About.com: Vishing, accessed March 22, 2007, en.wikipedia.org/wiki/voice_phishing
36. M. E. Kabay, "Kitty Porn," *InfoSec Perception* Website, September 28, 2012, <http://resources.infosecskills.com/perception/kitty-porn>
37. Patterson Clark and Robert O'Harrow, Jr., "Social Engineering: Using Social Media to Launch a Cyberattack." *The Washington Post*, September 26, 2012, www.washingtonpost.com/investigations/social-engineering-using-social-media-to-launch-a-cyberattack/2012/09/26/a282c6be-0837-11e2-a10c-fa5a255a9258_graphic.html
38. Julie Inouye, "Oh, What a Year It's Been!" LinkedIn Blog, December 21, 2012, <http://blog.linkedin.com/2012/12/21/oh-what-a-year-its-been/>
39. LinkedIn, "User Agreement," last revised May 13, 2013, www.linkedin.com/static?key=user_agreement
40. Stacy Cowley, "LinkedIn Is a Hacker's Dream Tool," *CNNMoney*, March 12, 2012, <http://money.cnn.com/2012/03/12/technology/linkedin-hackers/index.htm>
41. Sarah Granger, "Social Engineering Fundamentals, Part I: Hacker Tactics," Symantec Website, December 18, 2001, updated November 3, 2010, www.symantec.com/connect/articles/social-engineering-fundamentals-part-i-hacker-tactics
42. Microsoft Corporation, "How to Protect Insiders from Social Engineering Threats," *Microsoft Technet*, August 18, 2006, <http://technet.microsoft.com/library/Cc875841#EIXAE>
43. Bruce Schneier, "The Psychology of Security," Essays and Op Eds, 2007, www.schneier.com/essay-155.html
44. Richards J. Heuer, Jr., "What Are Cognitive Biases?" in *Psychology of Intelligence Analysis*, Part 3—Cognitive Biases, 1999, <https://www.cia.gov/library/center-for-the-study-of-intelligence/csi-publications/books-and-monographs/psychology-of-intelligence-analysis/art12.html>
45. R. Cialdini, *Influence: The Psychology of Persuasion* (New York: HarperCollins, 1998).
46. Mara Mathers, Eldar Shafir, and Marcia K. Johnson, "Misremembrance of Options Past: Source Monitoring and Choice," *Psychological Science* 11, No. 2, March 2000, <http://psych.princeton.edu/psychology/research/shafir/pubs/Misremembrance%202000.pdf>
47. Granger, "Hacker Tactics."
48. J. St. B. T. Evans, J. L. Barston, and P. Pollard, "On the Conflict between Logic and Belief in Syllogistic Reasoning," *Memory and Cognition*, 11 (1983): 295–306.

NOTES 19 · 25

49. Steven Soderbergh (dir.), *Ocean's Eleven*, 2001, www.imdb.com/title/tt0240772
50. R. B. Zajonc, "Attitudinal Effects of Mere Exposure," *Journal of Personality and Social Psychology* 9, No. 2 (1968): 1–27.
51. E. Grandjean, "A Prime Target for Social Engineering Malware," *McAfee Security Journal*, Fall (2008): 9–12.
52. A. Tversky and D. Kahneman, *Judgment under Uncertainty: Heuristics and Biases*, *Science*, 185 (1974): 1124–1130.
53. Joan Goodchild, "Social Engineering Stories," CSO, May 24, 2010, www.csoonline.com/article/594924/social-engineering-stories
54. Tim Wilson, "Eight Faces of a Hacker," *Dark Reading*, March 29, 2007, www.darkreading.com/security/news/208804443/eight-faces-of-a-hacker.html
55. Malcolm Allen, "Social Engineering, A Means to Violate a Computer System," SANS Institute Information Security Reading Room, June 2006, www.sans.org/reading_room/whitepapers/engineering/social-engineering-means-violate-computer-system_529
56. Radha Gulati, "The Threat of Social Engineering and Your Defense Against It," SANS Institute Information Security Reading Room, 2003, www.sans.org/reading_room/whitepapers/engineering/threat-social-engineering-defense_1232
57. Granger, "Hacker Tactics."
58. Kelly Jackson Higgins, "Social Engineering Gets Smarter," *Dark Reading*, June 16, 2006, www.darkreading.com/authentication/social-engineering-gets-smarter/208803849
59. P. Roberts, "Hacker Grrl Behind HBGary Attack," *threatpost*, March, 17, 2011, <http://threatpost.com/hacker-grrl-behind-hbgary-attack-031711/75037>
60. Ryan Groom, "Top 5 Social Engineering Techniques," About: Business Security, accessed March 17, 2007, <http://bizsecurity.about.com/od/physicalsecurity/a/topsocialengine.htm>
61. Michael Hoeschele and Marcus Rogers, "Detecting Social Engineering," in *Advances in Digital Forensics: IFIP International Conference on Digital Forensics* (February 2005): 67–71.
62. RSA, "Risk-Based Authentication," Information Security Glossary, RSA Website, accessed April 12, 2012, www.rsa.com/glossary/default.asp?id=1094
63. M. J. Edwards, "IE 7.0 and Firefox 2.0 Both Have New Antiphishing Technologies," WindowsITPro, October 26, 2006, <http://windowsitpro.com/networking/ie-70-and-firefox-20-both-have-new-antiphishing-technologies>
64. Dimensional Research, "The Risk of Social Engineering on Information Security: A Survey of IT Professionals," Check Point Software Technologies Website, September 2011, www.checkpoint.com/press/downloads/social-engineering-survey.pdf

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 20

SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

Stephen Cobb

20.1 UNWANTED EMAIL AND OTHER PESTS: A SECURITY ISSUE	20·1	20.1.1 Common Elements	20·2	20.4.5 Spam Filters	20·21
		20.1.2 Chapter Organization	20·3	20.4.6 Network Devices	20·24
				20.4.7 Email Authentication	20·25
				20.4.8 Industry Initiatives	20·26
				20.4.9 Legal Remedies	20·26
20.2 EMAIL: AN ANATOMY LESSON	20·3	20.2.1 Simple Mail Transport Protocol	20·3	20.5 PHISHING	20·27
		20.2.2 Heads-Up	20·4	20.5.1 What Phish Look Like	20·27
				20.5.2 Growth and Extent of Phishing	20·29
				20.5.3 Where Is the Threat?	20·30
				20.5.4 Fighting Phishing	20·31
20.3 3 SPAM DEFINED	20·6	20.3.1 Origins and Meaning of Spam (not SPAM®)	20·7	20.6 TROJAN CODE	20·31
		20.3.2 Digging into Spam	20·8	20.6.1 Classic and Recent Trojans	20·32
		20.3.3 Spam's Two-Sided Threat	20·12	20.6.2 Basic Anti-Trojan Tactics	20·34
				20.6.3 Lockdown and Quarantine	20·35
20.4 FIGHTING SPAM	20·18	20.4.1 Enter the Spam Fighters	20·18	20.7 CONCLUDING REMARKS	20·36
		20.4.2 A Good Reputation?	20·18	20.8 FURTHER READING	20·36
		20.4.3 Relaying Trouble	20·20	20.9 NOTES	20·36
		20.4.4 Black Holes and Block Lists	20·20		

20.1 UNWANTED EMAIL AND OTHER PESTS: A SECURITY ISSUE. Three oddly named threats to computer security are addressed in this chapter: *spam*, *phishing*, and *Trojan code*. Spam is unsolicited commercial email. Phishing is the use of deceptive unsolicited email to obtain—to fish electronically for—confidential information. Trojan code, a term derived from the Trojan horse, is software designed to achieve unauthorized access to systems by posing as legitimate applications. In this chapter, we

20 · 2 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

outline the threats posed by spam, phishing, and Trojans, as well as the mitigation of those threats.

These threats might have strange names, but they are no strangers to those whose actions undermine the benefits of information technology. Every year, for at least the last three years, the U.S. Internal Revenue Service (IRS) has had to warn the public about email scams that take the name of the IRS in vain, attempting to defraud taxpayers of their hard-earned money by aping the agency's look and feel in email messages. The fact that the IRS states that the agency never ever uses email, text messaging, or social-media channels to communicate with taxpayers¹ is a sad comment on our society, for there is no good reason why this should be the case. After all, the agency allows electronic filing of annual tax returns, and every day hundreds of millions of people around the world do their banking and bill paying online, in relative safety. The technology exists to make email safe and secure. In many ways, this chapter is a catalog of what has happened because we have not deployed the technology effectively.

20.1.1 Common Elements. Each of these threats is quite different from the other in some respects, but all three have some important elements in common; first, and most notably, they use deception. These threats prey on the gullibility of computer users and achieve their ends more readily when users are ill-trained and ill-informed (albeit aided and abetted, in some cases, by poor system design and poor management of services, such as broadband connectivity).

Second, all three attacks are enabled by system services that are widely used for legitimate purposes. Although the same might be said of computer viruses—they are code and computers are built to run code—the three threats that are the focus of this chapter typically operate at a higher level, the application layer of the Open Systems Interconnection (OSI) model (discussed in Section 6.5.1 in this *Handbook*). Indeed, this fact may contribute to the extent of their deployment—these threats can be carried out with relatively little technical ability, relying more on the skills associated with social engineering than with coding. For example, anyone with an Internet connection and an email program can send spam. That spam can spread a ready-made Trojan. Using a toolkit full of scripts, you can add a Website with an input form to your portfolio and go phishing for personal data to collect and abuse.

Another reason for considering these three phenomena together is the fact that they often are combined in real-world exploits. The same mass emailing techniques used to send spam may be employed to send out messages that spread Trojan code. As described in Chapter 15 in this *Handbook*, Trojan code may be used to aid phishing operations. Systems compromised by Trojan code may be used for spamming, and so on. What we see in these attacks today are the very harmful and costly result of combining relatively simple strategies and techniques with standards of behavior that range from the foolish and irresponsible to the unabashedly criminal.

These three threats also share the distinction of having been underestimated when they first emerged. All three have evolved and expanded in the 21st century. For example, for 2012, Symantec reported that Android-operating-system threats (virtually unknown a decade earlier) against mobile devices increased drastically between January 2010 and the end of 2012, from fewer than 100 unique variants in about 10 families to around 4,500 variants in about 170 families.² In addition, an increasing number of phishing scams in 2012 were being carried out through social media—a term not even mentioned in the corresponding report for the first half of 2006.³ (For more information about viruses and worms, spyware, rootkits, and other malware, see Chapter 16 in this *Handbook*.)

EMAIL: AN ANATOMY LESSON 20 · 3

One other factor unites these three threats: their association with the emergence of financial gain as a primary motivator for writing malicious code and abusing Internet connectivity (as discussed in Chapter 2 in this *Handbook*). The hope of making money is the primary driver of spam. Phishing is done to facilitate fraud for gain through theft of personal data and credentials, either for use by the thief or through resale in the underground economy. Trojan code is used to advance the goals of both spammers and perpetrators of phishing scams. In short, all three constitute a very real threat to computer security, the well-being of any computer-using organization, and the online economy.

20.1.2 Chapter Organization. After a brief email anatomy lesson, each of the three threats addressed by this chapter is examined in turn. The email anatomy lesson is provided because email plays such a central role in these threats, enabling spam and phishing and the spread of Trojan code. Responses to the three threats are discussed with respect to each other, along with consideration of broader responses that may be used to mitigate all three.

For an introduction to the general principles of social engineering, see Chapter 19 in this *Handbook*.

20.2 EMAIL: AN ANATOMY LESSON. Email plays a role in numerous threats to information and information systems. Not only does it enable spam and phishing, it is used to spread Trojan code, viruses, and worms. A basic understanding of how email works will help to understand these threats and the various countermeasures that have been developed.

20.2.1 Simple Mail Transport Protocol. All email transmitted across the Internet is sent using an agreed-on industry standard: the Simple Mail Transport Protocol (SMTP). Any server that speaks SMTP is able to act as a Mail Transfer Agent (MTA) and send mail to, and receive mail from, any other server that speaks SMTP. To understand how simple SMTP really is, Exhibit 20.1 presents an example of an SMTP transaction.⁴ The text that follows represents the *actual* data being sent and received by email servers (as viewed in a telnet session, with the words in CAPS, such as HELO and DATA, being the SMTP commands defined in the relevant standards, and the numbers being the standard telnet responses).

Developed at a time when computing resources were relatively expensive, and designed to operate even when a server was processing a dozen or more message connections per second, the SMTP conversation was kept very simple in order to be very brief. However, that simplicity is both a blessing and a curse. As you can see from the example, only two pieces of identity information are received before the mail is delivered: the identity of the sending server, in this case example.com, and the From address, in this case foo@example.com. SMTP has no process for verifying the validity of those identity assertions, so both of those identifiers can be trivially falsified. The remaining contents of the email, including the subject and other header information, are transmitted in the data block and are not considered a meaningful part of the SMTP conversation. In other words, no SMTP mechanism exists to verify assertions such as “this message is from your bank and concerns your account” or “this message contains the tracking number for your online order” or “here is the investment newsletter that you requested.”

As described in more detail later, some email services do perform whitelist or blacklist look-ups on the Internet Protocol (IP) address of the sending server during

20 · 4 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

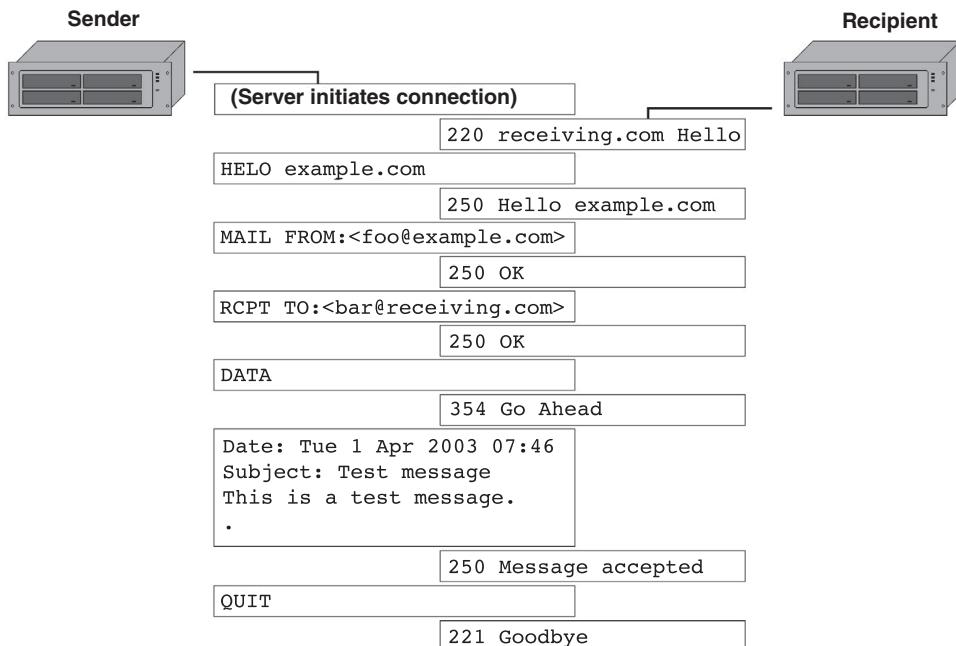


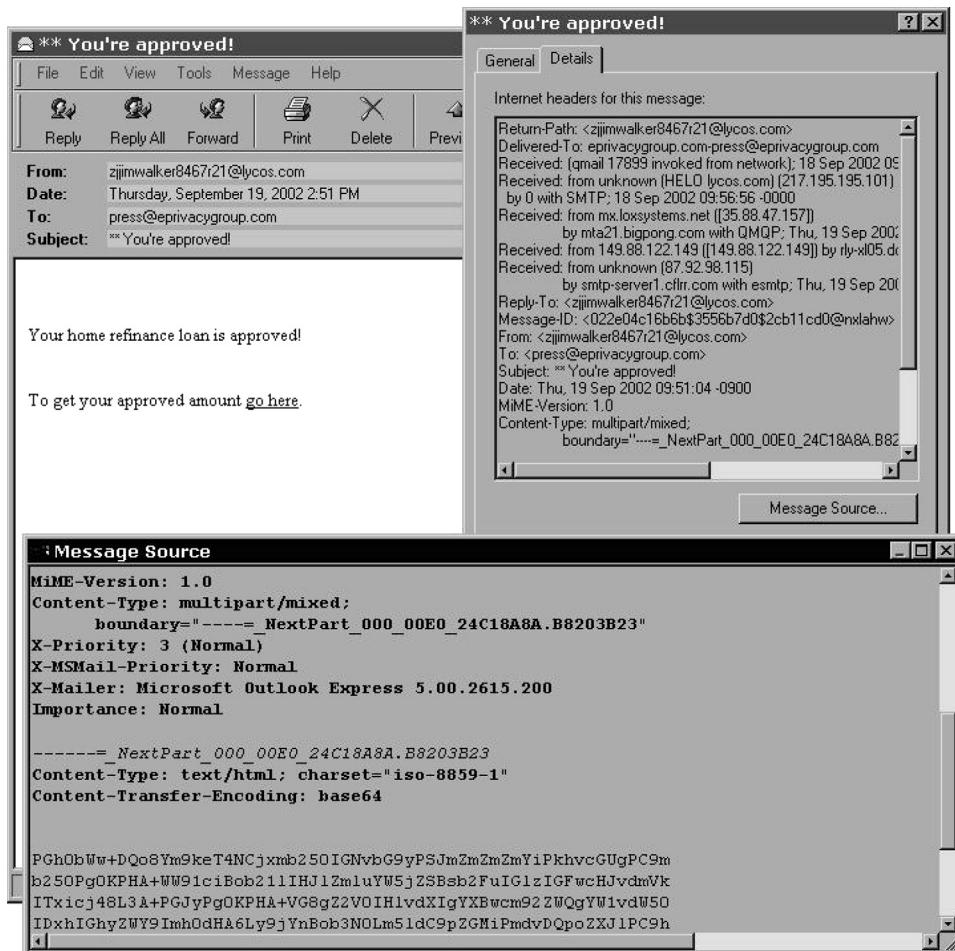
EXHIBIT 20.1 Basic Email Protocol

the SMTP conversation, but those inquiries can dramatically slow mail processing, requiring extra capacity to offset the loss of efficiency. A whitelist identifies email senders that are trusted; a blacklist identifies those that are not trusted. Maintenance of whitelists can be time consuming, and blacklists have a long history of inaccuracies and legal disputes. In short, the need for speed creates a system in which there are virtually no technical consequences for misrepresentations in mail delivery. And this is precisely why spammers have been, and continue to be, incredibly effective at getting unwanted email delivered.

SMTP is, as Winston Churchill might have put it, the worst way of doing email, except for all the others that have been tried. The reality is that SMTP works reliably and has been widely implemented. To supplant SMTP with anything better would mean a wholesale redesign of the entire global email infrastructure, a task that few in the industry have been willing to undertake. Some people have endeavored to develop solutions that can ride atop the existing SMTP infrastructure, allowing SMTP to continue functioning efficiently while giving those who use it the option of engaging more robust features that help differentiate legitimate mail from spam. There is more on these solutions later in the chapter.

20.2.2 Heads-Up. Email cannot be delivered without something called a header, and every email message has one, a section of the message that is not always displayed in the recipient's email program but is there nonetheless, describing where the message came from, how it was addressed, and how it was delivered. Examining the header can tell you a lot about a message. Consider how a message appears in Microsoft Outlook Express, as illustrated in Exhibit 20.2.

At the top you can see the From, the To, and the Subject. For example, you can see part of a message that appears to be from zjjimwalker8467r21@lycos.com

EMAIL: AN ANATOMY LESSON 20 · 5**EXHIBIT 20.2** Viewing Message Header Details in Outlook Express

to press@eprivacygroup.com with the subject: *You're approved! As you may have guessed, press@eprivacygroup.com is not a real person. This is just an address that appears on a company Website as a contact for the press, and of course, nobody actually used this address on a mortgage application. The address was harvested by a program that automatically scours the Web for email addresses.

When you open this message in email clients (e.g., Outlook or Outlook Express) and use the File/Properties command or equivalent, you can click on the Details tab to see how the message made its way through the Internet. The first thing you see is a box labeled "headers," as shown in Exhibit 20.2. Reading this will tell you that the message was routed through several different email servers. Qmail is the company's mail program, which is the first instance of Received. The next three below that are intermediaries, until you get to the last one, smtp-server1.cflrr.com. That is an email server in Central Florida (cfl) on the Road Runner (rr) cable modem network, which supplies high-speed Internet access to tens of thousands of households.

So who sent this message? That is very hard to say. As you might expect, there is no such address as zjjimwalker8467r21@lycos.com. The best way to determine who sent a spam message like this is to examine the content. Spam cannot reel in suckers

20 · 6 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

unless there is some way for the suckers to contact the spammer. Sometimes this is a phone number, but in this message it is a hyperlink to a Web page. However, clicking links in spam messages is a dangerous way to surf—a much safer technique to learn more about links in spam is message source inspection. Outlook Express provides a Message Source view, but this may only show the encoded content and not be readable ASCII. Some email clients (e.g., Outlook) allow one to set a default of conversion of all messages into plain text; if one wishes to see images, one can have them downloaded on a per-message basis. In addition, many email clients can display the underlying details of any hyperlink when one hovers the cursor over the label for the link.

One way to get at the content of such messages is to forward them to a different email client, for example, Qualcomm's Eudora, then open the mailbox file with a text editor such as TextPad. This reveals the http reference for the link that this spammer wants the recipient to click. In this case, the sucker who clicks on the link that says “To get your approved amount go here” is presented with a form that is not about mortgages, but about data gathering. To what use the data thus gathered will be put is impossible to say, but a little additional checking using the *ping* and *whois* commands reveals that the Web server collecting the data is in Beijing. The chances of finding out who set it up are slim. The standard operating procedure is to set up and take down a Website like this very quickly, gathering as much data as you can before someone starts to investigate. However, over the last 10 years, the resources devoted to investigating spammers have been minimal compared to the resources that have gone into sending spam. Even when lawmakers can agree on what spam is and how to declare it illegal, government funds are rarely allocated for spam fighting. There have been a few high-profile arrests and prosecutions, often driven by large companies, such as Microsoft and AOL, but spamming remains a relatively low risk form of computer abuse.

If one can identify a link in spam, the Better-Whois service (<http://betterwhois.com/>) may be able to provide details of who is officially registered as owners and administrators of the site—assuming they have not hidden themselves behind anonymizing companies such as GoDaddy.com, which can conceal the true owners while providing an anonymized email address through which to attempt to communicate with the owners.

20.3 3 SPAM DEFINED. The story of spam is a rags-to-riches saga of a mere nuisance that became a multibillion-dollar burden on the world’s computing resources. For all the puns and jokes one can make about spam, it may be the largest single thief of computer and telecommunication resources since computers and telecommunications were invented.

Perhaps it is no surprise that spam has become a costly problem, because spam was the first large-scale computer threat to be purely profit-driven. The goal of most spam is to make money, and it cannot make money if people do not receive it. In fact, spam software will not send spam to a mail server that has a slow response time; it simply moves on to other targets that can receive email at the high message-per-minute rate needed for spam to generate income (based on 1 person in perhaps 100,000 actually responding to the spam message). If spam does not generate income, it becomes pointless because it takes money to send spam. You need access to servers and bandwidth (which you either have to pay for or steal). Ironically, a few simple changes to current email standards could put an end to most spam by more reliably identifying and authenticating email senders, a subject addressed later in this chapter.

3 SPAM DEFINED 20 · 7

20.3.1 Origins and Meaning of Spam (not SPAM®). SPAM® has been a trademark of Hormel Foods for over 70 years.⁵ For reasons that will be discussed in a moment, the world has settled on *spam* as a term to describe unwanted commercial email. However, uppercase SPAM® is still a trademark, and associating SPAM®, or pictures of SPAM®, with something other than the Hormel product could be a serious violation of trademark law. Security professionals should take note of this. The word *Spam* is acceptable at the start of a sentence about unsolicited commercial email, but SPAM can be used for spam only if the rest of the text around it is uppercase, as in “TIRED OF SPAM CLUTTERING YOUR INBOX?”

The use of the word “spam” in the context of electronic messages is widely thought to stem from a comedy sketch in the twenty-fifth episode of the BBC television series *Monty Python’s Flying Circus*.⁶ First aired in 1968, before email was invented, the sketch featured a restaurant in which SPAM® dominated the menu. When a character called Mr. Bun asks the waitress, “Have you got anything without SPAM® in it?” she replies: “Well, there’s SPAM®, egg, sausage and SPAM®, that’s not got much SPAM® in it.” The banter continues in this vein until Mr. Bun’s exasperated wife screams, “I don’t like SPAM®!” She is further exasperated by an incongruous group of Viking diners singing a song, the lyrics of which consist almost entirely of the word “SPAM®.” A similar feeling of exasperation at having someone foist on you something that you do not want and did not ask for clearly helped to make “spam” a fitting term for the sort of unsolicited commercial email that can clutter inboxes and clog email servers.

In fact, the first use of spam as a term for abuse of networks did not involve email. The gory details of spam were carefully researched in 2003 by Brad Templeton, former chairman of the board of the Electronic Frontier Foundation.⁷ According to Templeton, the origins lie in annoying and repetitive behavior observed in multiuser dungeons (also known by the acronym MUDs, an early term for a real-time, multiperson, shared environment). From there the term *spam* migrated to bulletin boards—where it was used to describe automated, repetitive message content—thence to USENET, where it was applied to messages that were posted to multiple newsgroups.

A history of USENET and its relationship to the history of spam is salutary for several reasons. First of all, spam pretty much killed USENET, which was once a great way to meet and communicate with other Internet users who shared common interests, whether those happened to be political humor or HTML coding. Newsgroups got clogged with spam to the point where users sought alternative channels of communication. In other words, a valuable and useful means of communication and an early form of social networking was forever tainted by the bad behavior of a few individuals who were prepared to flaunt the prevailing standards of conduct.

The second spam–USENET connection is that spam migrated from USENET to email through the *harvesting* of email addresses, a technique that spammers—the people who distribute spam—then applied to Web pages and other potential sources of target addresses. Spammers found that software could easily automate the process of reading thousands of newsgroup postings and extracting, or harvesting, any email addresses that appeared in them. It is important to remember that, although it might seem naive today, the practice of including one’s email address in a message posted to a newsgroup was common until the mid-1990s. After all, if you posted a message looking for an answer, then including your email address made it easy for people to reply. It may be hard for some computer users to imagine a time when email addresses were so freely shared, but that time is worth remembering because it shows how easily the abuse of a system by a few can erode its value to the many. Harvesting, which resulted in receiving spam at an email address provided in a newsgroup posting for

20 · 8 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

the purposes of legitimate communication, was an enormous advance in abuse of the Internet and a notable harbinger of things to come. Elaborate countermeasures had to be developed, all of which impeded the free flow of communication. Much of the role fulfilled by newsgroups migrated to closed forums, where increasingly elaborate measures are used to prevent abuse.

20.3.2 Digging into Spam. Around 1996, as early adopters of email began experiencing increasing amounts of unsolicited email, efforts were made to define exactly what kind of message constituted spam. This was no academic exercise, and the stakes were surprisingly high. We have already employed one definition: unsolicited commercial email (sometimes referred to as UCE). This definition seems simple enough, capturing as it does two essential points: Spam is email that people did not ask to receive, and spam is somehow commercial in nature.

However, while UCE eventually became the most widely used definition of spam, it did not satisfy everyone. Critics point out that it does not address unsolicited messages of a political or noncommercial nature (such as a politician seeking your vote or a charity seeking donations). Furthermore, the term “unsolicited” is open to a lot of interpretation, a fact exploited by early mass emailings by mainstream companies whose marketing departments used the slimmest of pretexts to justify sending people email. As the universe of Internet users expanded in the late 1990s, encompassing more and more office workers and consumers, three things about spam became increasingly obvious:

1. People did not like spam.
2. Spammers could make money.
3. Legitimate companies were tempted to send unsolicited email.

20.3.2.1 Get-Rich-Quick Factor. How much money could a spammer make? Consider the business plan of one Arizona company, C.P. Direct, that proved very profitable until it was shut down in 2002 by the U.S. Customs Service and the Arizona Department of Public Safety. Here are some of the assets that authorities seized:

- Nearly \$3 million in cash plus a large amount of expensive jewelry
- More than \$20 million in bank accounts
- Twelve luxury imported automobiles (including eight Mercedes, plus assorted models from Lamborghini, Rolls-Royce, Ferrari, and Bentley)
- One office building and assorted luxury real estate in Paradise Valley and Scottsdale

These profits were derived from the sale of \$74 million worth of pills that promised to increase the dimensions of various parts of the male and female anatomy. The company had used spam to sell these products, but in fact it was not shuttered for sending spam, the legal status of which remains ambiguous to this day, having been defined differently by different jurisdictions (not least of which was the U.S. Congress, which arguably made a hash of it in 2004).

C.P. Direct crossed several lines, not least of which was the making of false promises about its products (none of which it ever tested and all of which turned out to contain the same ingredients, regardless of which part of the human anatomy they were promised

3 SPAM DEFINED 20 · 9

to enhance). The company compounded its problems by refusing to issue refunds when the products did not work as claimed. However, rather than discourage spammers, this case proved that spam can make you rich, quick. Indeed, the hope of getting rich quick remains the main driver of spam. The fact that a handful of people were prosecuted for conducting a dubious enterprise by means of spam was not perceived as a serious deterrent.

20.3.2.2 Crime and Punishment. Spammers keep on spamming because risks are perceived to be small relative to both the potential gains and the other options for getting rich quick. The two people at the center of the C.P. Direct case, Michael Consoli and his nephew and partner, Vincent Passafiume, admitted their guilt in plea agreements signed in August 2003; but they were out of jail before May 2004, and they seemed to have suffered very little in the shame department. Two years after their release, the pair asked the state Court of Appeals to overturn the convictions and give back whatever was left of the seized assets.

Consider what has happened with Jeremy Jaynes, named as one of the world's top 10 spammers in 2003 by Spamhaus, a spam-fighting organization discussed later in Section 20.4.1 of this chapter. When he was prosecuted for spamming, Jaynes was thought to be sending out 10 million emails a day. How much did he earn from this? Prosecutors claimed it was about \$750,000 a month. In 2004, Jaynes was convicted of sending unsolicited emails with forged headers, in violation of Virginia law. He was sentenced to nine years in prison. However, in September 2008, Jaynes, who had not served any time and remained free on bail (set at less than two months' worth of his spam earnings) had his conviction overturned by the Virginia Supreme Court, which ruled that the state antispam law was unconstitutional⁸—a ruling attacked as ill-founded by some critics specializing in antispam legislation.⁹

Another notorious spammer was Sanford "Spamford" Wallace, founder of Cyber Promotions in the 1990s, which actively used spam as a commercial service. In October 2009, Facebook won a civil suit against him for sending fraudulent messages to its users.¹⁰ Wallace was ordered to pay a \$711 million fine—which he was unlikely to pay because he filed for bankruptcy in June 2009.

In August 2011, he was indicted in the federal court in San Jose on

... multiple counts of fraud and related activity in connection with electronic mail. Wallace was also charged with three counts of intentional damage to a protected computer and two counts of criminal contempt. According to the indictment, from approximately November 2008 through March 2009, Wallace executed a scheme to send spam messages to Facebook users. Those messages compromised approximately 500,000 legitimate Facebook accounts, and resulted in more than 27 million spam messages being sent through Facebook's servers. The indictment alleges that Wallace sent spam messages to Facebook users during three time periods: First, on or about Nov. 5, 2008, and continuing until approximately Nov. 6, 2008, Wallace accessed Facebook's computer network to initiate the transmission of a program that resulted in more than 125,000 spam messages being sent to Facebook users; Second, on Dec. 28, 2008, Wallace accessed Facebook's computer network to initiate the transmission of a program that resulted in nearly 300,000 spam messages being sent to Facebook users; Third, on Feb. 17, 2009, Wallace accessed Facebook's computer network to initiate the transmission of a program that resulted in more than 125,000 spam messages being sent to Facebook users.¹¹

At the time of writing (May 2013), USA v. Sanford Wallace was scheduled to open on Monday June 3, 2013, in the court of Judge Edward J. Davila in the San Jose Courthouse.¹²

20 · 10 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

20.3.2.3 A Wasteful Game. In many ways, spam is the heir to the classic sucker's game played out in classified advertisements that promise to teach you how to "get cash in your mailbox." The trick is to get people to send you money to learn how to get cash in their mailbox. If a spammer sends out 25 million email messages touting a product, she *may* sell enough products to make a profit, but there are real production costs associated with a real product. In contrast, if she sends out enough messages touting a list of 25 million proven email addresses for \$79.95, she may reel in enough suckers willing to buy that list and make a significant profit, because the list costs essentially nothing to generate. The fact that most addresses on such lists turn out to be useless does not seem to stop people from buying or selling them.

One form of spam that does not rely on selling a product is the *pump-and-dump* stock scam. Sending out millions of messages telling people how hot an obscure stock is about to become can create a self-fulfilling prophecy. Here is how it works:

- Buy a lot of shares on margin or a lot of shares in a company that is trading at a few pennies a share.
- Spam millions of people with a message talking up the stock you just bought.
- Wait for the price of the shares to go up, and then sell your shares for a lot more than you paid for them.

Such a scheme breaks a variety of laws but can prove profitable if you are not caught (you can hide your identity as a sender of spam); nevertheless, regulators still can review large buy and sell orders for shares referred to in stock spams.

The underlying reasons for the continuing rise in spam volume can be found in the economics of the medium. Sending out millions of email messages costs the sender very little. An ordinary personal computer (PC) connected to the Internet via a \$10-per-month dial-up modem connection can pump out hundreds of thousands of messages a day; a small network of PCs connected via a \$50-per-month cable modem or digital subscriber line (DSL) can churn out millions. Obviously, the economic barrier to entry into get-rich-quick spam schemes is very low. The risk of running into trouble with the authorities is also very low. The costs of spam are borne downstream, in a number of ways, by several unwilling accomplices:

Email Recipient

- Spends time separating junk email from legitimate email. Unlike snail mail, which typically is delivered and sorted once per day, email arrives throughout the day and night. Every time you check it, you face the time-wasting distraction of having to sort out spam.
- Pays to receive email. There are no free Internet connections. When you connect to the Internet, somebody pays. The typical consumer at home pays a flat rate every month, but the level of that rate is determined, in part, by the volume of data that the Internet service provider (ISP) handles, and spam inflates that volume, thus inflating cost.

Enterprise

- Loses productivity because employees, many of whom must check email for business purposes, are spending time weeding spam out of their company email

3 SPAM DEFINED 20 · 11

inbox. Companies that allow employees to access personal email at work also pay for time lost to personal spam weeding.

- Wastes resources because spam inflates bandwidth consumption, processing cycles, and storage space.

Internet Service Providers (ISPs) and Email Service Providers (ESPs)¹³

- Wastes resources on handling spam that inflates bandwidth consumption, processing cycles, and storage space.
- Have to spend money on spam filtering, block list administration, spam-related customer complaints, and filter/block-related complaints.
- Have to devote resources to policing their users to avoid getting block-listed. (There is more about filters and block lists in the next section.)

Two other economic factors are at work in the rise of spam: hard times and delivery rates. When times are tough, more people are willing to believe that get-rich-schemes like spamming are worth trying, so there are more spammers (tough times affect the receiving end as well, with recipients more willing to believe fraudulent promises of lotteries won and easy money to be made). When delivery rates for spam go down—due to the use of spam filters and other techniques that are discussed in Section 20.4.5—spammers compensate by sending even more spam.

20.3.2.4 How Big Is the Spam Problem? There is some controversy about the proportion of email classified as spam. By 2006, spam was said to be consuming over 90 percent of email resources worldwide.¹⁴ This was a staggering level of system abuse by any standard. However, when a handful of security professionals had claimed, a decade earlier, that spam was a computer-security threat, they were met with considerable skepticism and some suspicion, perhaps in part because of the anticlimax of Y2K; also, there was doubtless an element of the recurrent suspicion that security professionals trumpet new threats to drum up business—a strange notion, given the perennial abundance of opportunities for experts in a field that persistently reports near-zero levels of unemployment.

A useful historical perspective on spam's depressing impact on email is provided by the reports freely available at the MessageLabs Website. Its annual reports provide conservative estimates of the growth in spam; for example, the 2007 annual report¹⁵ showed total spam hovering around 85 percent of all emails from 2005 through 2007, with new varieties (those previously unidentified by type or source) keeping fairly steady at around 75 percent of all email.

Symantec reported about spam as over 90 percent of total email traffic in 2009¹⁶ and 2010,¹⁷ but the proportion began dropping over the next few years. In June 2011, Symantec reported a spam rate of about 73 percent of total email¹⁸; by December 2011, they reported a further drop to about 70 percent.¹⁹ According to Kaspersky Lab in early 2013, "... the share of spam in email traffic decreased steadily throughout 2012 to hit a five-year low. The average for the year stood at 72.1%—8.2 percentage points less than in 2011. Such a prolonged and substantial decrease in spam levels is unprecedented."²⁰ A different study in January 2013 suggested that only about 60 percent of all email was spam in 2012.²¹ Estimates by other experts suggested that only about 15 percent of the total spam was getting through all the spam filters at ISP and application levels.²²

20 · 12 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

Although most companies and consumers probably will concur that spam has grown from a mere annoyance to a huge burden, some people will say spam is not a big problem. These include:

- Consumers who have not been using email for very long
- Office users who do not see the spam addressed to them due to some form of antispam device or service deployed by their company

These perceptions obscure the fact, noted earlier, that spam consumes vast amounts of resources that could be put to better use. The consumer might get cheaper, better Internet service if spam was not consuming so much bandwidth, server capacity, storage, and manpower. In the author's experience working with a regional ISP, growth in spam volume is directly reflected in server costs. The ISP had to keep adding servers to handle email. By the time it got to four servers, 75 percent of all the email was determined to be spam. The ISP had incurred server costs four times greater than needed to handle legitimate email. Furthermore, even with four servers, spikes in spam volumes were causing servers to crash, incurring the added cost of service calls in the middle of the night, not to mention the loss of subscribers annoyed by outages.

Spam has a direct effect on infrastructure spending. Storage is one of the biggest hardware and maintenance costs incurred by email. If even 70 percent of all email is unsolicited junk, then companies that process email are spending much more on storage than they would if all email was legitimate. The productivity hit for companies whose employees waste time weeding spam out of their inboxes is also enormous. Beyond these costs, consider the possibility that spam actually acts as a brake on Internet growth rates, having a negative effective on economies, such as America's, that derive considerable strength from Internet-related goods and services.

Although the Internet appears to grow at a healthy clip, it may be a case of doing so well that it is hard to know how badly we are doing. Perhaps the only reason such an effect has not yet been felt is that spam's impact on new users is limited. As noted earlier, when new users first get email addresses, they typically do not get a lot of spam. According to some studies, it can take six to twelve months for spammers to find an email address, but when they do, the volume of spam to that address can increase very quickly. This leads some people to cut back on their use of email and the Internet.

In "The Economics of Spam," published in the Summer 2012 issue of the *Journal of Economic Perspectives*, Justin M. Rao (Microsoft) and David H. Reiley (Google), both formerly employees of Yahoo! Research, discuss the *externality* of spam—the use of victims' resources to support profit for the criminals. They write,

We estimate that American firms and consumers experience costs of almost \$20 billion annually due to spam. Our figure is more conservative than the \$50 billion figure often cited by other authors, and we also note that the figure would be much higher if it were not for private investment in anti-spam technology by firms.... On the private-benefit side, based on the work of crafty computer scientists who have infiltrated and monitored spammers' activity... we estimate that spammers and spam-advertised merchants collect gross worldwide revenues on the order of \$200 million per year. Thus, the "externality ratio" of external costs to internal benefits for spam is around 100:1.²³

20.3.3 Spam's Two-Sided Threat. When mail servers slow down, falter, and finally crash under an onslaught of spam, the results include lost messages, service interruptions, and unanticipated helpdesk and tech support costs. Sales are missed.

3 SPAM DEFINED 20 · 13

Customers do not get the service they expect. The cost of keeping spam out of your inbox and your enterprise is one thing, the cost of preventing spam from impacting system availability and business operations is another; but there is another side of the spam threat, the temptation to become an abusive mass mailer, otherwise known as a spammer. This is something that spam has in common with competitive intelligence, otherwise known as industrial espionage. When done badly, mass mailings, like competitive intelligence, can tarnish a company's reputation.

Leaving aside the matter of the increasingly nasty payloads delivered with spam, things like Trojans and phishing attacks, even plain old male enhancement spam constitutes a threat to both network infrastructure and productivity. For a start, spam constitutes a theft of network resources. The inflationary effect on server budgets has already been mentioned. The negative impact on bandwidth may be less obvious but is definitely real. In 2002 and 2003, the author was involved in the beta testing of a prototype network-level antispam router. It was not unusual for a company installing this device to discover that spam had been consuming from two-thirds to three-quarters of its network bandwidth. Whether this impact was seen as performance degradation or cost inflation, very few companies were willing to remove this device once it had been installed. In other words, when companies see what their network performance and bandwidth cost is like when spam is taken out of the equation, they realize just what a negative impact spam has. This is something that might otherwise be hard to detect given that spam has risen in volume over time.

An even more dramatic illustration of the damage that spam can cause comes when a network is targeted by a really big spam cannon (a purpose-built configuration of MTA devices connected to a really big broadband connection; for example, a six-pack of optimized MTAs can fire off 3.6 million messages an hour). The effect can be to crash the receiving mail servers, with all of the attendant cost and risk that involves. One way to prevent this happening, besides deployment of something like an antispammer router, is to sign up for an antispam service, which intercepts all of your incoming email and screens out the spam. Such a solution addresses a number of email-related problems, but at considerable ongoing cost, which still constitutes a theft of resources by spammers.

The company Commtouch.com provides a spam cost calculator that produces some interesting numbers.²⁴ Consider these input values for a medium-size business:

Employees: 800

- Average annual salary: \$45,000
- Average number of daily emails per recipient: 75
- Average percentage of email that is spam: 80 percent

According to the calculator, total annual cost of spam to this organization, which is assumed to be deploying no antispam measures, is just over \$1 million. This is based on certain assumptions, such as the time taken to delete spam messages, but in general, it seems fairly realistic. Of course, the size of the productivity hit caused by the spam that makes it to an employee inbox has been hotly debated over the years, but it is clearly more than negligible and not the only hit. Even if antispam filtering is introduced, there will still be a need to review or adjust the decisions made by the filter to ensure that no important legitimate messages are erroneously quarantined or deleted. In other words, even if the company spends \$50,000 per year on antispam filtering, it will not reclaim all of the \$1 million wasted by spam.

20 · 14 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

20.3.3.1 Threat of Outbound Spam. Even as organizations like the Coalition Against Unsolicited Commercial Email (CAUCE) were trying to persuade companies that spamming was an ill-advised marketing technique that could backfire in the form of very annoyed recipients, and just as various government entities were trying to create rules outlawing spam, some companies were happy to stretch the rules and deluge consumers with email offers, whether those consumers had asked for them or not. This led some antispammers to condemn all companies in the same breath. The author's own experience, working with large companies that have respected brand names, was that none of them actually wanted to offend consumers. Big companies always struggle to rein in maverick marketing activities, and mass emailing is very tempting when a rogue employee or simply an uninformed one are under pressure to produce sales; but upper management is unlikely to condone anything that could be mistaken for spamming.

The motives for corporate responsibility in email are not purely altruistic. Smart companies can see that the perpetuation of disreputable email tactics only dilutes the tremendous potential of email as a business tool. Whatever one thinks of spam, there is no denying that, as a business tool, bulk email is powerful. It is also seductive. When one has a story to tell or a product to sell, and a big list of email addresses is just sitting there, bulk email can be very tempting. Naïve users can find themselves thinking "Where's the harm?" and "Who's going to object?" But unless they have documented permission to send their message to the people on that list, the smart business decision is to resist the temptation. Remember, it only takes one really ticked-off recipient of your unsolicited message to ruin your day.

20.3.3.2 Mass Email Precautions. One of the most basic business email precautions is this: Never send a message unless you are sure you know what it will look like to the person who receives it. This covers the formatting, the language you use, and, above all, the addressing. If you want to address the same message to more than one person at a time, you have three main options, each of which should be handled carefully:

1. Place the email addresses of all recipients in the To field or the Copy (Cc) field so all the recipients will be able to see the addresses of the other people to whom you sent the message. This is sometimes appropriate for communications within a small group of people.
2. If the number of people in the group exceeds about 20, or if you do not want everyone to know who is getting the message, move all but one of them to the Blind Copies (Bcc) field. The one address in the To field may be your own. If the disclosure of recipients is likely to cause any embarrassment whatsoever, do a test mailing first. Send a copy of the message to yourself and to at least one colleague outside the company, and then have the message looked at to make sure the Bcc entries were made correctly.
3. To handle large groups of recipients, or to personalize one message to many recipients, use a specialized application like Group Mail that can reliably build individual, customized messages to each person on a list. This neatly sidesteps errors related to the To and Cc fields. Group Mail stores email addresses in a database and builds messages on the fly using a merge feature like a word processor. You can insert database fields in the message. For example, a message can refer to the recipient by name. The program also offers extensive testing

3 SPAM DEFINED 20 · 15

of messages, so that you can see what recipients will see before you send any messages. And the program has the ability to send messages in small groups, spread over time, according to the capabilities of your Internet connection.

Whether you use a program like Group Mail for your legitimate bulk email, or something even more powerful, depends on several factors, such as the size of your organization, the number of messages you need to send, and your privacy policy. Privacy comes into play because some software used to send email, such as Group Mail, allows the user of the program access to the database containing the addresses to which the mail is being sent. This is generally not a problem in smaller companies, or when the database consists solely of names and addresses without any special context; but it can be an issue when the database contains sensitive information or the context is sensitive. For example, a list of names and addresses can be sensitive if the person handling the list knows, or can infer, that they belong to patients undergoing a certain kind of medical treatment.

You may not want to allow system operators or even programmers to have access to sensitive data simply because they are charged with sending or programming messages. Fortunately, mailing programs can be written that allow an operator to compose mail and send it without seeing the names and addresses of the people to whom it is being sent. Test data can, and should, be used to proof the mailing before it is executed.

Another basic email precaution is never to send any message that might offend any of the recipients. In choosing your wording, your design, your message, know your audience. Use particular caution when it comes to humor, politics, religion, sex, or any other sensitive subject. When using email for business email, it is better to stand accused of a lack of humor rather than a lack of judgment.

Respect people's preferences, if you know them, for content. If people have expressed a preference for text-only messages, do not send them HTML and hope they decide to change their minds. Ask first because the forgiveness-later path is not cost effective when you have to deal with thousands of unhappy recipients calling your switchboard. When you do use HTML content, still try to keep size to a minimum, unless you have recipients who specifically requested large, media-rich messages.

Antispam measures deployed by many ISPs, companies, and consumers sometimes produce false positives, flagging legitimate email as spam, potentially preventing your email from reaching the intended recipients, even when they have asked to receive your email. If your company's email is deemed to be particularly egregious spam, the server through which it is sent is likely to be blocked. If this is your company's general email server, blocking could affect delivery of a lot more than just the spam. If the servers through which your large mailing is sent belong to a service provider, and its servers get blocked, that could be a problem for you too. And if you use a service provider to execute the mailing for you but fail to choose wisely, your mail may be branded as spam just because of the bad reputation of the servers through which it passes. Be aware that using spam to advertise your products could place you in a violation of the CAN-SPAM Act of 2003, even if you do not send the spam yourself. (See Section 20.4.9 for more on CAN-SPAM.)

What can be done to prevent messages that are not spam from falling victim to antispam measures? Adherence to responsible email practices is a good first step. Responsible management of the company's email servers will also help, as will selection of reputable service providers. You should also consider tasking someone with tracking antispam measures, to make sure that your email is designed to avoid, as much as

20 · 16 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

possible, any elements of content or presentation that are currently being flagged as spam.

To make sure your company is associated with responsible email practices, become familiar with the “Six Resolutions for Responsible E-Mailers.” These were created by the Council for Responsible Email (CRE), which was formed under the aegis of the Association for Interactive Marketing (AIM), a subsidiary of the Direct Market Association (DMA). Some of the country’s largest companies, and largest legitimate users of email, belong to these organizations, and they have a vested interest in making sure that email is not abused. Here are the six resolutions:

1. Marketers must not falsify the sender’s domain name or use a nonresponsive IP address without implied permission from the recipient or transferred permission from the marketer.
2. Marketers must not purposely falsify the content of the subject line or mislead readers from the content of the email message.
3. All bulk email marketing messages must include an option for the recipient to unsubscribe (be removed from the list) from receiving future messages from that sender, list owner, or list manager.
4. Marketers must inform the respondent at the time of online collection of the email address for what marketing purpose the respondent’s email address will be used.
5. Marketers must not harvest email addresses with the intent to send bulk unsolicited commercial email without consumers’ knowledge or consent. (Harvest is defined as compiling or stealing email addresses through anonymous collection procedures such as via a Web spider, through chat rooms, or other publicly displayed areas listing personal or business email addresses.)
6. The CRE opposes sending bulk unsolicited commercial email to an email address without a prior business or personal relationship. (Business or personal relationship is defined as any previous recipient-initiated correspondence, transaction activity, customer service activity, third-party permission use, or proven off-line contact.)

These six resolutions may be considered a reasonable middle ground between antispam extremists and marketing-at-all-costs executives. If everyone abided by these resolutions, there would be no spam, at least according to most people’s definition of spam. After all, if nobody received more than one or two messages per week that were unwanted and irrelevant, antispam sentiment would cool significantly.

20.3.3.3 Appending and Permission Issues. Some privacy advocates have objected to the sixth resolution because it permits email *appending*. This is the practice of finding an email address for a customer who has not yet provided one. Companies such as Yesmail and AcquireNow will do this for a fee. For example, if you are a bank, you probably have physical addresses for all your customers, but you may not have email addresses for all of them. You can hire a firm to find email addresses for customers who have not provided one. However, these customers may not have given explicit permission for the bank to contact them via email, so some people would say that sending email messages to them is spamming.

Whether you agree with that assessment or not, several factors need to be assessed carefully if your company is considering using an email append service. First of all,

3 SPAM DEFINED 20 · 17

make sure that nothing in your privacy policy forbids it. Next, think hard about the possible reaction from customers, bearing in mind that email appending is not a perfect science. For more on how appending works, enter “email append” as a search term in Google—you will find a lot of companies offering to explain how they do it, matching data pulled from many different sources using complex algorithms.

Another concern that must be considered is that some messages will go to people who are not customers. For that reason, you probably want to make your first contact a tentative one, such as a polite request to make further contact. Then you can formulate the responses to build a genuine opt-in list. As you might expect, you can outsource this entire process to the append service, which will have its own, often automated, methods of dealing with bounced messages, complaints, and so on.

Do not include any sensitive personal information in the initial contact, since you have no guarantee that bob.jones@majorfreemail.net is the Robert Jones you have on your customer list. When Citibank did an append mailing in the summer of 2002, encouraging existing customers to use the bank’s online account services, it came in for some serious criticism. Although the bank was not providing immediate online access to appended customers, the mere perception of this, combined with a number of mistaken email identities, produced negative publicity. Here is what Citibank said in the message it sent to people it believed to be customers, even though these people had not provided their email address directly to the bank:

Citibank would like to send you email updates to keep you informed about your Citi Card, as well as special services and benefits... With the help of an email service provider, we have located an email address that we believe belongs to you.

Although this message is certainly polite, it clearly raised questions in the minds of recipients. Two questions could undermine appending as a business practice: Where exactly is this service provider looking for these email addresses? Why doesn’t the company that wants my email address just write and ask for it? The fact is, conversion rates from email contact are higher than from snail mail, so the argument for append services is that companies that use them move more quickly to the cheaper and better medium of email than those that do not. The counterpoint is that too many people will be offended in the process.

Consider to what lengths you are stretching the prior business relationship principle cited in the sixth responsible email resolution. A bank probably has a stronger case for appending email addresses to its account holder list than does a mail order company that wants to append a list of people who requested last year’s catalog. The extent to which privacy advocates accept or decry the concept of prior business relationship is largely dependent on how reasonable companies are in their interpretation of it.

Marketing to addresses that were not supplied with a clear understanding that they would be used for such purposes is not advisable. Depending on your privacy statement, it could be a violation of your company’s privacy policy. Going ahead with such a violation could not only annoy customers, but also draw the attention of industry regulators, such as the Federal Trade Commission (FTC). Of course, you will have to decide for yourself if it is your job or responsibility to point this out to management. And be sure to provide a simple way for recipients to opt out of any further mailings. (A link to a Web form is best for this. Avoid asking the recipient to reply to the message. If the email address to which you mailed the message is no longer their primary address, they may have trouble opting out.)

20 · 18 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

20.4 FIGHTING SPAM. The fight against incoming spam can be addressed both specifically, as in “What can I do to protect my systems from spam?”, and in general, “What can be done to prevent spam in general?” Of course, if the activity of spamming were to be eliminated, everyone would have one less threat to worry about.

20.4.1 Enter the Spam Fighters. In the five-year period from 1997 to 2002, a diverse grouping of interests and organizations fought to make the world aware of the problem of spam, and to encourage counter-measures. The efforts of CAUCE, which continue to this day, spurred lawmakers into passing antispam legislation and helped bring order to the blacklisting process by which spam-sending servers are identified. In 1998, the *Spamhaus Project*, a volunteer effort founded by Steve Linford, began tracking spammers and spam-related activity. (The name comes from a pseudo-German expression, coined by Linford, to describe any ISP or other firm that sends spam or willingly provides services to spammers.)

The Anti-Phishing Working Group (APWG) was founded in 2003 and is one of the most active and productive anti-phishing organizations today:

The APWG is a worldwide coalition unifying the global response to cybercrime across industry, government and law-enforcement sectors. APWG's membership of more than 2000 institutions worldwide is as global as its outlook, with its directors, managers and research fellows advising: national governments; global governance bodies like ICANN; hemispheric and global trade groups; and multilateral treaty organizations such as the European Commission, Council of Europe's Convention on Cybercrime, United Nations Office of Drugs and Crime, Organization for Security and Cooperation in Europe and the Organization of American States. Membership is open to financial institutions, retailers, solutions providers, ISPs, telcos, defense contractors, law enforcement agencies, trade groups, treaty organizations and government agencies. APWG member benefits include: clearinghouses of cybercrime event data; cybercrime response utilities for professionals from the private, public, and NGO sectors who combat cybercrime; community-building conferences for cybercrime management professionals; public education utilities for cybercrime prevention; standards development for cybercrime data exchange and programs for promotion of cybercrime research.²⁵

Commercial antispam products started to appear in the late 1990s, starting with spam filters that could be used by individuals. Filtering supplied as a service to enterprises appeared in the form of Brightmail, founded in 1998 and now owned by Symantec, and Postini, founded in 1999 and now owned by Google. A host of other solutions, some free and open source, others commercial, attempted to tackle spam from several different directions. Nevertheless, despite the concerted efforts of volunteers, lawmakers, and entrepreneurs, spam has continued to sap network resources while evolving into part of a subculture of system abuse that includes delivering payloads that do far worse things than spark moral outrage.

20.4.2 A Good Reputation? Since spam is created by humans, it is notoriously difficult, if not impossible, for computers to identify spam with 100 percent reliability. This fact led some spam fighters to consider an alternative approach: reliably identifying legitimate email. If an email recipient or a Mail Transfer Agent could verify that certain incoming messages originated from legitimate sources, all other email could be ignored. One form of this approach is the challenge-response system, perhaps the largest deployment being the one that Earthlink rolled out in 2003. When you send a message to someone at Earthlink who is using this system, and who has not received a message from you before, your message is not immediately delivered. Instead, Earthlink sends an email asking you to confirm your identity, in a way that

FIGHTING SPAM 20 · 19

would be difficult for a machine to fake. The recipient is then informed of your message and decides whether to accept it or not. Unfortunately, this approach can be problematic if the user does a lot of e-commerce involving email coming from many sources that are essentially automated responders (which cannot pass the challenge). Manually building the whitelist of allowed respondents can be tiresome, and failure to do so may mean some messages do not get through (e.g., if the sender is a machine that does not know how to handle the challenge). One solution to this problem is to compile an independent whitelist of legitimate emailers whose messages are allowed through without question. This is the reputational approach to spam fighting, and it works like this:

- Bank of America pledges never to spam its customers and send them only email they sign up for (be it online statement notification or occasional news of new bank services).
- Bank of America email is always fast-tracked through ISPs and not blocked as spam.
- Bank of America remains true to its pledge because its reputation as a legitimate mailer enables it to conduct email activities with greater efficiency.

There are considerable financial and logistical obstacles to making such a system work, and it tends to work better with bigger mailers. Adoption also faced skepticism from some privacy and antispam advocates who suspected that the goal of such systems was merely to legitimize mass mailings by companies that still did not understand the need to conduct permission-only mailings. The relentless assault of spam on the world's email systems eventually may lead all legitimate companies to foreswear unsolicited email and thus make a reputational system universally feasible. Tough economic times, though, may tempt formerly clean companies to turn to spam in a desperate effort to boost flagging sales.

Another and more universal method of excluding spam would be for ISPs to deliver, and consumers to accept, only those emails that were stamped with a verifiable cryptographic seal of some kind. A relatively simple automated system to accomplish this was developed in 2001 by a company called ePrivacy Group, and it proved very successful in real-world trials conducted by MSN and several other companies. If adopted universally, such a system could render spamming obsolete, but that goal proved to be unattainable. The success of this approach depended on widespread deployment, and the project ultimately was doomed by infighting between the larger ISPs, despite ePrivacy Group's willingness to release the underlying technology to the public domain.

The ultimate in reputation-based approaches to solving the spam problem may well be something that information security experts have urged for years: widespread use of email encryption. If something like Secure/Multipurpose Internal Mail Extensions (S/MIME) was universally implemented, everyone could ignore messages that were not signed by people from whom they were happy to receive email. Of course, the fact that email encryption can be made to work reliably is not proof that it always will, and the encryption approach does not, in itself, solve the fundamental barrier to any universal effort at outlawing spam: the willingness of everyone to participate. One of the qualities that led email to become the most widely used application on the Internet, the lack of central control, is a weakness when it comes to effecting any major change to the way it operates.

20 · 20 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

20.4.3 Relaying Trouble When considering the problem of spam in general, the question is why ISPs allow people to send spam. The fact is, many do not. Spamming is a violation of just about every ISP's terms of service. Accounts used for spamming are frequently closed. However, there are some exceptions. As you might imagine, a few ISPs allow spam as a way to get business (albeit business that few others want). Some of these ISPs use facilities in fringe countries to avoid regulatory oversight. Furthermore, some spammers find it easier, and cheaper, to steal service and use unauthorized access to other people's servers to send their messages. A phenomenon of email is mail relaying. According to Janusz Lukasiak of the University of Manchester, mail relaying occurs "when a mail server processes a mail message from an unauthorized external source with neither the sender nor the recipient being a local user. The mail server is an entirely unrelated third party to this transaction and the message should not be passed via the server."²⁶ The problem with unauthenticated third parties is that they can hide their identity.

In the early days of the Internet, many servers were left *open* so that people could conveniently relay their email through them at any time. It was quite acceptable for individuals to send their email through just about any mail server, since the impact on resources was minimal. Abuses by spammers, whose mass emailings *do* impact resources, led to ISPs instituting restrictions. Some ISPs require the use of port 587 for SMTP authentication. Others require logging in to the POP server to collect incoming mail before sending any out. Since that requires a user name and password, it prevents strangers from sending email through the server.

Open relays are now frowned on. However, relaying continues to occur, partly because configuring servers to prevent it requires effort. Also, spammers keep finding new ways to exploit resources that are not totally protected. For example, port 25 filtering, instituted to reduce spamming, can be bypassed with tricks like asynchronous routing and proxies. A relatively recent development is the use of botnets, groups of compromised computers, to deliver spam. (For more about botnets see Chapters 15, 16, 17, 30, 32, and 41 in this *Handbook*.)

A war is continually being waged between spammers and ISPs, and that war extends to the ISPs' ISP—most Internet service providers actually get their service from an even larger service provider, companies like AT&T, MCI, and Sprint, which are the backbone of the Internet. How much trouble do these companies have with spam? As far back as 2002, network abuse (spam) generated 350,000 trouble tickets each month at a single carrier. These companies are working hard to prevent things like mail relaying and to defeat the latest tricks that spammers have devised to get around their preventive measures.

20.4.4 Black Holes and Block Lists Black hole lists, or block lists, catalog server IP addresses from ISPs whose customers are deemed responsible for spam and from ISPs whose servers are hijacked for spam relaying. ISPs and organizations subscribe to these lists to find out which sending IP addresses should be blocked. The receiving end, such as the consumer recipient's ISP, checks the list for the connecting IP address. If the IP address matches one on the list, then the connection gets dropped before accepting any traffic. Other ISPs choose simply to ignore, or black hole, IP packets at their routers. Among the better known block lists are RBL, otherwise known as MAPS Realtime Blackhole List, Spamcop, and Spamhaus.²⁷

How does an entity's IP address get on these lists? If an ISP openly permits spam, or does not adequately protect its resources against abuse by spammers, it will likely be reported to the list by one or more recipients of such spam. Reports are filed by people

FIGHTING SPAM 20 · 21

who take the time to examine the spam's header, identify the culpable ISP, and make a nomination to a block list. Different block lists have different standards for verifying nominations. Some test the nominated server; others take into account the number of nominations. If an ISP, organization, or individual operating a mail server finds itself on the list by mistake, it can request to be removed, which usually involves a test by the organization operating the block list.

Note that none of this block-list policing of spam is official. All block lists are self-appointed and self-regulated. They set, and enforce, their own standards. The only recourse for entities that feel they have been unfairly blocked—and there have been plenty of these over the years—is legal action. Some block lists are operated outside of the United States, but if an overseas organization block-lists a server located in the United States, it probably can be sued in a U.S. court. However, it is important to note that the blocking is not done by the operator of the block list, it is done by ISPs that subscribe to, and are guided by, the lists.

20.4.5 Spam Filters Block-list systems filter out messages from certain domains or IP addresses, or a range of IP addresses; they do not examine the content of messages. Filtering out spam based on content can be done at several levels.

20.4.5.1 End User Filters Spam filtering probably began at the client level, and even today many email users perform manual negative filtering for spam, identifying spam by a process of elimination. This is easy to do with any email application that allows user-defined filters or rules to direct messages to different inboxes or folders. Many people have separate mailboxes into which they filter all of the messages they get from their usual correspondents, friends, family, colleagues, subscribed newsletters, and so on. This means that whatever is left in the in basket is likely to be spam, with the notable exception of messages from new correspondents for whom a separate mailbox and filter has not yet been created.

To perform a filter that positively identifies spam, elements common to spam messages need to be identified. Many products do this, and they usually come with a default set of filters that look for things like From addresses that contain a lot of numbers and Subject text that contains a lot of punctuation characters—spammers often add these in an attempt to defeat filters based on specific text, so the Subject line “You're Approved” might be randomly concatenated with special characters and spaces like this: “**You~re Ap proved!***”

In fact, the tricks and tweaks used by spammers in crafting messages designed to bypass filters are practically endless. Nevertheless, even the spam filtering built into some basic email programs offers a useful line of defense. For example, in Exhibit 20.3 you can see the control panel for spam filtering in the Eudora email application, which places mail in a Junk folder based on a score derived from common spam indicators, giving users control over how high they want to set the bar.

Unfortunately, some of the language found in spam also appears in legitimate messages, such as “You are receiving this message because you subscribed to this list” or “To unsubscribe, click here.” This means that the default setting in a spam filter is likely to block some messages that you want to get. The answer is to either weaken the filtering or create a whitelist of legitimate From addresses so that the spam filter will allow through anything from these addresses. Most personal spam filters, like the one shown in Exhibit 20.3, can read your address book, and add all of the entries to your personal whitelist. New correspondents can be added over time.

20 · 22 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

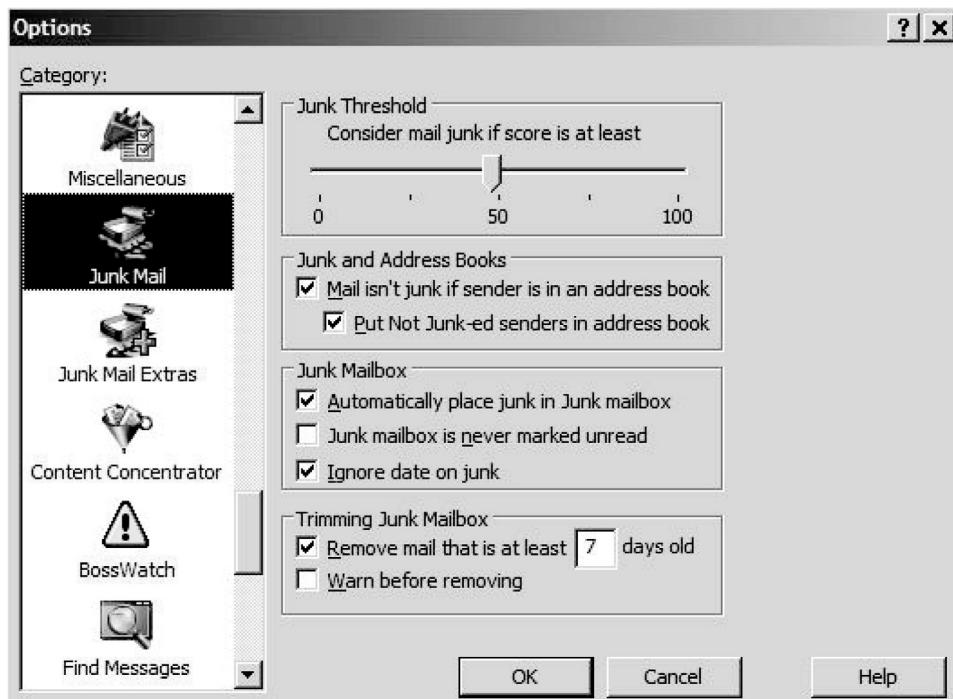


EXHIBIT 20.3 Junk Mail Filter Controls

Most personal spam filters direct the messages they identify as spam into a special folder or in basket where they can be reviewed by the user, a process sometimes referred to as quarantining. To avoid wasting hard drive space, the filtering software can be programmed to empty the quarantine folder at set intervals or simply to delete suspected spam older than a set number of days. This quarantine approach gives the user time to review and reclaim wrongly suspected spam.

20.4.5.2 ISP Filtering Faced with customer complaints about increased levels of spam around the turn of the century, ISPs began to institute spam filtering. However, they were hesitant to conduct content-based spam filtering, fearing that it might be construed as reading other people's email and lead to claims of privacy invasion. Yet ISPs must be able to read headers to route email, so filtering on the From address and Subject field was introduced (hence the increasingly inventive attempts by spammers to randomly vary Subject text). Some ISPs found that distaste for spam reached a level at which some users were prepared to accept revised terms of agreement, to allow machine reading of email content to perform more effective spam filtering. Sticky legal issues still exist for ISPs, however, especially since there is no generally accepted definition of spam and no consensus on the extent to which freedom of speech applies to email. For example, do political candidates have a right to send unsolicited email to constituents? Do ISPs have a right to block it? These are questions on which the courts of law and public opinion have yet to render a conclusive verdict.

One place where content-based spam filters are being deployed with little or no concern for legal challenges is the corporate network. Based on the fact that the company network belongs to the company, the right to control how it is used trumps

FIGHTING SPAM 20 · 23

concerns over privacy. Employees do not have the right to receive, at work, whatever kind of email they want to receive. And most companies would argue that there is virtually no corporate obligation to deliver email to employees. For more details of email policies, see Chapter 48 in this *Handbook*.

20.4.5.3 Filtering Services Companies like Brightmail and Postini arose in the late 1990s to offer filtering services at the enterprise level. Brightmail developed filters that sit at the gateway to the enterprise network and filter based on continuously updated rules derived from real-time research about the latest spam outbreaks. Postini actually directs all of a company's incoming email to its servers and filters it before sending it on. With specialization, filtering services can utilize a wide range of spam fighting techniques, including:

- Block lists
- Header analysis
- Content analysis
- Filtering against a database of known spam
- Heuristic filtering for spamlike attributes
- Whitelisting through reputational schemes

Because these systems constantly deal with huge amounts of spam, they are able to refine their filters, both positive and negative, quickly and relatively effectively. And they catch a lot of spam that would otherwise reach consumers. For example, by 2008, nine of the top 12 ISPs were using Brightmail, which claims a very low false positive rate and 95 percent catch rate. Unfortunately, that still means that some spam gets through, and when spammers can direct millions of messages an hour at a network, the volume of what gets through still can pose a threat.

Another approach uses the collective intelligence of subscribers to identify spam quickly and spread the word via network connections. Cloudfilter, for example, had several million subscribers at the time of writing (April 2008) paying about \$40 per year to receive nearly instantaneous updates, from servers that receive and categorize reports from members on any spam that gets through. A member's credibility score rises with every correct identification of spam and sinks with incorrect labeling of legitimate email as spam (e.g., newsletters the member has forgotten subscribing to). The reliability of reporters helps the system screen false information from spammers that might try to game the system by claiming that their own junk email was legitimate.

20.4.5.4 Collateral Damage There are two major drawbacks to spam filters. First, they allow spammers to continue spamming. In other words, because they must decide, on a message-by-message basis, which messages are spam and which are not, filters consume a lot of resources, in some cases more than if all the spam was allowed through. When block lists and whitelists work well, they can substantially reduce the amount of spam that reaches the filtering stage, but eventually all filters are serial, and thus resource constrained and resource intensive.

Second, filters sometimes err, in one of two ways. They sometimes produce a false positive, flagging a legitimate message as spam, preventing a much-needed message from getting to the recipient in a timely manner. And they sometimes produce false

20 · 24 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

negatives, allowing spam messages into the inbox. False negatives and false positives impose a drag on productivity.

20.4.5.5 Snowshoe Spamming It was inevitable that spammers would adapt to the filtering measures described above. By late 2012, an increasing number of spammers were distributing their spam among a large number of compromised servers, leading to the term *snowshoe* spamming. McAfee Labs summarized the problem as follows:

Snowshoe spamming is now one of the biggest spam problems. The issue has exploded over the past two years and will continue to increase sharply due to lack of exposure by law enforcement authorities and threats of lawsuits by companies using the illegal email lists.

The phenomenon is characterized by the following:

- Spammers blast out millions and millions of blatantly illegal spam messages every day from newly rented hosts until they get evicted from their subnetworks or move on.
- Recipients have their inboxes bombarded with these spam messages and are unable to opt out of them because they are not sent from a legitimate source.
- The result of snowshoe spamming is permanently blacklisted addresses and sometimes subnetworks.
- Because spamming is seen as simply annoying rather than malicious, authorities have largely ignored this problem, despite the growing volumes of unwanted email originating from these sources. Companies using these shady marketers have threatened to file defamation lawsuits when researchers have tried to expose this activity.²⁸

20.4.6 Network Devices The rapid rise of spam volumes in the late 1990s occurred at a time of massive email-enabled computer virus and worm attacks, arising from the classic malware motivation of bragging rights. Analysis of these problems, according to the classic factors of means, motive, and opportunity, revealed that spammers differed quite significantly from virus writers. Spammers are mostly motivated by money, not bragging rights. This insight opened up a new line of defense against spam, removing that motivation.²⁹

All that remained was to understand how spammers make money, that is, the economics of spam, then to find a way to disrupt those economics. The classic spam model is to send out a vast number of emails offering a product or service, relying on the fact that at least some of these offers will reach real people, at least some of whom will make a purchase. If enough sales are made to create a profit over and above the cost of doing business, the spammer will continue to operate. Although a good public education and awareness program can reduce the number of people who buy the spammers' offers, it is unlikely to reduce that number enough. Because filtering spam out of the message stream reduces the spammer's return on investment, it tends to make spammers more inventive in their attempts to beat filters and to send even more spam in the hopes that enough will get through. In fact, spammers were able to find companies willing to sell them bandwidth on a massive scale, including dedicated point-to-point T3 digital phone line connections.

Another way to attack the economics of spam is by following the money. Every deal has to be closed, typically via a Website. Is it possible to shut down the spammer's Websites? This line of inquiry led the author to an interesting finding: Spam goes stale very quickly. An examination of spam archives showed that most links in old spam were dead, sometimes because the Website was shut down by the hosting company,

FIGHTING SPAM 20 · 25

sometimes because the spammer did not want to risk being identified. Thus, the key to the economics of spam was revealed: time. If you slow spam down, spammers cannot generate enough responses before response sites are taken down.

Shortly after this realization, network security expert David Brussin devised a means of slowing down spam through TCP/IP traffic shaping. This became the technology at the heart of the antispam router. Instead of looking at each message in turn to determine if it is spam, the antispam router samples message traffic in real time and slows that traffic down if the sample suggests the traffic contains spam. To spammers, or rather the spamming software used by spammers, a network protected with an antispam router behaves as if it is on a 300-baud modem. In other words, the connection is too slow to deliver enough messages quickly enough for that one-in-a-million hit that the spam scheme needs in order to make money before the Website is shut down. The spamming software quickly drops the connection; no messages have been lost and no messages have been falsely labeled as spam. Delivery of legitimate email has not been affected. One or two spam messages have been delivered, but a lot less than are allowed through by the 95 percent catch rate of a typical spam filter.

There are two main tricks to this technology. One is tuning the TCP/IP traffic shaping, the other is tuning the sampling process. The latter needs regular updates about what spam looks like currently, so that it can identify a spammy connection as quickly as possible. These updates can be derived from services that constantly identify new spam. An application control panel is used to handle the tuning of the traffic shaping. Network administrators have found that a properly tuned antispam router can reduce bandwidth requirements by as much as 75 percent. Unfortunately, an antispam router works best on a high-volume connection, protecting MTAs at ISPs, larger companies, schools, and government agencies. There is no desktop version. Still, the technology has become a valuable part of the antispam arsenal, although the scourge of spam still continues.

20.4.7 Email Authentication On a technical level, what allows spam to continue to proliferate is the lack of sender authentication in the SMTP protocol. Several initiatives have sought to change this situation, either by changing the SMTP protocol or by adding another layer to it. One obvious way to authenticate email is via the sender's domain name, and numerous schemes to accomplish this have been proposed:

- **Sender Policy Framework (SPF).** An extension to SMTP that allows software to identify and reject forged addresses in the SMTP MAIL FROM (Return-Path) that are typically indicative of spam. SPF is defined in Experimental RFC 4408.³⁰
- **Certified Server Validation (CSV).** A technical method of email authentication that focuses on the SMTP HELO-identity of MTAs.
- **SenderID.** An antspoofing proposal from the former MARID IETF working group that joined Sender Policy Framework and Caller ID. Sender ID is defined primarily in Experimental RFC 4406.³¹
- **DomainKeys.** A method for email authentication that provides end-to-end integrity from a signing MTA to a verifying MTA acting on behalf of the recipient. It uses a signature header verified by retrieving and validating the sender's public key through the Domain Name System (DNS).
- **DomainKeys Identified Mail (DKIM).** A specification that merges DomainKeys and Identified Internet Mail, an email authentication scheme supported by Cisco, to create an enhanced implementation.

20 · 26 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

When properly implemented, all of these methods can be effective in stopping the kind of forgery that has enabled spam to dominate email. Each has varying pros and cons and requirements. SPF, CSV, and SenderID authenticate just a domain name. SPF and CSV can reject forgeries before message data is transferred. DomainKeys and DKIM use a digital signature to authenticate a domain name and the integrity of message content. In order for SenderID and DomainKeys to work, they must process the headers, and so the message must be transmitted.

Obviously a scheme like DomainKeys makes email more complex than SMTP. Additional data and processing effort are required. DNS records have to be expanded and maintained. However, the effort may be worth it. Although something like DomainKeys does not prevent email abuse, it does make abuse of email on protected domains easier to detect and track, and this is a deterrent to spamming now that several stiff sentences for spamming have been handed down.³² Since 2004, the email services of both Yahoo and Google have signed outgoing email with DomainKeys, but it is not yet clear if any one of the described schemes, or some derivative thereof, will become the standard for all email. And therein lies the rub. Spam could be effectively sidelined *if* all legitimate email was reliably authenticated; recipients could simply ignore all unauthenticated messages. However, that is a large *if*.

20.4.8 Industry Initiatives Back in 2002, at the urging of consumer advocates, government agencies, and many large corporate users of email, the major email service providers started to hold meetings to discuss a unified approach to improving email and eliminating spam. There were high hopes that the problem of spam would be solved by this group, which was known in the business as AMEY, for AOL, Microsoft, Earthlink, and Yahoo!. There was certainly plenty of encouragement. Early in 2003, the FTC, America's main consumer protection agency, convened a conference on spam, and the commissioners made it clear they wanted action. Indeed, by the end of 2003, Congress had passed the first federal antispam legislation (see Section 20.4.9).

In January 2004, then Microsoft chairman and chief executive Bill Gates announced to the World Economic Forum, “Two years from now, spam will be solved.” Two months later, the AMEY companies filed lawsuits against persons alleged to be major spammers, claiming that spam cost businesses in North America \$10 billion each year in lost productivity, network upgrades, and destroyed or lost data. However, by January 2006, spam constituted more than three-quarters of all email. What went wrong? Repeated efforts to get the AMEY companies to adopt a new, open, royalty-free standard have run aground on concerns over intellectual property rights. Both Microsoft and Yahoo! have asserted ownership over parts of the various mechanisms put forward to authenticate email and thus solve the spam problem. Ironically, a lack of trust between the large ISPs has resulted in a lack of trust of the email they deliver.

20.4.9 Legal Remedies Laws against spam have been passed by many countries, including the United States. The CAN-SPAM Act of 2003 (Controlling the Assault of Non-Solicited Pornography And Marketing Act, effective January 1, 2004) established certain requirements for anyone sending commercial email. The law also provides penalties for spammers and for companies whose products are advertised in spam. The law addresses “email whose primary purpose is advertising or promoting a commercial product or service, including content on a Website.”³³

Enforcement of the CAN-SPAM Act is primarily the job of the FTC, but the act also gives the Department of Justice authority to enforce the act’s criminal sanctions. Other federal and state agencies can enforce the law against organizations under their

PHISHING 20 · 27

jurisdiction, and ISPs can also sue CAN-SPAM violators. Here are the law's main provisions, as stated by the FTC:

- **Bans false or misleading header information.** The From, To, and routing information of your email—including the originating domain name and email address—must be accurate and identify the person who initiated the email.
- **Prohibits deceptive subject lines.** The subject line cannot mislead the recipient about the contents or subject matter of the message.
- **Requires that your email give recipients an opt-out method.** You must provide a return email address or another Internet-based response mechanism that allows a recipient to ask you not to send future email messages to that email address, and you must honor the requests. You may create a menu of choices to allow a recipient to opt out of certain types of messages, but you must include the option to end any commercial messages from the sender.

Any opt-out mechanism you offer must be able to process opt-out requests for at least 30 days after you send your commercial email. When you receive an opt-out request, the law gives you 10 business days to stop sending email to the requestor's email address. You cannot help another entity send email to that address or have another entity send email on your behalf to that address. Finally, it is illegal for you to sell or transfer the email addresses of people who choose not to receive your email, even in the form of a mailing list, unless you transfer the addresses so another entity can comply with the law.

- **Requires that commercial email be identified as an advertisement and include the sender's valid physical postal address.** Your message must contain clear and conspicuous notice that the message is an advertisement or solicitation and that the recipient can opt out of receiving more commercial email from you. It also must include your valid physical postal address.

Note that messages which are transactional or relationship based—that is, email that “facilitates an agreed-upon transaction or updates a customer in an existing business relationship”—are also covered in the sense that these messages may not contain false or misleading routing information.

Readers will find extensive discussion about whether CAN-SPAM has been—and could ever be—effective by searching the Web with any search engine. For a critique of the law published in February 2004, see “Can CAN-SPAM can spam?”³⁴

20.5 PHISHING. Phishing is the use of unsolicited commercial email to obtain—electronically fish for—information about you, information that you would not normally disclose to a stranger, such as your bank account number, personal identification number (PIN), and other personal identifiers such as Social Security number. There was virtually no phishing activity before 2002 and relatively little in 2003, but by 2004 phishing had become an everyday threat and has continued unabated ever since. This is one category of threat that could have been squelched, along with the rest of spam, if industry leaders had chosen to cooperate rather than to compete for customers based on promising “better antispam than the other guys.”

20.5.1 What Phish Look Like. Most of us first became aware of *phishing* when we received an email about a problem with an account, typically a bank account, but possibly an eBay, PayPal, Amazon, or other online account. A typical phishing

20 · 28 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

message is designed to look as if it had been sent by a large enterprise, such as the Bank of America, complete with accurate copies of the company logo, typeface, and terminology. At first glance, such messages can be quite convincing. If you receive a message like this that references an institution at which you *do* have an account, you might be tempted to read it. You might even be tempted to follow the instructions, and these are likely to lead you to a Website that asks for confidential information. Consider the text of a typical phishing message:

You are receiving this message, due to your protection, Our Online Technical Security Service Foreign IP Spy recently detected that your online account was recently logged on from am 77.32.11.84 without am International Access Code (I.A.C) and from an unregistered computer, which was not verified by the Our Online Service Department.

If you last logged in your online account on Monday May 5th 2007, by the time 6:45 pm from an Foreign Ip their is no need for you to panic, but if you did log in your account on the above Date and Time, kindly take 2-3 minute of your online banking experience to verify and register your computer now to avoid identity theft, your protection is our future medal.

What looks, at first glance, to be a very official and technical message, turns out to be full of errors. The best thing to do with these messages is delete them. The worst thing to do is respond to them. Until there is a major overhaul of email security, no reputable institution will be using email to request changes or updates to confidential account information.

If you received an email like the one discussed above but did not have a Bank of America account, you might have been confused. The fact is, the people sending these phishing messages usually have no idea whether the recipients have accounts at the institution named in the message. Indeed, this type of mismatch is the easiest way to spot some phishing messages. However, if the phisher, who has likely sent out millions of copies of the same email, gets lucky and you do happen to have an account at the named institution, or if the email is generic, then things get a little trickier. You are more likely to open a message that appears, often quite convincingly, to come from your bank. If you cannot resist looking at email about an account problem, here are some of the clues that the message is bogus. (Note that we are not implying that messages lacking these clues are therefore legitimate.)

- **Deceptive link.** Most phishing messages make an effort to look like they are legitimate, for example, by using logos and graphics stolen from the Website of the targeted institution. All phishing messages we have seen also include a link to a Website, where you are asked to provide the data the phisher is trying to steal. However, this link typically is disguised. For example, the link might be long and complicated and include the name of the bank but actually not take you to the bank's Website. Alternatively, the link may appear to be plain and simple text but in fact it is HTML-coded to go somewhere else. Some email programs, such as Eudora, will warn you of this deception and show you the real link when you place your mouse over the link text prior to clicking.

The design of the link in a phishing scheme can be critical to its success. An increasing number of users are appropriately wary of long, complex, or merely numerical IP address links in email (e.g., "Click <http://123.212.192.68>"). By the use of HTML coding, messages typically obscure the destination address. However, the link is also relevant in the next stage of the attack, when the victim clicks the link. If the URL of the phishing site appears as a numerical address in the

PHISHING 20 · 29

URL field of the browser, the victim may become suspicious. Various techniques are used to make this address look plausible.

- **Change PIN/Password.** Emails that ask you to change your account access credentials are highly suspect. Legitimate companies do not make such respects via email precisely because email is so unreliable. No security update to a reputable banking Website is going to ask that you log in to your account to reset it or to prevent its being suspended. And why would a legitimate message ask you to use a password *that has supposedly been compromised?* No government agency is going to ask for your credentials or personal information via email. And no lottery on Earth uses email to notify winners. Ignore these messages, or report them as spam, and move on.
- **Bad spelling, grammar, and logic.** Whoever thought those tedious grammar lessons could be so useful? Bad grammar, spelling, and even faulty logic can be the fastest way to spot bogus email. Consider this example: “Therefore, if you are the rightful holder of the account please fill in the form below so that we can check your identy.[sic]” There is a telling typo here (*identy* for *identity*), and the logic is hopeless. Think about it: Why would a bank send an email to someone if it was not sure the person was the rightful holder of the account? Again, this is not how real companies do business today, so just move on.
- **Generic phish.** Generic account warnings are particularly nasty. They are one way that phishing attacks try to get around the problem of not knowing where the victim (you) has an account. For example, all bank accounts in America are insured by an institution called the Federal Deposit Insurance Corporation (FDIC). In 2004 someone created a particularly nasty attack that preyed on this fact, potentially snagging anyone with a bank account. This was one of the first phishing attacks to fake the linked URL using a vulnerability in Microsoft Internet Explorer to mask the actual URL. If you clicked the link in the message, the site you went to really looked as if it were www.fdic.gov but actually led to the phisher’s site.

20.5.2 Growth and Extent of Phishing. Any new form of computer abuse can be said to have arrived when a toolkit facilitating the abuse becomes available. Phishing toolkits first appeared in 2006, offering scripts that enabled attackers to automatically set up phishing Websites spoofing the legitimate Websites of different brands (including the illegal appropriation of images and logos that consumers associate with those brands). These scripts helpfully generate corresponding phishing email messages. Criminals continue the evolution of their tools without respite; a January 2013 report found that spear phishers have applied a new tool called Bouncer that adapts their URLs to include unique identifiers for their intended victims; attempting to access the criminals’ pages without a valid identifier in the URL results in a 404 (no such page) error, thus interfering with researchers’ analysis of the phishing pages.³⁵

In February 2013, analysts found that “When security experts looked into some of the highest profile hacks in recent years—one particular criminal group kept on coming to their attention. The Comment Group, which industry insiders say is based in China, offer hacking for hire—be it for individuals, corporations or governments...” They research individual companies or organizations to locate detailed information that allows highly specific topics and even content in the phishing messages. For example, a Coca-Cola executive reportedly opened a phishing email supposedly from his own boss; the link he clicked on downloaded spyware into his computer and allowed Chinese

20 · 30 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

industrial spies to extract information which stymied the acquisition of China's largest soft-drinks company.³⁶

In April 2013, the U.S. Department of Homeland Security (DHS) warned "organizations that post a lot of business and personal information on public web pages and social media sites" not to do so. In October 2012, phishers harvested detailed information about employees from a public posting by an energy company listing attendees at a conference. In addition to spear-phishing attacks on named individuals, "Malicious emails that appeared to be from one of the attendees were sent to others on the list informing them of a change in the sender's email address. Recipients were politely asked to click on an attached link that promptly took them to a site containing malware."³⁷

The APWG is a valuable repository of statistical information about phishing.³⁸ Their report on phishing during the fourth quarter of 2013 includes the following findings (these are direct quotations formatted as bullet points with page references removed):

- Phishing attacks against online game players saw a massive increase, from 2.7 percent of all phishing attacks in Q3 to 14.7 percent in Q4.
- Financial services continued to be the most targeted industry sector in the fourth quarter of 2012, with payment services close behind.
- Online gaming credentials are valuable to certain criminals, who sell them on the black market. In-game items held in those accounts can also be sold by phishers for real-world cash. Victims can even have their real-life identities stolen.
- Attacks against social media doubled to 6 percent, up from 3 percent in the third quarter.
- During Q4, about 30 percent of personal computers worldwide were infected with malware.
- More than 57 percent of PCs in China may have been infected, while PCs in European nations were infected least often.
- Except for October 2012, the number of phishing sites declined every month from April 2012 through December 2012.
- April 2012 saw 63,253 unique phishing sites detected, falling to 45,628 in December 2012.
- The APWG received reports of 28,195 unique phishing sites in December. December's total was 31 percent lower than the high of 40,621 reports in August 2009.
- Use of crimeware dipped slightly in this quarter from the previous, as did the use of data-stealing malware.
- The use of other malware has increased by a statistically significant amount from the previous quarter.³⁹

20.5.3 Where Is the Threat? Phishing appears at first to be a personal problem, a matter of consumer computer security rather than enterprise information security. Indeed, the federal agency charged with consumer protection, the FTC, has been active in educating consumers about the problem and the ways to avoid getting scammed by these messages.⁴⁰ However, phishing is also a computer-based threat to the information and welfare of companies as well as a threat to e-commerce in general.

TROJAN CODE 20 · 31

Phishing attacks collect usernames and passwords. Many people use the same credentials for work-related systems as they do for personal systems (including the head of a highly secret government facility who was found to be using her bank ATM PIN as her top-secret network password). Criminal hackers are known to have penetrated systems by harvesting personal credentials and applying them to the target's business login. So a company computer security awareness program would do well to include warnings about phishing attacks.

Beyond the penetration threat, phishing can undermine consumer trust in a company. Although it hardly seems fair for consumers to resent Bank of America because a criminal has attempted to defraud them by abusing the bank's identity, such resentments do exist, and they need to be understood in a business context. Automation improves the profitability of banking, which is why banks encourage customers to use online services and offer incentives to drop paper statements and notifications. However, if banks are perceived as providing automation on the cheap, with insufficient safeguards built in to protect customer privacy and account access, some consumers will object, slowing down the pace of automation and jeopardizing productivity increases, potentially impacting the bottom line.

20.5.4 Fighting Phishing Some clever technological defenses specific to phishing have been put forward, such as methods for enabling browsers to verify URLs, but the best defense is a twofold approach that is entirely obvious. The first step is education, at the consumer and corporate level, teaching people how to spot phishing attacks and avoid falling victim to them. The second step is fixing the fundamentals of email in the ways outlined earlier with respect to spam. From a technical perspective, we already have the technology to do this. All that is missing is a willingness on the part of the leading service providers to take concerted action. Perhaps they could be encouraged to do so by the banks and other institutions that stand to gain a great deal if email is made more secure.

Failure to act on email security has already cost billions of dollars in wasted resources and lost productivity, but the effects go further than that, as evidenced by the emerging relationship among phishing, spamming, criminal hacking, Trojan code, personal data harvesting, and account compromise. The application of skills such as virus and worm writing to commercial ends has been enabled by the willingness of spammers to pay for compromised hosts with which to launch attacks and collect dollars and/or data. Spamming techniques have enabled the growth of phishing, which in turn has solidified the black market in compromised hosts and purloined personal data. New types of attack are constantly emerging from this unholy alliance between coders and criminals. Email addresses harvested by spam may be used for spear phishing attacks, where a specific company is targeted through emails sent to known customers or employees (the Department of Defense faced a rash of spear-phishing attacks in 2006). Trojan code distributed by spam methods has been used to corrupt local domain name servers and produce the same effect as a phishing message: Users are misdirected to malicious Websites simply by clicking on an apparently valid URL. (DNS attacks of this type are referred to as pharming.)

20.6 TROJAN CODE. Like the original Trojan horse, deployed by the Greeks to defeat the Trojans who were protected by the impregnable walls of the city of Troy, a computer Trojan is a bad thing disguised as a good thing (where Troy is your computer and the Greeks are the people who would like to gain unauthorized access). The technology has come a long way from the "steed of monstrous height" described

20 · 32 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

EXHIBIT 20.4 Screensaver Trojan

in Virgil's *Aeneid*, but the goal of Trojan code is the same as that of the original Trojan horse: Fool the defender of a protected place into granting access to outsiders. (See Section 15.3.2.17 for more on the origins of the name.)

Almost as soon as the Internet provided a way to distribute executable code, either as download links or email attachments, some people decided to exploit this ability to distribute their code without explicit permission. In other words, an innocent party might be tempted to download an apparently innocuous executable, which was actually something else. The person doing the downloading does so intentionally but is ignorant of the intentions of the person who crafted the Trojan code within the executable.

20.6.1 Classic and Recent Trojans. Despite the fact that the screens of today's computers do not need saving, screensavers continue to be a source of Trojan code. Some users apparently find the promise of animated waterfalls and fish tanks hard to resist. This leads users to download screensavers that they encounter on Websites or open screensaver files that they receive in email. In Exhibit 20.4 you can see an example of one screensaver.

The creator of this Trojan screensaver, discovered in February 2007, was apparently unhappy with a file-sharing network popular in Japan and named Winny. When this screensaver is opened, the Trojan displays an image warning the computer operator not to use Winny. This tactic is reminiscent of some of the very earliest viruses, which sought to spread relatively innocuous messages rather than cause damage. In fact, code does not have to have a malicious intent to be considered a Trojan—it merely needs the intent to get installed without permission. Unfortunately, in this particular case, the Trojan does go on to destroy data by overwriting certain files (e.g., those with .txt and .jpg extensions).

Another example of a screensaver Trojan is the one reported to be circulating as an email attachment called bsaver.zip in July 2007. When the file within the ZIP attachment is opened, the system is infected with the Agent-FZB Trojan horse, which then drops two rootkits to evade detection by security software and to make the system accessible to unauthorized users. (See Chapter 15 in this *Handbook* for more on rootkits.) This

TROJAN CODE 20 · 33

screensaver was distributed by a spam campaign with message subject lines such as *Life is beautiful*, *Life will be better*, *Good summer*, and *help you*. The message text included phrases like “Good morning/evening, man! Really [sic] cool screensaver in your attachment!”

More recent reports about Trojans include the following cases and studies:

- The BackDoor.Wirenet.1 Trojan was identified in August 2012; the malware “is the first Trojan Horse program that works on the Mac OS X and Linux platforms that is ‘designed to steal passwords stored by a number of popular Internet applications.’”⁴¹
- The “PandaLabs Q1 Report” for 2013 found that “Trojans set a new record, causing nearly 80 percent of all computer infections worldwide. Despite their inability to replicate, Trojans are capable of triggering massive infections through compromised Web sites that exploit vulnerabilities in browser plug-ins like Java, Adobe Reader, etc. This attack method allows hackers to infect thousands of computers in just a few minutes with the same Trojan or different ones, as attackers have the ability to change the Trojan they use based on multiple parameters such as the victim’s location, the operating system used, etc.”⁴²
- In March 2013, Kaspersky Labs reported that a new spyware attack on Tibetan freedom activists used a Trojan designed for the Android mobile-phone operating system.⁴³
- The Flashback Trojan had infected more than 600,000 Macintosh computers by early April 2013 by exploiting a flaw in Java.⁴⁴
- In April 2013, criminals sent out invitations to watch video footage of the Boston Marathon bombings. A remote-access Trojan was installed as a “WinPcap Packet Driver (NPF)” to evade notice.⁴⁵
- In May 2013, Graham Cluley of Sophos reported on a Trojan (Mal/BredoZp-B) circulated in an email supposedly from Tiffany & Co. that claimed to include details of an export license and a payment invoice.⁴⁶

A Trojan tries to look useful or interesting so that users will install it, but buried within is unauthorized, undocumented code for unauthorized functions, just some of which are listed here:

- Deleting files or folders or entire drives
- Changing data, in subtle or dramatic ways
- Encrypting data (possibly for purposes of extortion)
- Copying data to other computers (possibly for purposes of industrial espionage)
- Downloading files without user consent (possibly for purposes of illicit e-commerce, illegal file sharing, pornography site hosting, cheap storage)
- Corrupting browser software and network files to redirect user from legitimate Websites to fake, bogus, malicious sites
- Enabling remote access to the compromised system (sometimes referred to as a RAT, for Remote Access Trojan), often used to turn computers into zombies that can be aggregated into botnets (for the purposes of spamming, phishing, and executing denial-of-service attacks)
- Aiding the spread of viruses using dropper code to cause infections

20 · 34 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

- Disabling antivirus and firewall programs
- Disabling competing forms of malware
- Logging keystrokes to obtain data such as passwords and credit card numbers
- Reporting browsing activity (enabling spyware)
- Harvesting email addresses for the purpose of spamming and other activities such as phishing

The motives behind Trojan code typically fall into one or more of three categories: malice, bragging rights, and financial gain. The gains can be realized in several ways. Purloined data can be sold. Access to compromised machines can be sold. Encrypted or stolen data can be ransomed. Denial-of-service threats can be used for extortion. Some companies have even used Trojans to try to increase the sales and installed base of their software.

20.6.2 Basic Anti-Trojan Tactics. Unfortunately, Trojans can be tough to combat. The first line of defense is well-educated users who know better than to execute code of dubious origin. Warnings to this effect should be part of any computer security awareness program, and a computer security awareness program should be part of every enterprise security model. Technical measures can be used to combat Trojans, but because none of these is perfect, it makes no sense to rely on them to the neglect of end user education. In fact, thinking of end users as *computer operators* rather than *computer users* might be a good place to start, given that the role of user implies a relatively passive role in maintaining the integrity and health of a computer, whereas operator reflects more accurately the level of responsibility required of anyone employing a computer in either work or recreation. (A computer user might be likened to someone who merely drives a car and never checks the oil or the tires; an operator is closer to a commercial driver who knows that getting safely and reliably from A to B takes a lot more than just holding the steering wheel and pressing the correct pedal.)

Technical measures against Trojans start with keeping all operating system and application patches up to date, universal use of memory-resident antivirus software that scans incoming files, and regular scanning of the entire system against a regularly updated virus database. At the same time, it helps to run memory-resident antbot software, such as Norton AntiBot, designed to recognize and thwart activity indicative of botnets. This software does not directly prevent Trojans from getting onto a computer, but it minimizes the damage that a Trojan can cause.

A good antispam solution is also necessary because spam is a major vector for Trojan attacks. If Jim in accounting never gets that message about the cool screensaver, he will not ever be tempted to open the attached file. It is not just screensavers that can be tempting. In April 2007, spam techniques were used to distribute, on a massive scale, messages with subject headings like “Worm Alert!” and “Worm Detected.” The messages came with a ZIP file attachment that posed as a patch that would prevent the bogus attack. When recipients, alarmed into action, went to open the ZIP file, they found that it was password protected, and some took that as a sign it was genuine. (The password was included in the message.) Proceeding with installation led to a rootkit, disabling of security software, theft of confidential information from the affected machine, and enrollment in a botnet of compromised computers. In one 24-hour period, Postini counted nearly 5 million copies of this Trojan spam directed at users of its antispam service. The company calculated that this one spam accounted for 87 percent of all malware being spread through email during that time.

TROJAN CODE 20 · 35

20.6.3 Lockdown and Quarantine. More drastic anti-Trojan measures include preventing unauthorized changes or additions to executables and preventing systems from connecting to networks until they have been vetted. The idea of defeating malicious executables by freezing the operating system, and authorized applications in a known good state, goes back a long way, at least to the early days of antivirus efforts. However, the complexity of most operating systems and application code today makes such an approach challenging at best. Consider the emergence of Patch Tuesday as a standard means of maintaining the world's most widely used operating system. Indeed, for many years now, a lot of software development has been predicated on the assumption that patches and updates can always be pushed out if needed, arguably leading to far less rigor in coding practices than when production code was burned to disk, at considerable expense, and changes required further shipments of disks.

The idea of quarantining computers when they attempt to connect to a network, using some form of network access control, offers a slightly different approach to the idea of locking down machines so that unauthorized code cannot be installed. Some enterprise network security managers have pursued this approach because it is increasingly difficult to control some network end points, such as the company laptop that travels with the employee to client sites and conferences, hotels, and WiFi hot spots, even spending time in the employee's car and home. Laptops used in this manner face any number of attack vectors. Traditional defense mechanisms can be applied to these end points, including password protection, biometric authentication, disk encryption, and antivirus programs. However, these may not be enough, as many enterprises have discovered to their cost. So why not prevent these machines from connecting to the enterprise network until they have been scanned for unauthorized executables, like a routine health check? This approach holds promise but is not easy to implement. Furthermore, it leaves machines open to attack while they are away from the corporate network. A spam-distributed phishing exercise or screensaver Trojan could still lobotomize an end point between checkups and compromise confidential personal or corporate data. As you can see from Exhibit 20.5, a table compiled by Symantec in 2011 to show prices paid for various forms of data in the underground economy in 2009 and 2010, the financial incentives for such activities are real.

EXHIBIT 20.5 Table of Prices Paid for Data Traded in the Underground Economy⁴⁷

Overall Rank			Percentage			2010 Price Ranges
2010	2009	Item	2010	2009		
1	1	Credit card information	22%	19%	\$0.07-\$100	
2	2	Bank account credentials	16%	19%	\$10-\$900	
3	3	Email accounts	10%	7%	\$1-\$18	
4	13	Attack tools	7%	2%	\$5-\$650	
5	4	Email addresses	5%	7%	\$1/MB-\$20/MB	
6	7	Credit card dumps	5%	5%	\$0.50-\$120	
7	6	Full identities	5%	5%	\$0.50-\$20	
8	14	Scam hosting	4%	2%	\$10-\$150	
9	5	Shell scripts	4%	6%	\$2-\$7	
10	9	Cash-out services	3%	4%	\$200-\$500 or 50%-70% of total value	

20 · 36 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

20.7 CONCLUDING REMARKS. The threats addressed in this chapter are still evolving, as are the countermeasures discussed here, although some of those mentioned may be rendered obsolete by new developments. If there is one lesson security professionals can learn from studying the threats discussed herein, it is that continued failure to improve the underlying security of Internet email will prolong the onslaught of spam, phishing attacks, and Trojans. The second lesson is that we must continue to educate computer users in order to prevent them falling prey to the deceptions of computer abusers. A third lesson might well be that more effective prosecution of computer criminals is needed, along with stiffer sentences.

20.8 FURTHER READING

- Brown, B. C. *The Complete Guide to Email Marketing: How to Create Successful, Spam-Free Campaigns to Reach Your Target Audience and Increase Sales*. Ocala, FL: Atlantic Publishing Group, 2007.
- Goodman, D. *Spam Wars: Our Last Best Chance to Defeat Spammers, Scammers & Hackers*. New York, NY: Select Books, 2004.
- Haskins, R., and D. Nielsen. *Slamming Spam: A Guide for System Administrators*. New York, NY: Addison Wesley Professional, 2004.
- Jackobsson, M., and S. Myers, eds. *Phishing and Countermeasures: Understanding the Increasing Problem of Electronic Identity Theft*. Hoboken, NJ: John Wiley & Sons, 2006.
- James, L. *Phishing Exposed*. Rockland, MA: Syngress, 2005.
- Lininger, R., and R. D. Vines. *Phishing: Cutting the Identity Theft Line*. Indianapolis, IN: John Wiley & Sons, 2005.
- McWilliams, B. S. *Spam Kings: The Real Story behind the High-Rolling Hucksters Pushing Porn, Pills, and %*@# Enlargements*. Sebastopol, CA: O'Reilly, 2004.
- Schryen, G. *Anti-Spam Measures: Analysis and Design*. Berlin: Springer, 2007.
- Silver Lake, eds. *Scams & Swindles: Phishing, Spoofing, Spyware, Nigerian Prisoner, ID Theft*. Aberdeen, WA: Silver Lake Publishing, 2006.
- Spammer-X. *Inside the SPAM Cartel: Trade Secrets from the Dark Side*. Rockland, MA: Syngress, 2004.
- Wolfe, P., C. Scott, and M. Erwin. *Anti-Spam Tool Kit*. Emeryville, CA: McGraw-Hill Osborne Media, 2004.
- Zdziarski, J. *Ending Spam: Bayesian Content Filtering and the Art of Statistical Language Classification*. San Francisco, CA: No Starch Press, 2005.

20.9 NOTES

1. U.S. Internal Revenue Service, “Report Phishing,” IRS Website, May 8, 2013, www.irs.gov/uac/Report-Phishing
2. Symantec, “2013 Internet Security Threat Report, Volume 18,” *Symantec | Security Response Publications*, April 15, 2013, www.symantec.com/content/en/us/enterprise/other_resources/b-istr_main_report_v18_2012_21291018.en-us.pdf
3. Dean Turner, “Symantec Internet Security Threat Report, Volume X,” *Symantec | White Papers*, September 22, 2006, http://eval.symantec.com/mktginfo/enterprise/white_papers/ent-whitepaper_symantec_internet_security_threat_report_x_09_2006.en-us.pdf.⁴⁷ The diagram in Exhibit 20.1 and accompanying text originally appeared in *Trusted Email Open Standard: A Comprehensive Policy and Technology Proposal for Email Reform*, a white paper published in 2003 by a company called

NOTES 20 · 37

ePrivacy Group; contributors to the paper included Vincent Schiavone, David Brussin, James Koenig, and Ray Everett-Church.

4. The diagram in Exhibit 20.1 and accompanying text originally appeared in *Trusted Email Open Standard: A Comprehensive Policy and Technology Proposal for Email Reform*, a white paper published in 2003 by a company called ePrivacy Group; contributors to the paper included Vincent Schiavone, David Brussin, James Koenig, and Ray Everett-Church.
5. Karl Taro Greenfeld, “How Spam Meat Has Survived Spam E-Mail,” *Bloomberg Businessweek | Lifestyle*, May 17, 2012, www.businessweek.com/articles/2012-05-17/how-spam-meat-has-survived-spam-e-mail
6. Monty Python, “Spam,” *Youtube | Monty Python Channel*, uploaded January 13, 2009, www.youtube.com/watch?v=M_eYSuPKP3Y&list=UUGm3CO6LPcN-Y7HIuyE0Rew&index=32
7. See online article at www.templetons.com/brad/spamterm.html
8. Austin Modine, “Virginia De-convicts AOL Junk Mailer Jeremy Jaynes: Overturns Anti-Spam Law, Invokes Founding Fathers.” *The Register*. September 13, 2008, www.theregister.co.uk/2008/09/13/virginia_overturns_antispam_conviction
9. Grant Gross, “Court Overturns Virginia Spam Law, Conviction,” *PCWorld*, September 12, 2008, www.pcworld.com/article/151014/spam_law_overturned.html
10. Associated Press, “Sanford Wallace: Facebook Wins \$711 Million in Case against ‘Spam King,’ ” *Huffington Post | Tech*, October 30, 2009, www.huffingtonpost.com/2009/10/30/sanford-wallace-facebook_n_339703.html
11. FBI, “Sanford Wallace Indicted for Spaming Facebook Users: Self-proclaimed ‘Spam King’ Sent More Than 27 Million Spam Messages,” *FBI | San Francisco Division | US Attorney’s Office | Northern District of California*, August 04, 2011, www.fbi.gov/sanfrancisco/press-releases/2011/sanford-wallace-indicted-for-spamming-facebook-users
12. San Jose Federal Court, “Calendar for Judge Edward J. Davila,” *U.S. Courts*, May 23, 2013, www.cand.uscourts.gov/CEO/cfd.aspx?7143
13. The term “Email Service Provider” applies to companies that provide email addresses but not necessarily the Internet connection used to access those addresses. Google is a prime example with its Gmail; other examples include AOL, MSN, and Yahoo!, although all three of these have at times offered Internet connectivity as well as email and thus operated as ISPs as well as Email Service Providers.
14. Spam constituted approximately 90 percent of all email sent in 2007, according to numbers compiled by a variety of email service providers, including Postini, which is owned by Google.
15. See www.message-labs.com/intelligence.aspx for the list of monthly and annual reports; the 2007 report was available from www.message-labs.com/mlireport/MLI_2007_Annual_Security_Report.pdf
16. Robert McMillan, “90 percent of e-mail is spam, Symantec says,” *Computerworld | Applications*, May 26, 2009, www.computerworld.com/s/article/9133526/90_percent_of_e_mail_is_spam_Symantec_says
17. Amar Toor, “Symantec Says 92-Percent of All E-mail Is Spam, Phishing Attacks Declining,” *SWITCHED*, August 13, 2010, www.switched.com/2010/08/13/symantec-says-92-percent-of-all-e-mail-is-spam-phishing-attacks

20 · 38 SPAM, PHISHING, AND TROJANS: ATTACKS MEANT TO FOOL

18. Nicole Henderson, "Symantec Report Finds Spam Accounts for 73 Percent of June Email," *Web Host Industry Review*, June 28, 2011, www.thewhir.com/web-hosting-news/symantec-report-finds-spam-accounts-for-73-percent-of-june-email
19. Lance Whitney, "Spam Sinks to Lowest Level in Almost Three Years, Says Symantec: The Amount of Spam around the Globe Now Accounts for 70 Percent of All E-mail, a Sharp Decline from 2009 When It Accounted for 90 Percent," *c|net | News | Security & Privacy*, December 7, 2011, http://news.cnet.com/8301-1009_3-57338317-83/spam-sinks-to-lowest-level-in-almost-three-years-says-symantec
20. Kaspersky Lab, "Spam in 2012: Continued Decline Sees Spam Levels Hit 5-year Low," Kaspersky Lab, January 13, 2013, www.kaspersky.com/about/news/spam/2013/Spam_in_2012_Continued_Decline_Sees_Spam_Levels_Hit_5_year_Low
21. Darya Gudkova and Tatyana Shcherbakova, "Spam in January 2013," *Securelist | Analysis*, February 21, 2013, www.securelist.com/en/analysis/204792282/Spam_in_January_2013
22. Sara Radicati and Quoc Hoang, "Email Statistics Report, 2012–2016: Executive Summary," Radicati Group, April 23, 2012, www.radicati.com/wp/wp-content/uploads/2012/08/Email-Statistics-Report-2012-2016-Executive-Summary.pdf
23. Justin M. Rao and David H. Reiley, "The Economics of Spam," *Journal of Economic Perspectives* 26, no. 3 (2012): 87–100, www.aeaweb.org/articles.php?hs=1&fnd=s&doi=10.1257/jep.26.3.87
24. Located at www.networkworld.com/spam/index.jsp
25. Anti-Phishing Working Group, "About APWG," APWG Website, 2013, <http://apwg.org/about-APWG>
26. Janusz Lukasiak, "Blocking Spam Relaying and Junk Mail," Jisc Website, October 1999, www.jisc.ac.uk/publications/publications/blockingspamfinalreport.aspx
27. The exact terminology for block lists and black holes is a subject of ongoing debate involving possible issues of trademark infringement and political correctness.
28. McAfee, "Snowshoe Spaming Emerges as Threat to Email Security," *McAfee | Business Home | Security Awareness*, December 27, 2012, www.mcafee.com/us/security-awareness/articles/snowshoe-spamming-biggest-problem.aspx
29. David Brussin, Stephen Cobb, Ray Everett-Church, and Vincent J. Schiavone, "Network Resource Theft Prevention: Destroying the Economics of Spam," ePrivacy Group, November 2003, http://cobbassociates.com/library/spamsquelcher_wp2.pdf
30. See www.rfc-archive.org/getrfc.php?rfc=4408
31. See www.rfc-archive.org/getrfc.php?rfc=4406
32. Sentences for spam convictions have been rising steadily. Although many of those convicted have appealed, at least half a dozen have now served jail time.
33. Federal Trade Commission, "The CAN-SPAM Act: Requirements for Commercial Emailers," www.business.ftc.gov/documents/bus61-can-spam-act-compliance-guide-business
34. M. E. Kabay, "Can CAN-SPAM can spam?" *Networkworld*, February 02, 2004, www.networkworld.com/newsletters/sec/2004/0202sec1.html
35. Paul, "New Phishing Toolkit Uses Whitelisting To Keep Scams Alive," *VERA-CODE | the security ledger*, January 16, 2013, <https://securityledger.com/new-phishing-toolkit-uses-whitelisting-to-keep-scams-alive>

NOTES 20 · 39

36. Dave Lee, “The Comment Group: The Hackers Hunting for Clues about You,” *BBC News | Business*, February 11, 2013, www.bbc.co.uk/news/business-21371608
37. Jaikumar Vijayan, “DHS Warns of Spear-Phishing Campaign against Energy Companies: Attackers Used Information from Company Website to Craft Attacks,” *Computerworld | Security | Malware and Vulnerabilities*, April 5, 2013, www.computerworld.com/s/article/9238190/DHS.warns.of.spear.phishing.campaign.against.energy.companies
38. Anti-Phishing Working Group, “APWG Phishing Attack Trends Reports,” Anti-Phishing Working Group Website, April 24, 2013, <http://apwg.org/resources/apwg-reports>
39. Greg Aaron, “Phishing Activity Trends Report: 4th Quarter 2012,” Anti-Phishing Working Group Website, April 24, 2013, http://docs.apwg.org/reports/apwg-trends_report_Q4_2012.pdf
40. See www.ftc.gov and <http://onguardonline.gov/phishing.html>
41. Anthony Wing Kosner, “New Trojan Backdoor Malware Targets Mac OS X and Linux, Steals Passwords and Keystrokes,” *Forbes | Tech*, August 31, 2012, www.forbes.com/sites/anthonykosner/2012/08/31/new-trojan-backdoor-malware-targets-mac-os-x-and-linux-steals-passwords-and-keystrokes
42. PandaLabs, “PandaLabs Q1 Report: Trojans Account for 80% of Malware Infections, Set New Record,” *Panda Security | Press Room | News*, May 3, 2013, <http://press.pandasecurity.com/news/pandalabs-q1-report-trojans-account-for-80-of-malware-infections-set-new-record>
43. Sean Gallagher, “First Targeted Attack to Use Android Malware Discovered: Kaspersky Uncovers Trojan Spread by “Spear-phish” to Tibet Activists,” *ars technica*, May 26, 2013, <http://arstechnica.com/security/2013/03/first-targeted-attack-to-use-android-malware-discovered>
44. Dwight Silverman, “More than Half a Million Macs Infected with Flashback Trojan Malware,” *Chron | TechBlog*, April 5, 2012, <http://blog.chron.com/techblog/2012/04/more-than-half-a-million-macs-infected-with-flashback-trojan-malware>
45. Graham Cluley, “Sick Malware Authors Exploit Boston Marathon Bombing with Trojan Attack,” *Sophos | nakedsecurity*, April 17, 2013, <http://nakedsecurity.sophos.com/2013/04/17/malware-boston-marathon-bombing>
46. Graham Cluley, “Breakfast Malware at Tiffany’s? Trojan Horses Spammed out Widely,” *Sophos | nakedsecurity*, May 22, 2013, <http://nakedsecurity.sophos.com/2013/05/22/tiffany-malware>
47. Symantec, “Fraud Activity Trends,” *Symantec | Enterprise | Security Response | Internet Security Threat Report*, 2010, www.symantec.com/threatreport/topic.jsp?id=fraud_activity_trends&aid=underground_economy_servers

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 21

WEB-BASED VULNERABILITIES

**Anup K. Ghosh, Kurt Baumgarten,
Jennifer Hadley, and Steven Lovaas**

21.1 INTRODUCTION	21·1	21.6.2 Network Protocol Risks	21·10
21.2 BREAKING E-COMMERCE SYSTEMS	21·1	21.6.3 Business Application Logic	21·12
21.3 CASE STUDY OF BREAKING AN E-BUSINESS	21·2	21.6.4 CGI Script Vulnerabilities	21·14
21.4 WEB APPLICATION SYSTEM SECURITY	21·5	21.6.5 Application Subversion	21·15
21.5 PROTECTING WEB APPLICATIONS	21·6	21.6.6 Web Server Exploits	21·16
21.6 COMPONENTS AND VULNERABILITIES IN E-COMMERCE SYSTEMS	21·8	21.6.7 Database Security	21·19
21.6.1 Client-Side Risks	21·8	21.6.8 Platform Security	21·21
		21.7 SUMMARY	21·22
		21.8 FURTHER READING	21·23
		21.9 NOTES	21·23

21.1 INTRODUCTION. This chapter systematically reviews the primary software components that make up Web applications, with a primary focus on e-commerce, and provides an overview of the risks to each of these components.¹ The goal of this chapter is to point out that every system will have risks to its security and privacy that need to be systematically analyzed and ultimately addressed.

21.2 BREAKING E-COMMERCE SYSTEMS. To make a system more secure, it may be advisable to break it. Finding the vulnerabilities in a system is necessary in order to strengthen it, but breaking an e-commerce system requires a different mind-set from that of the programmers who developed it. Instead of thinking about developing within a specification, a criminal or hacker looks outside the specifications.

Hackers believe that rules exist only to be broken, and they always use a system in unexpected ways. In doing so, they usually follow the path of least resistance. Those areas perceived to provide the strongest security, or the most resistance to hacking, will likely be ignored. For example, if a system uses Secure Sockets Layer (SSL) to encrypt Web sessions between Web clients and the Web server, a hacker will not try to

21 · 2 WEB-BASED VULNERABILITIES

break the encryption stream but instead will look for an easier way to get at the data after they are decrypted and stored in the clear.

Hackers go where the money is—sometimes literally, sometimes not. They typically try to hack into a site only if there is some reward for their effort. Sometimes hackers are motivated by money, but as often the motivation is the lure of fame, notoriety, or acceptance by a peer group. The level of protection should be commensurate with the value of the resource being protected. For instance, a Website that publishes the menus for local restaurants may not be seen as a target for a denial of service or any other type attack. Such a Website simply is not as attractive a target as a bank's online Website, where a hacker can certainly gain in notoriety and even profit financially. Similarly, most people do not bother encrypting email messages due to lack of security awareness, and because most potential snoopers are not interested in ordinary personal email. Sensitive email from a high-profile organization should be encrypted, however, in order to protect valuable intellectual capital.

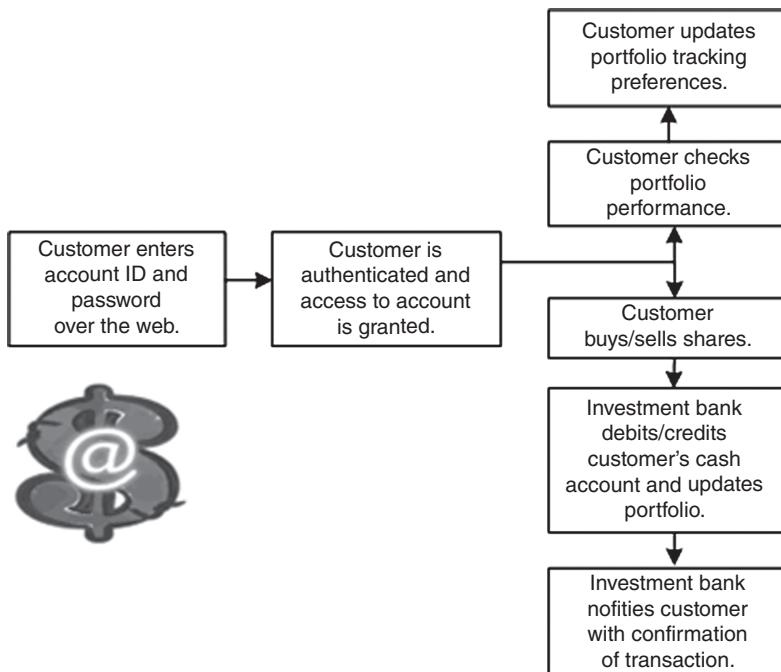
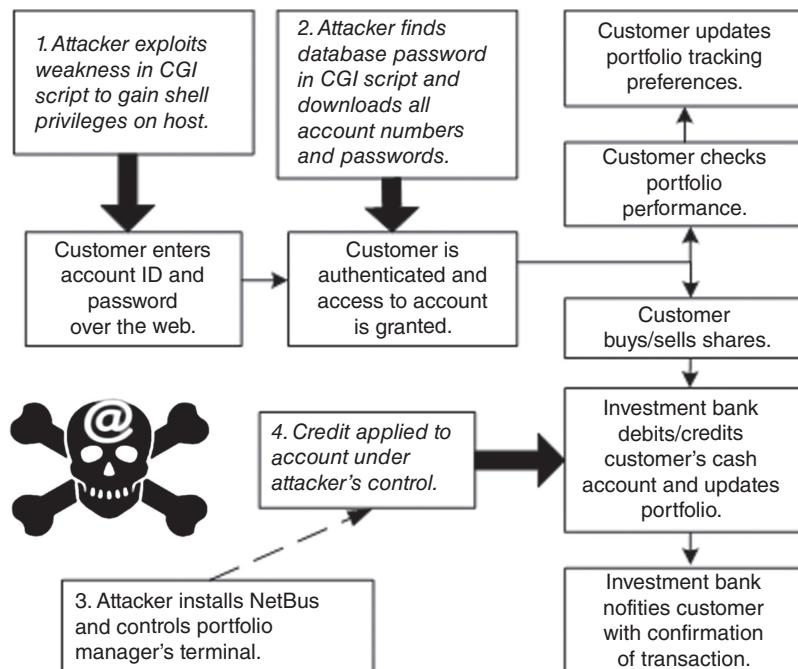
It is important to remember that the e-commerce system is a chain that is only as strong as its weakest link.² Hackers naturally attempt to attack that point of least resistance. This explains why a site may deliberately set up a honeypot (a sacrificial system) with appealingly vulnerable services in order to track and monitor potential hackers. In e-commerce, or any systems with a Web presence, cryptography often is perceived to provide the strongest level of protection; thus, hackers generally attack host-side services and client-side content.

Maintaining strong host security, both inside and outside the perimeters, is critically important. One unfortunate side effect of corporate perimeters is that system administrators tend to overlook host-based security, which may leave the host completely vulnerable. The result is that once hackers make it through or around perimeter devices and policies (e.g., firewalls, routers, or even receptionists), they can leverage the internal trust relationships to compromise many resources, including workstations, servers, and phone systems. The prudent administrator will exercise equal concern both at the entry to systems and within those systems.

21.3 CASE STUDY OF BREAKING AN E-BUSINESS. Consider an online investing e-business application and how a hacker might go about disassembling its security for malicious or financial gain. Online investing is very popular for several reasons. Rather than waiting for quarterly statements in the mail or dealing with a phone menu, customers can quickly view the status, balances, and current value of investment holdings by visiting the Web pages of their portfolio managers. If they wish to buy and sell equity shares on demand, they can establish online Web-enabled brokerage accounts. Exhibit 21.1 shows a simplified workflow diagram of an online investing application that enables users with established accounts to view portfolio holdings and cash balances, to update the stocks tracked, and to conduct online trades.

To see how this application can be broken, it is helpful to look at a sample network architecture that implements the online application. Exhibit 21.2 shows the network architecture of the system that implements the online investing application, along with example exploits. The system consists of the end users' client machines, the Internet, routers, firewall, front-end Web and email servers, application servers, databases, and workstations.

There are many ways a hacker could break this online application. Exhibit 21.2 shows one possible scenario. In step 1 of the attack, the hacker uses the Internet and a Web browser to misuse one of the Common Gateway Interface (CGI) scripts that invoke the application on a server. The CGI script could be a development CGI script

CASE STUDY OF BREAKING AN E-BUSINESS 21 · 3**EXHIBIT 21.1** Online Investing Application**EXHIBIT 21.2** Breaking an E-Business

21 · 4 WEB-BASED VULNERABILITIES

inadvertently left on the server before going into production, a default CGI script from an application server distribution, or a script that implements flawed logic in the online investment application. Exploiting CGI scripts via cross-site scripting or other common exploit methods allow hackers to gain shell access to Web servers. CGI script vulnerabilities are discussed later in this chapter.

The vulnerability need not be in a CGI script. Application servers can be implemented in Java, C, C++, Perl, or Python in various application server frameworks. The difficulty lies not in which language the business application logic is developed; more important are the vulnerabilities introduced by the complex logic at this middleware layer. One of the key problems in the development of application middleware is poor input sanity checking; that is, the developers fail to impose limits on acceptable input. The hacker can exploit the lack of input sanity checking to feed the application server unexpected input used in system commands. This technique can gain shell privileges on the machine.

Although application server misuse is a common way of breaking into systems, there are many other ways to gain the initial access in step 1 of the attack. For instance, the Web and mail servers may be running any of several network services, such as *FTP* and *BIND/DNS*, which may be misconfigured and thus “poisoned.” The Web and mail server software themselves may be vulnerable to attack. Many popular commercial Web and mail servers have been vulnerable to buffer overflow attacks that often permit full system root privileges on the host.³ Once attackers gain system privileges on an internal host, they can exploit the inherent internal trust often woven among machines through network policies in order to gain access to other systems on the network. This strategy is precisely what the attacker follows in step 2 of the attack illustrated in Exhibit 21.2

Once attackers have access to the various file systems on the application server, they can view source code of CGI scripts or other application middleware to discover customer account numbers, passwords, and even database administrator passwords for accessing the back-end databases. From there, they can download important and confidential client information stored in the database. In step 3 of the attack, the attacker leverages the internal privileges gained to plant backdoors into the system unnoticed. A suite of software, commonly known as a rootkit and available to hackers, allows them not only to get into a system unnoticed but also to erase their tracks in audit logs. In the example shown, the hacker installs a rogue remote administration program known as Back Orifice, which provides the ability to remotely administer the host and network with the same privileges and power as an authentic system administrator.

At this point in the attack, the hacker has assumed total control of the systems and any Web-facing or commerce-related processes, with many options including:

- Stealing customer or company information for the purpose of financial gain
- Defacing the Web pages for notoriety or to publicize an agenda
- Working in a stealthy manner to uncover proprietary business information and other confidential intellectual capital (espionage)
- Blackmailing the business with threats of discrediting it
- Subverting the application for any other personal gain

Step 4 of the attack illustrates the last case, where the attacker credits a personal cash account. The hacker must move quickly enough to withdraw these funds before traditional back-end auditing discovers the discrepancy. There have been many defaced

WEB APPLICATION SYSTEM SECURITY 21 · 5

Web pages of government agencies, and other important sites, and many reported instances of the other cases. Of course, some companies are understandably reluctant to publicize events that might lessen customer confidence, yet in many cases, legislation now requires the disclosure of a breach.

Unfortunately, it takes only a single flawed, unpatched computer, or overlooked vulnerability, for a hacker to compromise a system as a whole. Although *defense-in-depth* (using multiple forms of security, such as firewalls on the perimeter and intrusion detection inside the network) is a popular strategy, often multiple layers of defense fall like a house of cards when a single hole is exploited. For example, an attacker who gains *root* capability can disable all other security measures. The problem is known as an *asymmetric* attack because it is much more difficult and costly to defend against such an attack than to launch one. Although a committed hacker may be capable of spending as much time as is needed to break into a system, a company cannot spend as much money as needed to defend itself against all possibilities.

The number of flaws that can be exploited is staggering, considering all the different platforms and devices that make up current information technology (IT) infrastructures. Compounding the problem is the fact that a hacker can work in relative anonymity using a \$500 computer and modem to launch attacks. Even worse, hackers can work from any number of Internet kiosks available in airports, malls, cafés, and even laundries. As hackers get more sophisticated, and as more easily utilized scripts become available, attacks will be launched from mobile devices that can roam in and out of different geographic zones and then be discarded—making tracking of the attacker next to impossible.

21.4 WEB APPLICATION SYSTEM SECURITY. In spite of the fairly bleak picture painted here, organizations can effectively manage their risk from hackers. As in many other security domains, the security posture or stance assumed by the organization is critical for deterring and thwarting hackers. To use a physical-world analogy, consider burglars who intend to break into homes in a nice neighborhood. As the burglars scope out potential targets, they will notice some houses with burglar alarms—complete with conspicuous signs of the alarm systems—and some without. In all likelihood, the burglars will bypass the houses with the burglar alarms and move on to the other, less well-protected targets. Thus, the security stance assumed by the owner plays an important role.

Every organization must first determine its desired security stance as documented in its security policy. System administrators use the security policy to configure the systems, routers, firewalls, and remote access solutions. Without an explicit security policy, there is no way to determine what the security stance of the organization is, how to configure its systems, or even if a nonobvious security breach has occurred. Once the security policy is developed, the actual security implementation must be assessed. That is, the system must be tested and evaluated to determine how well it meets its security policy. Usually there is a difference between the desired stance and the actual stance. This difference is the security gap between where the organization would like to be (secure against particular threats) and where it actually is in practice.

The process of developing a security policy and evaluating the organization's systems against that policy will identify not only the gaps between the actual security stance and the desired posture but also weaknesses in the security policy itself. It is important to have an independent party, preferably an outside party, evaluate the security stance of the organization. A third party can fairly assess whether the organization's system upholds the security policy. If the group that develops the security policy or the system

21 · 6 WEB-BASED VULNERABILITIES

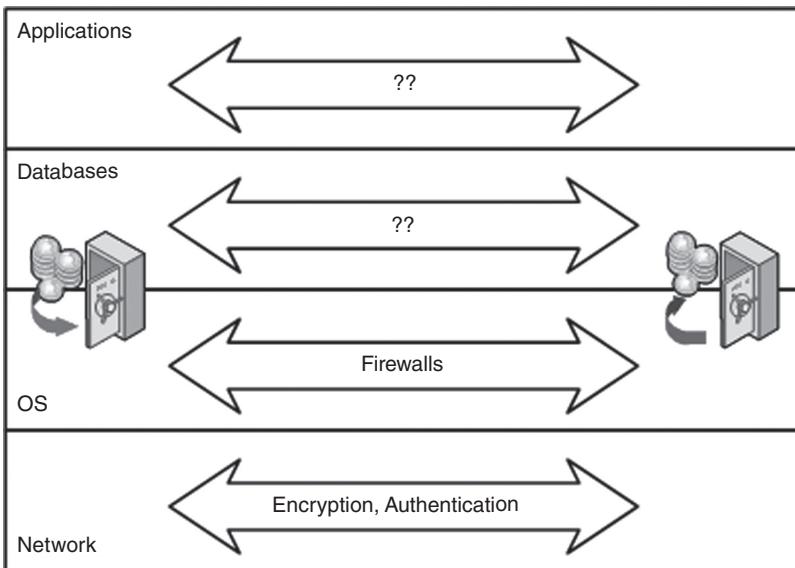


Exhibit 21.3 Layered View of an E-Business Application

configuration is also responsible for evaluating it, the evaluation may be biased, and potential vulnerabilities may be overlooked. This may occur not through ill will or dishonesty but rather because of the difficulty of being objective about one's own work.

21.5 PROTECTING WEB APPLICATIONS. Web application security can be understood from different views. First, consider the view of e-businesses in Exhibit 21.3. The diagram shows two e-businesses communicating over the Internet, perhaps performing business-to-business types of transactions.

In this view, the lowest layer of the e-business is the *networking layer*. At the networking layer, there are concerns about the reliability, integrity, and confidentiality of the data that runs over the communications channel. This layer is of particular concern because the Internet is a public switched network, meaning that any number of third parties may have access to the data that traverses the nodes of the Internet on the way from the data source to its destination. Also, the Internet Protocol (IP) is a connectionless protocol, which means there is no dedicated circuit between source and destination. As a result, packets sent during a given session may take different routes to their destination, depending on traffic congestion and routing algorithms. Because IP is an unreliable datagram protocol (i.e., IP uses independent packets forwarded from node to node, but there is no guarantee of successful transmission), the networking layer includes a connection-oriented reliable transmission layer such as Transmission Control Protocol (TCP) that ensures that dropped or lost packets are retransmitted and that bit flips that may have occurred during transmission (e.g., over wireless networks) are corrected.

Although TCP/IP provides for more reliable delivery of Internet packets, IP version 4 does not provide *secure* connection services. Typically, this means that there is no *guarantee* of confidentiality, identification, or even delivery of packets sent from one Internet host to another. Because packets often traverse several Internet nodes from source to destination, packet contents can be intercepted by third parties, copied, sub-

PROTECTING WEB APPLICATIONS 21 · 7

stituted, or even destroyed. This is the risk that most people citing risks of e-commerce have decried; they have overlooked the more substantive risks of e-commerce dealing with server- and client-side security and privacy. IPv6 holds promise to bring many security improvements to the network layer, but it has not been deployed widely in the United States. Fortunately, even in IPv4, we have good solutions to the data confidentiality problem. Cryptographic techniques can provide strong guarantees for data confidentiality, authentication of parties, and integrity of data sent during the transmission. Furthermore, digital signatures can be used to sign received mail in a “return receipt” application that provides guarantees of delivery of email. Thus, as shown in Exhibit 21.3, we can use encryption services to protect data transmitted over the network.

The operating system (OS), or *platform*, that hosts the Web applications lives on the networking layer. In a layered model, the services of one layer use the services of the lower layer and provide services to upper layers. The network layer often is thought of as a core portion of the operating system; however, from a layered services point of view, the OS software runs on top of the network layer.

Operating systems are notoriously rife with software flaws that affect system security. Operating systems are vulnerable because commercial OSs today are immensely complex; for instance, the Windows Vista operating system is purported to have more than 50 million lines of source code. It is impossible to catch all software design and programming errors that may have security consequences in a platform this complex. Even though UNIX operating systems have been in use for the better part of 30 years, new flaws in OS utilities are found on a weekly basis across all the different UNIX platform variants.⁴

Security holes in the platform are critical by nature. That is, if the OS itself is vulnerable to exploitation, security provided by the application can be compromised by holes in the platform. The OS is always the foundation on which applications are built, so cracks in the foundation make for weak security at the application layer. As Exhibit 21.3 suggests, firewalls provide protection against some operating system flaws. One of the key roles of firewalls is their ability to shut down services offered to *logical domain addresses*.

Using Internet domain addresses, the firewall administrator can partition Internet addresses into *trusted* and *untrusted* domain ranges. For instance, any Internet address outside the company’s domain can be considered untrusted. As a result, all OS services, such as remote logins, can be shut down to everyone outside of the company’s domain. Even within a company, the domains can be partitioned so that certain *subdomains* are trusted for access to certain machines, but others are not. The key benefit then of firewalls is their ability to restrict access to the platform through *offered services* (i.e., specific functions that pass data through the firewall). As a result, firewalls can make it easy to hide OS flaws from untrusted entities.

Even so, firewalls are vulnerable to data- or code-driven attacks through offered services. These are sometimes called “allowed path” vulnerabilities. For instance, an attack through SMTP (mail) or HTTP (Web) will not be stopped by a firewall if the firewall is configured to let email and Web services through, as is necessary for e-commerce. Firewalls also will not stop OS exploits from insiders or from the trusted entities that are granted access to the platform. However, it is important to realize that firewalls, if properly configured, can close down exposure to a significant number of platform vulnerabilities simply by denying untrusted outsiders access to platform utilities and services. Exhibit 21.3 illustrates reasonable protection from network- and platform-based attacks but not from application and database attacks. The database

21 · 8 WEB-BASED VULNERABILITIES

layer is shown separately in the diagram because of the importance of its role in e-commerce; however, database attacks usually can be considered as a type of application attack.

Application-based attacks represent a critical issue that is gaining more attention in recent years. There is no simple solution, but there is a growing sense that the security provided by a firewall, an SSL-enabled Website, and encryption, together with digital certificates and signatures, is not enough. These measures are necessary but not sufficient. Applications, above all, *are* the online business. The Payment Card Industry Data Security Standards (PCI DSS) puts any organization processing credit card transactions on notice that Web application vulnerability protection was required as of June 2008. The PCI DSS specifies either “having all custom application code reviewed for common vulnerabilities by an organization that specializes in application security” or “installing an application layer firewall in front of Web-facing applications.”⁵

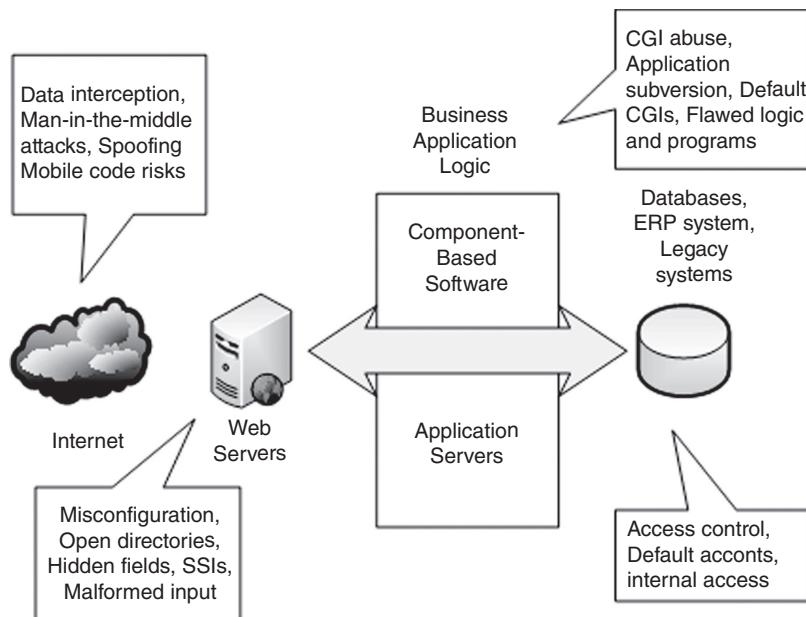
Online applications are increasingly sophisticated, and the software that implements the application logic has become highly complex, requiring component-based and object-oriented paradigms such as Enterprise Java Beans, CORBA, and DCOM/COM services. Collectively, these are known as *application servers*. The key point, however, is that because the application logic is custom and complex, it is often rife with errors in implementation or logic that can be and often are exploited by hackers.

Application security must not be confused with marketing claims. A secure online application is one that is resistant to attacks. It is not simply one that authenticates end users, encrypts transaction data, provides nonrepudiation of transactions, and guarantees service. These are all matters of importance that address characteristics of the transaction, not properties of the software. The remainder of this chapter addresses the software problem in some detail. The next section provides a different view of e-commerce systems from the layered view discussed previously. It identifies vulnerabilities in the different software components and strategies for managing the risks.

21.6 COMPONENTS AND VULNERABILITIES IN E-COMMERCE SYSTEMS.

Exhibit 21.4 shows a generic multilayer architecture of an e-business, together with a summary of the types of vulnerabilities and risks to each of the major components. Using the Internet, Web clients (running on PCs or handheld devices) interface with a front-end Web server, a middleware layer of business application logic, back-end databases, an ERP system, supply-chain management software, and even some legacy systems that are now brought to the Internet.

21.6.1 Client-Side Risks. Most e-commerce is performed using standard Web browsers and mail clients. Increasingly, e-commerce is being performed on handheld mobile devices such as personal digital assistants (PDAs) and mobile phones. The security risks particular to wireless devices are covered in Chapter 33 of this *Handbook*. Client-side security risks are mainly from malicious mobile code such as Web scripts, ActiveX controls, and hostile Java applets.⁶ Another major risk in client-side software is loss of privacy.⁷ Each computer, with its related software, receives and transmits a great deal of personal identifying information (PII). For instance, browsers may convey information about the computers (name, IP address, browser type, version, company name) and sometimes about the users themselves, particularly if automatic form-filling features have been enabled. Browsers also are used to track movements through the Web. For instance, every Website the browser visits typically gets a record of the previous site from which the user entered. Banner ads in Web pages also track which sites have been visited in order to create a profile of Web usage, and cookies can be

COMPONENTS AND VULNERABILITIES IN E-COMMERCE SYSTEMS 21 · 9**EXHIBIT 21.4** Multitier Architecture of an E-Business

placed on end-user systems for the purpose of tracking and at times re-creating user input.

A class of more insidious programs, known as *spyware*, can send information about computer usage out to specific sites, often without the user's knowledge or approval. One of the key risks with client-side software is that users simply do not know what the programs are revealing about themselves, or to whom. A simple principle is that a program should not be sending out any data that the user does not know about. An audit of the programs on any given machine probably would reveal that quite a few are in violation of that principle. Although most spyware programs are written to provide marketing data for software vendors or to redirect a user to predetermined sites, many manipulate and target users more effectively by spying on usage activity and even installing additional malicious code. For example, SpectorSoft markets a spyware tool known as eBlaster to people who suspect their spouses are engaging in illicit online affairs. This type of spyware program spies on user activity, from keystrokes to screen shots, and sends a log of the activity out over the network to a predetermined logging site. Spyware can be identified using diagnostic programs such as Ad-Aware from Lavasoft (www.lavasoft.com), but because of the dynamic nature of spyware technology, it may be difficult to identify all possibilities of infection.

A final client-side risk that businesses need to be especially concerned about is the risk of malicious executables that run on their user workstations. The desktop machine is like a petri dish for software: It is constantly changing and growing with new software executables—some of an unsavory nature. Malicious software, or *malware* as it is now known, finds many ways of infecting machines. For instance, one common way of disseminating malicious software is via email attachments. Another is by masquerading as legitimate software on a Web page available for download. Users often upload and download software to and from internal network file shares. In addition, USB “thumb drives” and even old-fashioned floppy disks are still a viable way of transmitting

21 · 10 WEB-BASED VULNERABILITIES

malicious software. The Back Orifice “remote administration kit” is a well-known example of malicious software that, when installed, will allow anyone—including a hacker—to administer and control a machine remotely. Some malware products, running on an internal machine, may compromise the entire network because a remote hacker can control a trusted resource, possibly with enterprise administration rights. Data leakage and theft offered by the initial spyware infection can result in incidents that might undermine the integrity, authenticity, availability, or auditability of an entire data infrastructure. It is essential, therefore, that corporations carefully filter and closely monitor the application software that is downloaded and run (with or without the user’s knowledge) on all systems from the e-commerce perimeter to the individual internal hosts.

21.6.2 Network Protocol Risks. Network risks primarily arise from sending confidential data over the Internet—a public, packet-switching network. Many good protocols address the risks of sending confidential data over the Internet.⁸ In fact, a few years ago, the list included:

- SET
- SSL
- S/HTTP
- S/MIME
- CyberCash

Although some of these protocols are still around in one form or another, the industry has generally accepted SSL as the protocol of choice for secure Web browser transactions. The objective of most secure network protocols is to layer security properties on top of the TCP/IP network layers. Although TCP/IP provides reliable and robust delivery of datagrams over the Internet, it does not provide confidentiality, authentication, or strong message-integrity services. These are the properties that secure protocols provide. Some go even further. For instance, SET-compliant protocols leave the credit card number encrypted even at the merchant site. Since the merchant does not need to know the consumer’s credit card number, by hiding the number from the merchant, a significant portion of credit card fraud can be eliminated. Rather than decrypting the credit card number at the merchant site, it is passed in encrypted form from the merchant to the credit-issuing bank. There it is decrypted, and the merchant’s account is credited the amount of the purchase. The protocol details of SET, SSL, and other e-commerce protocols are described in Chapter 3 of this *Handbook*, in *E-Commerce Security: Weak Links, Best Defenses*, and in other books as well.⁹

Depending on the needs of an online business application, there is a requirement for more or less of the security properties afforded by secure protocols. For most Web browsing, a secure protocol is unnecessary; the standard Web protocol, HTTP, suffices. However, when customers send confidential or personal information to a site, a secure protocol that encrypts the data is preferable. The de facto secure protocol standard is SSL, now implemented in every standard Web browser. SSL will not only negotiate a secret session key between the Website and a client to encrypt the data, but it also will authenticate the Website. The Website must have a valid certificate endorsed by a Certificate Authority, which the users implicitly trust. Using a list of trusted Certificate Authorities maintained within the client browser, access can be prevented to other, untrusted sites. Once the connection is established, the user can verify that the Website

COMPONENTS AND VULNERABILITIES IN E-COMMERCE SYSTEMS 21 · 11

is, in fact, the one intended by examining the site certificate. Users rarely do this in practice, but the certificate is available and should be more widely utilized.

These secure properties—that is, encrypted sessions and host site authentication—serve the purpose for most online commerce applications. Some applications, though, may demand even stronger security services. For instance, banking and investment applications often transmit highly confidential information, with the possibility of huge financial losses through inadvertent error or intentional fraud. These types of transactions require not only confidentiality of the data, but also authentication of the client. A financial institution must never give access to account information or permit transactions without authenticating the user.

Common identification schemes on the Internet include simple user name and password authentication. A much more secure solution is to require strong client authentication using client certificates. SSL supports client certificates, although sites rarely use this capability because it involves requiring customers to obtain a certificate from a Certification Authority.

In the future, e-commerce protocols will need to be increasingly sophisticated if they are to meet the stringent security and privacy requirements of new Internet applications. For example, as criminal, medical, patent, and other important databases migrate to the Internet, the protocols developed for accessing them need to consider the security and privacy needs of the database owners and maintainers, as well as those of the client or requester. Today, the progress in genomics is producing much information about the likelihood of developing deadly disease based on a genetic DNA sequence. This knowledge raises moral and ethical questions as to how much information about potential diseases should be revealed to doctors and patients, and also raises the specter of such information getting into the wrong hands once it is accessible on the Internet.

Consider the case of a doctor querying an online genetic disease database with a patient's DNA sequence. The online application attempts to match the DNA sequence with diseases that might develop in the future. If the database were maintained by a commercial entity such as an insurance provider, the patient almost certainly would not want the company to know of any disease that might be returned as the result of a query, because that information could be used to deny both insurance and employment.

Likewise, the database maintainer probably would not want to know of the query or of its result, as such knowledge might put it at risk for lawsuits, should the information be leaked. Furthermore, the database maintainer would want the rest of the database to remain inaccessible except for specific results returned to an approved inquiry. Preventing access to any other information in the database would help to protect the commercial interests of the company, because then the database could not be duplicated easily. Nor could queries be made without a cost-tracking mechanism. To support this dual model of secure and private information access, e-commerce protocols need to be developed and commercialized that not only encrypt data in transmission but also consider the total security and privacy needs of both parties to the transaction.

Another e-commerce application area that will require better security and privacy protocols involves applications that accept e-cash or digital coin payments. Currently, most online payment schemes use either credit or debit cards, with payments made from the buyer's checking account at a bank. Payments are made either with online verification of funds or with off-line batch payments at the end of the day. A number of applications, particularly those involving payments of a few dollars or even pennies, are being created that cannot support the costs of a bank-based transaction.

Many commercial services and products, such as vending machines and parking meters, are coin activated and do not require customers to have an account or a line

21 · 12 WEB-BASED VULNERABILITIES

of credit. Although efforts are being made to convert such services to computerized devices, with micropayment capabilities, it will be many years before this is actually accomplished. In any case, there will always be customers who want to pay with cash or its electronic equivalent, so that the transaction cannot be tracked by a third party, such as a bank or a creditor.

Newer online applications for micropayments may include collecting fees for downloading music, data, weather reports, stock quotations, articles from magazines, and pages from books. Many of these applications are provided today without charge and supported by banner advertising, but a concern for profits is motivating most Websites to seek additional sources of income. Whatever the application, there is a need for cash-based alternatives to the current account-based system for making payments. The key security and privacy concerns are with ensuring that e-cash is properly spent and accounted for and that anonymity is preserved throughout. Several protocols have been developed with these goals in mind, but none has reached commercial success or adoption by the vendor community. As mobile e-commerce begins to drive more traditionally cash-based transactions (e.g., parking meters, vending machines, and ticket booths), wireless vendors may adopt these new digital cash-based protocols.

Regardless of the network protocol used in an e-commerce application, the key concern is for those attackers who will attempt to breach the easiest obstacle in their quest to obtain system privileges and unauthorized access to data. If the security provided by the network protocol is perceived to be strong, attackers will look for alternatives that bypass the network security. For instance, these types of standard attacks from the hacker's toolkit will bypass the security provided by most e-commerce protocols:

- **Man-in-the-middle attacks.** Capturing transmissions in transit for eavesdropping or forgery
- **DNS attacks.** Altering records in the worldwide Domain Name System to misdirect connections to the wrong addresses
- **War dialing.** Automated testing of every phone number in a block of numbers
- **Exploiting software vulnerabilities in network services.** Such as FTP, Bind, SMTP, and HTTP servers
- **Internal access.** Improper use of authorized access by an insider (employee, contractor)
- **Leveraging trusted hosts.** Attacking from another system that has a privileged relationship with the target system
- **Brute-force crypto attacks.** Automated testing of all possible decryption keys to decipher a ciphertext

In summary, it is important not only to select the appropriate network protocol for each online application but also to avoid a false sense of security that might arise from use of a "secure" network protocol. Good security engineering will consider vulnerabilities in other components of the system that are more attractive targets to determined hackers.

21.6.3 Business Application Logic. The business application logic pictured in Exhibit 21.4 represents one of the key areas of vulnerability in e-commerce systems. The program logic encodes what the online business is offering in terms of products,

COMPONENTS AND VULNERABILITIES IN E-COMMERCE SYSTEMS 21 · 13

services, and customer convenience. It also defines the look and feel of the Website and provides all of the interactive features, such as dynamic Web pages, personalized Web pages, and online transaction capabilities. Because each application is unique, the software that implements the logic must be, in many cases, custom-developed for each particular site.

In contrast, most of the other software components of a Website are commercial off-the-shelf (COTS) software. For instance, the Web server, back-end databases, and supply-chain logistics software are often purchased off the shelf from software vendors. With COTS software, the end user has no control over the code and therefore is not responsible for coding bug fixes. When software bugs in COTS software are discovered, the software vendor usually issues a patch or incorporates a bug fix in the next release version. Software vendors can fix discovered bugs, but they depend on customer sites to actually apply the patches or upgrades to the software. In practice, this occurs less often than desired and is a significant reason why many Internet-based systems are vulnerable.¹⁰ The most important task in securing COTS software systems is to make sure that: (1) they are properly configured for a secure installation according to the site's security policy, (2) the software is properly updated to the current version and patch level, and (3) administrators are aware of the risks associated with bugs in the software.

Because many business applications are custom-developed, either by an in-house staff or by outsourcing to an e-business developer, the programs represent a key risk for several reasons. The dynamic, interactive nature of e-businesses, coupled with increasingly sophisticated online services, requires a significant amount of development to code the logic. As a result, the application programs tend to be very complex pieces of software, likely to contain flaws, and susceptible to the kinds of attacks launched against Websites. In practice, errors in design and implementation of business application logic often compromise the security of an e-business.

Traditionally, the middle tier of software is implemented on Web servers using CGI and more recently PHP: Hypertext Processor (PHP). CGI and PHP scripts are programs that run on the Web server machine as separate processes from the Web server software. The Web server invokes these general-purpose programs in response to user requests. The CGI/PHP script's main function is to process user input and to perform some service, such as retrieving data or dynamically creating a Web page for the end user. Because CGI and PHP scripts process untrusted user input, the security risks associated with them and other forms of middle-tier software are extremely high. Many attacks against Web-based systems are implemented by exploited CGI/PHP scripts. And while CGI and PHP scripts can be written in any general-purpose programming language, they are written most often in Perl, C, Tcl, and Python.

More recently, component-based software (CBS) is making inroads in e-commerce applications as well as in standard Web applications. The purpose of CBS is to develop, purchase, and reuse proven software in order to implement application logic quickly, easily, and with high quality. Two of the more popular component frameworks for e-commerce applications are Enterprise JavaBeans (EJB) using XML and Java 2 Enterprise Edition (J2EE), which supports component-based Java. Other component models include the Object Management Group's (OMG) Common Object Request Broker Architecture (CORBA) and Microsoft's Common Object Model (COM) and Distributed COM (DCOM). These component frameworks are the glue that enables software components to use standard infrastructure services while hiding the details of the implementation by using well-defined interfaces.

Business application logic, when coded in CBS systems, usually runs on application servers with particular component models, such as EJB, CORBA, COM, and DCOM.

21 · 14 WEB-BASED VULNERABILITIES

CBS also provides an interface to back-end services such as database management, enterprise resource planning (ERP), and legacy software systems.

In addition to supporting traditional CGI functions, component-based software is expected to enable distributed, business-to-business applications over the Internet. The component-based software paradigm also supports good software engineering, as described later. The Unified Modeling Language (UML) facilitates object-oriented analysis and design for component-based frameworks. In addition, as the market for component-based software expands, many standard business application components will be available for purchase off the shelf.

Although the benefits of component-based software are numerous, they pose security hazards similar to those of CGI scripts. Component-based software enables development in general-purpose programming languages such as Java, C, and C++, which can execute with all the rights and privileges of server processes. Like CGI, they process untrusted user input, and because component-based software can be used to build sophisticated, large-scale applications, the likelihood for errors may be even greater than for simple CGI scripts. Regardless of the implementation—CGI or application servers—the security risks of server-side software are great, and therefore server-side software must be designed and implemented carefully.

The key risks in the middleware layer of e-commerce sites are:

- Misconfiguration of the CGI
- Default and development CGI scripts being left on the production server
- CGI misuse
- Application subversion
- Flawed logic
- Programming errors

21.6.4 CGI Script Vulnerabilities. CGI scripts are frequent targets of attackers because they are often misconfigured and vulnerable to misuse.¹¹ When designing CGI scripts, it is prudent to expect the unexpected, particularly the malicious attack. Although the Web designer has control over the content of CGI scripts, there is no control over what end users are going to send to them. Also, often overlooked are vulnerabilities of CGI scripts that exist on the server as part of the distribution but that are not actually used in the application. Some CGI scripts, included as part of the Web server distribution, have well-known flaws that can be exploited to obtain unauthorized access to the server. Even if the default CGI scripts are not used as part of the Web server pages, anyone can access them by simply knowing the script names.

One of the most common—yet easily preventable—security hazards is misconfiguration of software, especially CGI scripts. One feature supported by many Web servers is the ability of individuals throughout an organization to write CGI scripts and have them execute from their own directories. Although useful for prettying up personal Web pages, this feature also can introduce system security hazards. In Web-based applications, the Web server should be configured to prevent CGI scripts from executing anywhere but in a single CGI directory under control of the system administrator.

The script-aliased CGI mode for Web servers ensures that CGI scripts will execute only from an explicitly named directory in the server configuration file. In addition, the CGI script path is not named in the Uniform Resource Locator (URL) to the CGI. Rather, the server “aliases” the explicit path to the CGI script to a chosen name, such

COMPONENTS AND VULNERABILITIES IN E-COMMERCE SYSTEMS 21 · 15

as *cgi-bin*. Thus, running the server in script-aliased CGI mode prevents rogue CGI scripts from executing while it also hides the explicit path to the CGI scripts.

The CGI script directories also should be properly configured using OS access controls. For instance, if CGI scripts are written in a compiled language such as C, the script sources should be excluded from the document root of the Web server, so that they cannot be accessed via the Web. They should be accessible to the system administrator or Web content development group only, and inaccessible to everyone else in the organization. If the script sources fall into the hands of malicious perpetrators, the source code can be inspected for flaws, making the perpetrator's job even easier. Access to the CGI executables directory, frequently called the *cgi-bin*, should be properly controlled as well. Only the Web server and administrator need access to this directory. Liberal access permissions to the CGI executables directory give those with malicious intent the opportunity to place their own scripts within the site.

Most CGI scripts are written in scripting languages such as Perl, JavaScript, and Python. Scripting languages are useful for rapidly prototyping systems, but they also let the developer write dangerous code very easily. For instance, it is easy to construct system commands with user input, a potentially dangerous situation. Writing the same system functionality requires several lines of programming code and knowledge of system libraries. The easy accessibility of scripting languages makes them appealing, but also threatening to security-critical applications. It is also important to prohibit access to interpreters from the Web server. For instance, system administrators may be tempted to include the Perl interpreter in CGI script directories; however, doing so provides direct Web access to interactively execute Perl commands—an extremely dangerous configuration.

Finally, administrators should account for every CGI program on the server in terms of its purpose, origin, and modifications. Remove CGI scripts that do not serve a business function, and view with suspicion CGI scripts that are distributed with operating systems and Web servers, downloaded from the Internet, or purchased commercially. These steps will eliminate most of the potentially dangerous CGI scripts. Once a stable set of CGI programs is established, make a digital hash of the program executables (e.g., using MD5 or SHA-1) to enable future integrity checks.

21.6.5 Application Subversion. Application subversion attacks are not discussed often in relation to e-businesses, but they represent a significant threat to most online applications. Application subversion is a form of program misuse. Unlike buffer overflow attacks, application subversion attacks exploit the program logic without violating program integrity, in order to elevate user privileges and gain unauthorized access to data. It is the very complexity of the target program that gives the attacker the means to gain unauthorized access. Application subversion attacks use programs in ways that the program's designers and developers did not anticipate. Typically, these attacks are not scripted, but rather developed from online interactive use and subsequent abuse.

Referring to Exhibit 21.1, an application subversion attack will attempt to discover ways of short-circuiting paths in the workflow. For instance, there may be a hidden path that lets the user gain access to account information without being authenticated to that account. Many such attacks work on the premise that access to confidential information is not properly authenticated.

Another common attack sends malformed input to a program. Many Web pages use forms extensively to drive the application, while the data input on the form is checked using client-side scripting. An attacker can take advantage of the fact that many online

21 · 16 WEB-BASED VULNERABILITIES

application developers assume that the client is going to use the form properly and that the scripts will check all input sent to the site. The attacker can examine the data stream sent by the form and then, rather than using the form, send a modified data stream via URL within the browser or incorporating it into a custom command of the attacker's choosing. An attacker may be able to obtain access to the application by placing system commands in the input stream. If the input stream is subsequently used in a *system()* call by the online application, the end user may force the execution of system commands on the attacker's behalf.

Some application developers rely heavily on hidden fields in the HTML document. Hidden fields allow the Web developer to include information on the page that is not displayed, although the end user can see the hidden field data simply by viewing the HTML source. The mistake application developers make is, first, in believing that the end user cannot see the hidden fields, and, second, in relying on the integrity of the hidden field data for making online decisions. Some online merchants have made the mistake of including pricing information for items in the hidden fields and using those prices to determine the cost of the online transaction. The end user can simply change the pricing in the hidden fields and send lower prices back to the merchants, for a discounted purchase.

Another misuse of hidden fields is to redirect application output. For instance, some Websites include file system path information in hidden fields on their Web pages. This information is used by a server-side script to determine where to read or write transaction information. Attackers, simply by changing the hidden field, can overwrite files or read files to which they should not have access. In some cases, it may be possible to store a program entered in a form field to be used later as a means of running a privileged shell on the remote system.

In summary, rigorous software quality assurance is necessary throughout the design and development of all Web-enabled and e-business applications, including front-end Web pages, application middleware, and operating systems. Once the software is believed to be immune to application misuse and subversion attacks, the system administrator must perform other activities to ensure the security of the e-business middleware:

- All unnecessary scripts or application server programs must be eliminated from the production server.
- Source code of application middleware must be carefully guarded against download or unauthorized access.
- Proper configuration of the CGI and application middleware is necessary to ensure executable access only to the correct application middleware, with the lowest practical privilege level. Sanity checking of inputs to application middleware must be done to ensure that only well-formed input is accepted.
- Testing of application code with the use of programs such as WebInspect by SPI Dynamics (www.spidynamics.com) that are specifically designed to discover flaws within Web-based applications.

21.6.6 Web Server Exploits. Web server security has been written about and covered in detail, including in *E-Commerce Security: Weak Links, Best Defenses*¹² and *The Web Security Source Book*,¹³ among other titles. Here we highlight some of the common exploits of Web servers used against many e-businesses.

COMPONENTS AND VULNERABILITIES IN E-COMMERCE SYSTEMS 21 · 17

21.6.6.1 Configuration. The key to Web server security is its configuration. Like other complex pieces of software, Web servers are highly configurable to meet the needs of any given site. By default, most software vendors configure the software application for maximum functionality but minimum security. Thus, by default, when the server is first started, it is likely to be more permissive than any given company's security policy would like. The principal premise for configurability is in the variations that exist among sites.

A correctly configured Web server is the result of a policy that defines what access is allowed, to which individuals, for each resource. This policy, in turn, is used to configure routers, firewalls, and all public servers, such as the Web server. Configuring the Web server, while necessary, is by no means sufficient to secure the system. The discussion of application server exploits in Section 21.6.5 demonstrates this principle.

21.6.6.2 HTML Coding and Server-Side Includes. Once the Web server is configured as securely as possible, it is important to ensure that the Web pages themselves do not open holes in the security. Many Web page developers fall into some common pitfalls that may compromise the site's security. The preceding section mentioned the problem of relying on hidden fields in HTML for security or business-critical data. Users can abuse the hidden field data to subvert an application.

HTML offers other potential vulnerabilities. One that is most often criticized is *Server Side Includes* (SSI). The SSI configuration option, if enabled, allows directives to be embedded in HTML that the server will execute. For instance, if the next statement were embedded in an HTML document, it would direct the server to display the contents of the system password file:

```
<!--#exec /bin/cat /etc/passwd -->
```

Certainly, Web pages should not be written with SSIs without a compelling reason. Although access to the HTML code is normally under control of the site, there are many ways an attacker might get an SSI into an HTML page. First, the attacker may have found another way into the system (e.g., by a CGI script exploit) but may want to provide either an easier backdoor or a redundant backdoor in case the CGI script vulnerability is found and closed. Once the CGI script is exploited, the attacker may implant an SSI directive within one of the site's HTML pages. Another way for an attacker to gain access is to have the server generate a Web page with the SSI of choice embedded in the HTML. How can an attacker do this? One approach exploits a server that generates dynamic HTML depending on the end user's data, preferences, or history. If the Web server ends up using some of the attacker's data to generate the HTML page, the attacker may be able to insert an SSI directive in the HTML. In summary, a better solution than continuously monitoring the HTML pages is simply to disable the SSI. In that event, even if SSIs were embedded, the server would not execute them. This is a configuration option, and like all system configuration files, the Web server configuration file should be protected by both file permission protection and file integrity checks, to ensure that it cannot be tampered with.

Although SSIs are often highlighted, a more common risk results from keeping documents or files in a publicly accessible portion of the Web server. The accessible portion of the Web server is called the *document root*. This root specifies the portion of the file system that the Web server can read and display to a Web client if requested. The document root can be a superset of the Web pages that are actually displayed when the user clicks through a Website. There may be other documents in the document root

21 · 18 WEB-BASED VULNERABILITIES

that are not linked to or from Web pages. This does not mean they are not accessible, however. Simply giving the correct address for any document will result in either displaying or downloading the document to any Web client. Therein lies the problem.

It is worth noting that HTML content is no longer found only in files served by the Web. On the local hard disk, Microsoft Windows applications often install Help documents comprised of CHM (“chum”) files and binaries. An attacker could exploit this vulnerability, leading to elevated privileges to the Windows XP “My computer” zone; unfortunately, local files are often overlooked and considered “trusted” during penetration tests.¹⁴ Recognizing these risks, Microsoft Corp. began use of Microsoft Assistance Markup Language with the Windows Vista OS.

21.6.6.3 Private Documents in Public Directories. Private documents inadvertently placed in a public directory can result in a compromise of confidential information and loss of privacy. For example, if a database file of credit card numbers (say, *cardnumbers.mdb*) were stored in the document root of the Web server for a fictitious company, *mycompany.com*, this URL address typed in a Web browser could download the file: www.mycompany.com/cardnumbers.mdb

This risk is even greater if *directory browsing*, another configurable feature, is enabled. Directory browsing allows an end user to view the contents of the file system at a given directory level if a Web page does not exist with the same name. Directory browsing is really a way to explore the file system of another site using a Web browser. Users may view this feature if they go back one level from a given Web page by deleting the right-hand section of a page’s URL and viewing its higher-level directory (e.g., in <http://a.b.com/d/e.htm>, one would remove the “e.htm” thus attempting to browse <http://a.b.com/d/>). Attackers can learn a lot of valuable information from viewing the contents of a directory including private files. Furthermore, the browser itself provides for clicking on a file name in the directory structure, which causes the file to be downloaded. Again, directory browsing, if enabled, is an open vulnerability and an unfortunately easy way to download private or confidential information.

21.6.6.4 Cookies and Other Client-Side Risks. Another potential vulnerability for an e-business is the use of *cookies*. Because HTTP is a *stateless protocol*, each new Web page that is visited has no memory of the last Web page that was visited by that user. Cookies are used to “keep state” between different Web pages visited in a given session. Cookies can make an e-business transaction appear to be seamless, sequential, and coordinated.

Most people, when discussing the risks of cookies, focus on the client-side privacy risks. Although these certainly exist, cookies also pose risks to the businesses that employ them. If the information contained in cookies is trusted, much the same way that the content in hidden fields is trusted, then the e-business may be vulnerable to cookie exploits called *cookie poisoning*. Some Websites use cookies to carry authentication information for a given user who traverses its pages. Once users have authenticated themselves, their token of authentication may be carried with them via cookies from one Web page to all subsequent pages at that site. Using cookies is a fairly weak form of authentication. The cookie can be stolen easily by someone snooping a user’s local area network or the Internet and then, with the information gained, to access the user’s personal pages on the Website. Secure protocols such as SSL should be employed to mitigate this risk.

Cookies also are used for other purposes that can introduce new vulnerabilities into critical transactions. Because cookies are under the control of end users, they can be

COMPONENTS AND VULNERABILITIES IN E-COMMERCE SYSTEMS 21 · 19

changed in whatever manner a user chooses. If cookies are designed to instruct a Web server where to write a customer-specific file, by changing the cookie data, an end user might overwrite other customer files or even replace critical system files. Similarly, if cookies are used for carrying order information, as is common in electronic shopping carts, then changing the contents of the cookies would corrupt the transaction. This could result in an unauthorized deep discount to the customer. This example and other cross-site scripting vulnerabilities are quite prevalent on today's Web. Browser exploits, HTTP header injection, cross-site request forgeries (CSRFs), and numerous flavors of phishing are all examples of Web client/user attack.

Regardless of the technology used, it is important to examine the features and functions of Web servers from a security-critical viewpoint. Dependence on a specific technology for critical transactions demands that the technology be trustworthy, so that it does not provide vulnerable points of attack.

21.6.6.5 Embedded Identifiers in URLs. Sometimes registration or unsubscribe messages include personalized URLs; for example, "To unsubscribe, click on <http://www.some-company.com/unsubscribe/?=12345>" The problem is that the identifying code (12345 in the example) can be replaced by any other number—and some of the generated numbers will unsubscribe other subscribers. Another unfortunate use of such a technique occurs in email messages for participation in Webinars; sometimes the coded URL in the email goes directly to a registration page with information filled in—information such as the precise name, email address, and employer information drawn from the database of previous contacts.

In general, it is poor practice to provide such short identifiers that a user or a simple program or batch file could be used to delete valid records in a database or to harvest confidential data from a Website. For example, a simple html file in Adobe Acrobat can be used to force the program to access every URL in the file; if a malefactor has created, say, 20,000 unique identifiers and all or most of them are affiliated with real people, the malefactor could easily affect the accounts or collect the information shown on the customer-specific pages.

To avoid such problems, either create an address space with far more possible identifier values than a 1:1 mapping of real accounts to identifiers; e.g., for a 100,000-person list, use-generated values using a string of 20 positions with 94 values per position available in uppercase, lowercase, numbers, and special characters, resulting in $94^{20} = 10^{39}$ keyspace, far beyond anything manageable by an amateur.

21.6.7 Database Security. Databases traditionally have represented an important intellectual property of information-based companies. As a result, they have almost always been company proprietary and unavailable to public access. In the new model of business, however, many of these proprietary databases are made available on the Internet, often without careful consideration of the risks involved. The Web browser becomes the database query interface, often to unknown and untrusted entities.

Although there has been much research in database security over the last two decades, the commercial sector has adopted only two key tenets: authenticating and authorizing principals to certain objects. Access to databases thus is controlled by properly authenticating the credentials of the requesting principal and then verifying which objects the authenticated principal is authorized to access. Any online database application must perform these two functions rigorously to protect the most valuable assets of e-business as well as customer privacy.

21 · 20 WEB-BASED VULNERABILITIES

Although many vendors claim secure channel access from the Web server to the database, there are many pitfalls. To start with, the fundamental reason why databases are vulnerable is that Web interfaces are commonly appended to what once may have been a closed and proprietary interface, without concern for overall security. Second, unsecured middleware programs such as CGI scripts or application servers usually mediate access from the Web server to the database. Web servers can provide client authentication from simple user name and password entry to strong certificate authentication. The needs of the business and the size of the accessing community will dictate which solution is feasible.

Despite obvious security advantages, most users do not store encrypted information in their databases, primarily for performance reasons. Encrypting and decrypting information on the fly during search, retrieve, and store operations can be too slow for real-time transactions. Also, even encrypting the data in a storage unit would not provide complete protection, as the online application must be able to read from and write to the database in clear text. Application-based attacks would still be able to get at the data while it was in plain, unencrypted text format.

Another key vulnerability of online databases also arises from application-based attacks. As described earlier, attacks that exploit vulnerabilities in the business-application logic often can provide unrestricted access to a database.

Attacks that exploit a buffer-overflow vulnerability in an application server program usually will be able to get command shell access on the remote server. From there, the attacker usually is able to find source code for the application server programs, such as Perl scripts or even C code, that are used to access the database. Because these programs *need* access to the database, they also must know the passwords used to access the various data partitions. If the programmers have foolishly *hard-coded* the passwords, then simply reviewing the source code may be enough to discover these passwords. With passwords in hand, the attacker can use the application server program via the Web interface to gain unauthorized access to the database. More directly, with a password it is possible to query a database from the command shell, using SQL commands or commands from the database language of choice.

Finally, like the other complex programs that run e-businesses, databases must be securely configured. Basic steps that the database administrator (DBA) needs to take include:

- Enforcing Web client authentication to the database
- Enforcing Web client authorization for access to database records
- Eliminating default database and database platform accounts
- Ensuring that passwords are read from encrypted files, not stored in program code
- Changing easily guessed passwords
- Configuring and maintaining internal access controls
- Auditing log files for suspicious activity

Account maintenance can be a key vulnerability in database management. Often the database software vendor will create a DBA account with an easily guessed password. Worse, DBAs may use the default account and password distributed with the database installation. This permits an attacker who has knowledge of the default passwords to gain access to all portions of the database by assuming the identity of the database administrator.

COMPONENTS AND VULNERABILITIES IN E-COMMERCE SYSTEMS 21 · 21

21.6.8 Platform Security. One area alluded to earlier in this chapter concerns the platforms that host components of an e-business. The platform, or operating system, represents the foundation of the e-business, but it is a potentially weak link in security. If there are cracks in the foundation, there is very little that even strong application software can do to keep the business secure. Therefore, it is imperative that system administrators properly patch platform vulnerabilities and maintain the security of the platform itself. As mentioned earlier, firewalls can go a long way toward blocking access to platform vulnerabilities by unauthorized outsiders. However, authorized but unprivileged users within the firewall can exploit known platform vulnerabilities to yield root privileges on a business-critical machine. Outsiders able to gain user privileges on the platform through any means also may be able to penetrate platform holes into severe security breaches.

Some key steps necessary to maintain platform security include:

- Eliminating default accounts generally installed with the operating system
- Prohibiting easily guessed passwords
- Enforcing password expiration
- Deactivating any unnecessary services that may be running by default
- Regularly applying security patches to the operating system
- Updating the operating system to its most recent release
- Ensuring that file access permissions are properly enforced, so as to prevent unnecessary access to critical files
- Enabling audit logging with intrusion monitoring
- Running system file integrity checks regularly

Some administrators believe that deploying a firewall is an acceptable substitute for configuring their platforms securely. Furthermore, with the plethora of different platforms running their enterprises, many system administrators give up on installing the latest OS patches and on securely configuring their platforms. They mistakenly assume that firewalls will protect them from all threats, even without maintenance. Unfortunately, relaxing host security can make the job of hackers easy. By hopping from machine to machine, they usually can find the valuable information they are looking for, or if that is their goal, they can wreak maximum damage.

An increasingly important platform for considerations of Website security is smart phones. By the end of 2012, there were 2.1 billion mobile Web users in the world at the end of 2012.

According to estimates by The ITU (June 2012), there were 2.1 billion active mobile-broadband subscriptions in the world. That is 29.5 percent of the global population.

- Mobile-broadband subscriptions have grown 40 percent annually over the last three years.
- Mobile-broadband subscription outnumber fixed broadband subscriptions 3:1.
- In developed countries mobile-broadband users often also have access to a fixed-broadband connection, but in developing countries mobile broadband is often the only access method available to people.

21 · 22 WEB-BASED VULNERABILITIES

- Ericsson (November 2012) forecasts that global mobile broadband subscriptions will reach 1.5 billion at the end of 2012, and 6.5 billion in 2018. The mobile phone will continue to be the dominant mobile broadband access device.¹⁵

Tablets also have their own versions of browsers. Web designers are increasingly having to cope with significant differences among platforms accessing their code.

Craig Smith included the following points in a January 2013 article entitled “Optimizing Ecommerce for Tablets and Smartphones”:

- **Improving navigation and usability:** Users should be able to select options easily and correctly; avoid drop-down menus and buttons that are close together.
- **Leveraging responsive design:** Determine how to present the Website according to what kind of device is accessing it.
- **Determining the purpose of the access:** Distinguish between the types of queries that are most common—and different—across platforms. Optimize design for the most common types for each platform.¹⁶

Another question is whether e-commerce sites should depend on Web access or develop applications (*apps*) specifically for smartphone and tablet operating systems.¹⁷ A major advantage of apps is that they can be programmed to avoid the vulnerabilities common to mobile devices in these early years of widespread adoption.

21.7 SUMMARY. This chapter has provided an overview of the weakest links in Web applications, including Web clients, network protocols, front-end Web servers, back-end databases, application servers, and the platforms on which they run. Secure network protocols are necessary but certainly not sufficient for securing e-commerce. The vulnerabilities described here are largely based on software flaws that exist in the application layer of e-commerce transactions.

Perhaps the most common vulnerability in e-commerce systems is misconfiguration of software. Because the responsibility for software configuration lies with the user, a security policy must be implemented and enforced. Once a system is configured, it is important to subject it to third-party validation and testing. A third-party audit can ensure that the configured system, including routers, firewalls, servers, and databases, meets the specifications of the security policy. The system also should be tested periodically against well-known, common attacks as well as against newer threats as they arise.

Like any software system, commercial off-the-shelf software has flaws, many of which are security-critical. It is imperative that system administrators stay current with vendor Websites, with hacker sites, with newsgroups, and with wherever threats and patches to their software are released. Both the security and hacker communities are constantly at work finding flaws in COTS that software vendors and security firms generally correct quickly. Software consumers—those who buy and use the software—must do their part by applying all relevant patches to their vulnerable software.

For custom-developed software, which includes front-end Web pages, application servers, CGI scripts, and mobile content such as ActiveX or Java, developers must do everything possible to ensure that their software is not vulnerable to attack and that it does not infringe on users’ privacy.

NOTES 21 · 23**21.8 FURTHER READING**

- Anderson, R. *Security Engineering: A Guide to Building Dependable Distributed Systems*, 2nd ed. Wiley, 2008.
- Bhargav, A. and B. V. Kumar. *Secure Java: For Web Application Development*. CRC Press, 2010.
- Cobb, Stephen. *Privacy for Business: Web Sites and Email*. Dreva Hill, 2002.
- Elliott, E. *Programming JavaScript Applications: Robust Web Architecture With Node, HTML5, and Modern JS Libraries*. O'Reilly Media, 2013.
- Gabarro, S. A. *Web Application Design and Implementation: Apache 2, PHP5, MySQL, JavaScript, and Linux/UNIX*, 2nd ed. Wiley, 2013.
- Gookin, D. *Android Phones for Dummies*. For Dummies, 2013.
- Gunasekera, S. *Android Apps Security*. Apress, 2012.
- Harwani, B. M. *Android Programming Unleashed*. Sams, 2012.
- Kissoon, T. *Securing Web Applications*. Auerbach, 2013.
- Long, F., D. Mohindra, R. C. Seacord, D. F. Sutherland, and D. Svoboda. *The CERT Oracle Secure Coding Standard for Java*. Addison-Wesley Professional, 2011.
- Mark, D., A. Horovitz, K. Kim, and J. LaMarche. *More iOS 6 Development: Further Explorations of the iOS SDK*. Apress, 2012.
- McGraw, G., and E. Felten. *Securing Java: Getting Down to Business with Mobile Code*, 2nd ed. Wiley, 1999.
- Musciano, C., B. Kennedy, and E. Weyl. *HTML5: The Definitive Guide*, 7th ed. O'Reilly Media, 2014.
- Scambray, J., Shema, M., and C. Sima. *Hacking Exposed Web Applications*, 2nd ed. San Francisco, CA: McGraw-Hill Osborne Media, 2006.
- Six, J. *Application Security for the Android Platform: Processes, Permissions, and Other Safeguards*. O'Reilly Media, 2011.
- Sullivan, B. *Web Application Security, A Beginner's Guide*. McGraw-Hill Osborne Media, 2011.
- Welling, L., and L. Thomson. *PHP and MySQL Web Development*, 5th ed. Addison-Wesley Professional, 2013.
- Zalewski, M. *The Tangled Web: A Guide to Securing Modern Web Applications*. No Starch Press, 2011.

21.9 NOTES

1. This original version of this chapter, entitled “E-Commerce Vulnerabilities,” was adapted from Anup K. Ghosh, *Security and Privacy for E-Business* (Hoboken, NJ: John Wiley & Sons, 2001), chapter 4; adapted by permission.
2. Anup K. Ghosh, *E-Commerce Security: Weak Links, Best Defenses* (New York: John Wiley & Sons, 1998).
3. C. Cowan, P. Wagle, C. Pu, S. Beattie, and J. Walpole, “Buffer Overflows: Attacks and Defenses for the Vulnerability of the Decade,” paper presented at DISCEX 2000, January 25–27, 2000, Hilton Head, S.C. Proceedings of the DARPA Information Survivability Conference and Exposition (Los Alamitos, CA: IEEE Computer Society Press, 2000).
4. Cowan, et al. “Buffer Overflows.”
5. PCI Security Standards Council, Payment Card Industry Data Security Standard, v1.1. (September 2006), retrieved from www.pcisecuritystandards.org/pdfs/pci_dss_v1-1.pdf

21 · 24 WEB-BASED VULNERABILITIES

6. G. McGraw and G. Morrisett, "Attacking Malicious Code: A Report to the Infosec Research Council," *IEEE Software* 17, No. 5 (September/October 2000): 33–41.
7. Ghosh, *Security and Privacy for E-Business*, chapter 7.
8. B. Schneier, *Applied Cryptography: Protocols, Algorithms, and Source Code in C*, 2nd ed. (New York: John Wiley & Sons, 1995).
9. Ghosh, *E-Commerce Security*.
10. W. A. Arbaugh, W. L. Fithen, and J. McHugh, "Windows of Vulnerability: A Case Study," *IEEE Computer* 33, No. 12 (December 2000): 52–59.
11. L. Stein, *Web Security: A Step-by-Step Reference Guide* (Reading, MA: Addison-Wesley, 1998).
12. Ghosh, *Security and Privacy for E-Business*.
13. A. Rubin, D. Geer, and M. Ranum, *The Web Security Sourcebook* (New York: John Wiley & Sons, 1997).
14. Michael Howard and David Leblanc. *Writing Secure Code*, 2nd ed. (Redmond, WA: Microsoft Press, 2002).
15. mobiThinking, "Global Mobile Statistics 2013 Part B: Mobile Web; Mobile Broadband Penetration; 3G/4G Subscribers and Networks," <http://mobithinking.com/mobile-marketing-tools/latest-mobile-stats/b#mobilebroadband>
16. C. Smith, "Optimizing Ecommerce for Tablets and Smartphones," *Practical ecommerce*, January 15, 2013, www.practicalecommerce.com/articles/3869-Optimizing-Ecommerce-for-Tablets-and-Smartphones
17. D. Traxler, "Mobile Commerce: Website or App?" *Practical ecommerce*, January 1, 2013, www.practicalecommerce.com/articles/3862-Mobile-Commerce-Website-or-App-

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER **22**

PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

Franklin Platt

22.1 INTRODUCTION	22·2	22.3.8 Completing the Threat Assessment Report	22·14
22.2 BACKGROUND AND PERSPECTIVE	22·3	22.4 GENERAL THREATS	22·14
22.2.1 Today's Risks Are Greater	22·3	22.4.1 Natural Hazards	22·15
22.2.2 Likely Targets	22·4	22.4.2 Other Natural Hazards	22·16
22.2.3 Productivity Issues	22·4	22.4.3 Health Threats	22·16
22.2.4 Terrorism and Violence Are Now Serious Threats	22·4	22.4.4 Man-Made Threats	22·16
22.2.5 Costs of Damaged IS Infrastructure	22·5	22.4.5 Wiretaps	22·18
22.2.6 Who Must Be Involved	22·5	22.4.6 High-Energy Radio-Frequency Threats	22·19
22.2.7 Liability Issues	22·6		
22.2.8 Definitions and Terms	22·7		
22.2.9 Uniform, Comprehensive Planning Process	22·8	22.5 WORKPLACE VIOLENCE AND TERRORISM	22·21
22.3 THREAT ASSESSMENT PROCESS	22·9	22.6 OTHER THREAT SITUATIONS	22·22
22.3.1 Set Up a Steering Committee	22·9	22.6.1 Leaks, Temperature, and Humidity	22·22
22.3.2 Identify All Possible Threats	22·10	22.6.2 Off-Hour Visitors	22·22
22.3.3 Sources of Information and Assistance	22·11	22.6.3 Cleaning and Maintenance Threats	22·23
22.3.4 Determine the Likelihood of Each Threat	22·12	22.6.4 Storage-Room Threats	22·23
22.3.5 Approximate the Impact Costs	22·12	22.6.5 Medical Emergencies	22·23
22.3.6 Costs of Cascading Events	22·13	22.6.6 Illicit Workstation	22·24
22.3.7 Determine the Vulnerability to Each Threat	22·13	22.6.7 Other Local Threats	22·24
		22.7 CONFIDENTIAL THREAT INFORMATION	22·25
		22.8 SUMMARY	22·26
		22.9 FURTHER READING	22·27
		22.10 NOTES	22·27

22 · 2 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

22.1 INTRODUCTION. This chapter describes the wide array of possible physical threats that can affect information systems (IS) infrastructure. A *threat* is any credible situation—whether actual, imminent, predicted, or possible—with the potential for causing harm, damage, or disruption. The terms *threat* and *hazard* are used synonymously in this chapter and have the same meaning. *Risk* refers to the probability of occurrence of a threat. For a structured terminology for discussing threats and risks, see Chapter 8 in this *Handbook*.

The infrastructure affected can be any component of a computer system or communications network; any of the cables, wiring, or devices that transfer power or data; or any of the support services or utilities needed to sustain full IS performance. In addition, this chapter considers threats against employees as a significant component of physical security.

The speed and accuracy of any system is dependent on the performance of a long chain of physical components as well as on the productivity of all of the people who utilize or maintain each component. Anything less than full system-wide performance can be costly.

A physical threat is any event that can degrade the performance of an information system—whether such an event is actually occurring or imminent, or credibly likely or possible, or completely unexpected and without warning. The list of possible threats is a long one, beginning with natural and man-made external events that can degrade IS performance. Many other possible threats are internal, resulting from accidents or misuse or deliberate attack. Internal threats also include reliability failures, installation or maintenance issues, or lack of proper testing. Threats can also be caused by situations within the facility, building, or complex, or as a consequence of local or regional events, or, increasingly, events worldwide and even from outer space. Threats happen and can come from almost anywhere, so it is prudent and far less costly to try to anticipate all likely and significant threats.

Although in this chapter the terminology often refers to “all possible threats,” readers must understand the implicit assumption that they are considering only reasonable categories of threat. Planning to cope with the results of a world-destroying impact with a continent-sized asteroid is useless; trying to define a strategy for responding to an attack by extraterrestrial invaders armed with ray-guns is impossible and useless.

This chapter begins a threat-assessment process, which starts by identifying all reasonably defined threat situations. This process must be comprehensive and rigorous and involve all stakeholders. It is not enough simply to assign likelihood and an impact for the usual or the obvious threat situations. There must be a full risk analysis, so that risk, vulnerability, and response and recovery costs can be quantified, as explained in Chapter 23 in this *Handbook*. Otherwise, the whole process is valueless. Absent comprehensive planning, security becomes an irrational selection of vendors and solutions, which can only add cost and increase liability risks. Alternatively, a strategic risk-management process can maximize future profits and add value, protect morale and productivity, enhance goodwill, and improve relations with customers and communities.

The events of September 11, 2001, sounded a shattering wake-up call that unexpected threats can actually happen and that even the best security practices, products, and services are of little value without proper planning, implementation, and support. Hurricane Katrina, four years later in 2005, demonstrated again that both business and government are still unprepared, even though such an event was predicted many times.

This chapter suggests a comprehensive perspective on physical security that can add value and help avert disaster. Chapter 23 in this *Handbook* goes on to suggest physical

BACKGROUND AND PERSPECTIVE 22 · 3

security implementation and management that can protect people and property while optimizing morale, productivity, and profitability.

22.2 BACKGROUND AND PERSPECTIVE. Historic data and statistics are of limited value in predicting future threats. For a discussion of the unreliability of computer-crime statistics, see Chapter 10 in this *Handbook*. The past is no longer prologue because many new threats are emerging and many types of incidents are increasingly severe, widespread, complex, damaging, and costly. Proliferating and increasingly dispersed infrastructure system components are often vulnerable and hard to protect. And hybrid configurations of new and legacy systems vastly complicate the protection process.

There are few threat statistics that can reliably predict the infrastructure future. Computer-crime reports mostly cover logical security, while general crime statistics often do not relate directly to computer security. Historic information is further flawed because many incidents are never detected and far more are not reported for fear of embarrassment, liability, or loss of business. Many security incidents are masked as quality-control problems for the same reasons, or are misdiagnosed because no one had time to determine the true cause(s). Lacking reliable precedents, predicting future threats is especially difficult—yet increasingly necessary. Chapter 10 discusses computer crime statistics in more detail; Chapter 58 and Chapter 62 discuss risk-management methodologies.

Most incidents happen suddenly, without warning, and often where least expected. Many threats once thought to be unlikely now occur widely and strike with surprising intensity and devastation. Other contributing factors include poor risk management due to inexperience, denial, or complacency that can quickly turn routine threats into costly incidents. Businesses with good security preparedness can usually survive, while many others will not.

22.2.1 Today's Risks Are Greater. Today's threats are increasingly sophisticated, unpredictable, potentially serious, and increasingly commonplace as well. Disruptive incidents can result from mistakes or accidents, snoopers or hacking, vandalism, disgruntled or disruptive persons, labor disputes, demonstrations, civil unrest, extremists of many stripes, and, increasingly, both domestic and international terrorism. Although violent-crime incidence has decreased in recent years in some parts of the world, these data may be misleading because they rarely include workplace-related events. Violence in the workplace is now common and often without warning. Incidents can include harassment, bomb scares, robbery, hostage situations, shooting, or arson. And each and every one of these threats can seriously affect IS performance.

Physical threats can extend beyond direct attacks. Many scares do not result in actual violence, and disruptive incidents can occur outside of the workplace. However and wherever they happen, violence-related incidents are increasingly commonplace, disruptive, and costly.¹

Those who threaten the IS infrastructure must at least gain access to some physical part of it, but often this can be done inconspicuously, by trickery, deception, or simply forced entry. A physical attack may be the best way to compromise an information system—much more effective and less likely to be detected than a logical attack. Often, many vital system components are vulnerable, exposed, or easily accessible. These components include wiring and cable runs, connection and junction points, system and network equipment, and the utilities that support them. Attacks or spying by physical means are often easy, fast, safe, and sure.

22 · 4 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

22.2.2 Likely Targets. Businesses and organizations are increasingly likely to be the targets of hackers, disgruntled employees, competitors, disturbed persons, demonstrators, hate groups, extremists, and even terrorists. Motives may include a conspicuous, tempting, or challenging target, rage or revenge, an opportunity to cause reputational damage or at least adverse publicity, to make a political statement, for extortion or blackmail, or for personal gain or profit. Often, there are no discernable motives. And beyond being likely targets, most businesses are convenient, easy, and safe targets as well, because most are unaware and unprepared for today's threats, much less for future ones. Although government facilities remain preferred targets, many are now better protected than most businesses, thanks to new procedures that are covered in Chapter 23 in this *Handbook*.

Another likely threat arises from the need of extremists and terrorists to finance their activities. Many groups and all independent, self-directed cells are dependent on crime to finance their operations. Today's crimes can include robbery and holdups, counterfeit currency and credit cards, Internet phishing and scams, theft, extortion, blackmail, and selling pirated software and other knock-offs. In addition, many foreign governments, businesses, and criminal organizations are actively engaged in spying. Although these are mainly corporate security and logical security problems, they also represent potential physical threats that must be deterred.

22.2.3 Productivity Issues. Good security can be directly correlated with high productivity, which, in turn, can improve performance, customer satisfaction, and goodwill. Good security is measurable and can strategically enhance each of these factors to add both value and profit. Anything less than good security invites wasted time and money.

People who do not feel safe will not be productive. This applies to employees, visitors, vendors, and others on premises, as well as to customers, vendors, stockholders, and other stakeholders at remote locations. Everyone using any information system must be comfortable that physical safety and privacy are insured and that the system is uninterruptible, secure, and operating at full performance. Everyone concerned must be involved in the planning process, generally understand the threats, and support the security procedures. Otherwise, performance inevitably suffers.

Whenever a security incident occurs, morale and productivity are likely to plummet and can remain low for many weeks, or even months. Whether the infrastructure is actually affected or not, significant disruption of operations is likely. Even when there is no injury or damage, the perception of a potential event can be costly; it can disrupt productivity, lose business and customers, and jeopardize goodwill.

No one knows how often productivity-related events occur. Businesses generally do not report them, and neither do the media. But these things do happen and probably with considerable frequency. Yet all threats can be mitigated to some extent and many more prevented at far less cost than the consequences by effective technical preparations, adequate employee awareness and training, and well-planned and well-rehearsed responses.

22.2.4 Terrorism and Violence Are Now Serious Threats. Acts of terrorism and violence are now a reality that can occur anywhere in the world. September 11, 2001, and the events that followed have brought home the stark reality that violence can happen anywhere and can cause massive damage and disruption. And most major disasters can disrupt information systems far removed from the actual incidents.

BACKGROUND AND PERSPECTIVE 22 · 5

Workplace violence is also happening with increasing frequency and often at facilities assumed to be safe.² Bomb and biological or chemical scares, personal threats, harassment, hostage situations, and shootings are all happening with increasing frequency. Whether actual violence occurs or not, the threats alone, the many rumors generated, and the imagined proximity to danger are all productivity-related events that can seriously disrupt the performance of information systems for a long time. Therefore, these become infrastructure security issues that require special planning and should not be left to premises security personnel to prevent.

There are other serious threats from foreign intelligence, terrorists, and domestic groups generally unknown to the public. Some of these are explained well in the Project Megiddo report published by the Federal Bureau of Investigation (FBI) in 1999 in anticipation of the millennium. The report provides “an FBI strategic assessment of the potential for domestic terrorism in the United States undertaken in anticipation of or response to the arrival of the new millennium.”³ The threats cited then are basically unchanged today, except that more previously unknown threats have since been added. There are consultants with FBI contacts who can offer valuable advice.

Attempted violence is now a serious threat to all IS infrastructures. However, thorough security planning can do much to avoid trouble and needless expenses.

22.2.5 Costs of Damaged IS Infrastructure. Direct costs of system downtime can exceed many thousands of dollars per hour. The losses include the slack-time costs of people who cannot use the systems, costs of support and maintenance people diverted to restoring operations, recovery expenses, overtime, and often lodging, food, and travel expenses during recovery. Usually, many outside resources are needed for response and recovery. Everything quickly becomes expensive. Often, to further compound the costs, needed resources are just not available immediately.

Indirect costs can also be significant. Reestablishing and keeping good public relations can be expensive. Often, public announcements, news releases, and briefings to the news and financial media and to stockholders are needed to neutralize public embarrassment and control rumors. Key customers must be contacted and reassured, and pending orders rescheduled. Still more costs include lost business or market share and dropping stock prices. Competitors often will take as much advantage as they can, which necessitates further costs defending brands and reputation.

Any number of such costs can devastate an enterprise unless strong security measures are deployed effectively and quickly. Ad-libbed responses are often disastrous. In reality, an infrastructure outage of more than a few hours is often fatal, and the enterprise can never fully recover. Good security usually can prevent disastrous costs.

22.2.6 Who Must Be Involved. Many threats to information systems also involve corporate or premises security personnel, whose role is to protect people and property within and surrounding the workplace. The CEO and the CFO are also involved, because they must now comply with the applicable laws and regulations. Compliance requires that events that might materially affect financial performance must be included within an organization’s financial filings and statements. Major security incidents are clearly such events whose impact must now be predicted—as described in this chapter.

Premises security often includes little more than guards, access control, and some surveillance. Their understanding of the IS infrastructure’s special security needs is often minimal, lumped in with overall physical security procedures, and rarely extends much beyond the immediate premises. Yet, IS security requires additional knowledge,

22 · 6 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

experience, protection, and support. Good security must be strong, fast-acting, focused on specific targets, and closely monitored. Effective early warning systems are necessary in order to prevent threats from happening. In reality, each security function will have its own needs and priorities and use its own resources. During a serious incident, security for the premises, occupants, information systems, and the infrastructure must all coordinate efficiently and effectively. They must also work smoothly with local fire and police departments, with other emergency responders, and with many outside resources.

Who, then, should manage the process of determining the threats to the IS infrastructure? And who are the stakeholders who should be involved in this process?

The best person to manage the physical security of information systems is one who knows a great deal about possible threats and about the IS infrastructure. The office manager, facilities manager, or corporate security director or their staffs are usually ill-equipped to determine or manage IS security. Often, too, the chief information officer (CIO), chief information security officer (CISO), and IS security officers (ISSOs) deal with protecting data and data processing, and are not the best persons to understand IS physical security, especially in a large installation. Therefore, the best person is a trusted individual with the right knowledge and experience and with enough time to manage the process well. Planning and managing infrastructure security, implementing and testing it, training and security awareness, monitoring, and periodic updating of the system and procedures are a full-time job in most organizations, and require a support staff in larger organizations.

Chapter 65 in this *Handbook* discusses the role of the CISO in detail.

Another consideration is that no one person should know all the secrets, which is a rule that has become especially important in dealing with the IS infrastructure. When trouble comes, many experts must mobilize very quickly, efficiently, and effectively. To do this, the responders must be familiar with all the information systems and have fast access to the infrastructure. No matter how well trusted, no one person should know everything about the physical and logical defenses. (If need be, such information should be safely stored with strong access controls.) It is wise to divide the secrets so that no one group knows or has access to them all. Having done this, multiple persons can then share each portion of the secrets and observe each other, so that no one person is indispensable. Chapter 45 in this *Handbook* presents details of employee management and security.

But only the security vulnerabilities and defensive implementations should be secrets. There is no point in letting others know about the defenses and where the organization is vulnerable. However, the process of determining the underlying threats should be common knowledge among all the stakeholders involved. By identifying a large number of threats that have been assessed—but keeping the likelihood, vulnerability, and impact information secret—we can discourage would-be troublemakers by at least providing some idea of the *number* of threats that have been considered. With any luck, they will attack less-prepared sites.

22.2.7 Liability Issues. Aside from the need for efficient emergency response, another issue that is becoming increasingly important, is potentially very costly, and is often overlooked. This is the issue of liability. Every organization has a legal and fiduciary duty to protect the people and property within and surrounding its premises. If any injury or damage occurs in or around a workplace—even long afterward—allegations of negligence will likely follow. The motive is often that damages awarded by the courts can be very large and lawyers often take on these cases on speculation, in hopes of

BACKGROUND AND PERSPECTIVE 22 · 7

receiving very large fees. And whether the organization was actually at fault or not, the resultant legal fees, the time and costs needed to defend the organization, bad publicity and possible loss of business, and the eventual fines and awards can be devastating.

If negligence is alleged, the issue is whether the organization was properly prepared for the emergency and whether its response was effective. The question is simply: Did management perform its duty to protect the organization? An affirmative answer would require at least:

- Evidence of a thorough threat assessment process
- Good security plans, policies, and procedures that have actually been implemented
- Proper training and current security awareness
- Periodic drills, exercises, security reviews, and feedback from known events
- Periodic updates to assure that security remains effective

If it can be demonstrated that all these measures were not taken, or that there were deviations from generally accepted standards, the result could be punitive as well as compensatory damage awards. Gross negligence will probably be alleged as well, in which case insurance may not defend the accused organization or individuals, who would then be personally liable. Even with insurance in effect, it may not be sufficient to cover the very large penalties that juries commonly award.

Chapter 60 in this *Handbook* discusses insurance for responding to damages in IS.

Even when all the correct answers are given, the accused often find they are considered to be guilty until they can prove themselves innocent. And this can be a long, painful, and costly process.

However, a nearly iron-clad defense against liability is that federal procedures were followed. The most comprehensive of these are from the Department of Homeland Security (DHS) and its Federal Emergency Management Agency (FEMA) Directorate.⁴ Of all of the other security planning and management procedures, which are many and varied, only the DHS/FEMA methodology is likely to become generally accepted. The procedures are already well known, uniform, and comprehensive, and they consider all threat possibilities and all response resources. It is unlikely that a plaintiff's attorney will ever allege that these procedures are deficient or ever bring a case once it appears that an organization is compliant. This is not to remove all liability, merely to reduce the scope to situations usually covered by insurance. Chapter 23 in this *Handbook* outlines how to plan, implement, and manage the DHS/FEMA procedures.

22.2.8 Definitions and Terms. Information infrastructure security is simply a means of IS performance assurance. Good security precludes any IS disruption that might degrade performance in any way. Any slowdown or loss of productivity, loss of data, breach of privacy, or disruption of the systems, networks, and utilities that support any information system diminishes performance. Good security assures that all systems remain fully operational, robust, and accurate, and that all their data remain private and cannot be compromised.

There are three elements of information systems security. Each element must be especially designed and maintained to protect against different threats of varying scope and intensity. The three elements include:

1. **Logical security**, which is also known as information systems security, protects only the integrity of data and of information processing.

22 · 8 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

2. **Physical security**, which is also called infrastructure security, protects the rest of the information systems and all of the people who use, operate, and maintain the systems. Physical security also must prevent any type of physical access or intrusion that can compromise logical security.
3. **Premises security**, which is also known as corporate or facilities security, protects the people and property within an entire area, facility, or building(s), and is usually required by codes, regulations, and fiduciary obligations. Premises security protects broad areas. It often provides perimeter security, access control, smoke and fire detection, fire suppression, some environmental protection, and usually surveillance systems, alarms, watchmen, and guards. Premises security is often an extension of law enforcement.

Clearly, there is much overlap between each element, and a threat to one is often a threat to the others as well. However, each element is likely to perceive and handle each threat differently. Although the remainder of this chapter deals with physical security, some threats that are usually considered the realm of logical or premises security must also be covered.

22.2.9 Uniform, Comprehensive Planning Process. The threat assessment process can go by a number of names, which are all functionally similar. Among the terms are emergency, disaster, operations, contingency and crisis-response planning, and damage control. Some processes are proprietary and do not share a common language or standardized procedures that can be understood by everyone involved. Many regulated industries must develop emergency response plans using terms and formats dictated by the regulating agency, which makes these terms and formats proprietary as well. There are many diverse examples of regulated industries that use hazardous or nuclear materials or operate dams.

Many organizations that are involved in emergency response—such as hospitals and emergency medical services, schools, the American Red Cross, and many volunteer agencies, National Guard and military units—may still use other terms and models. Government agencies at the local, state, and federal levels still use a wide variety of security plans and procedures, even though there is now one standard methodology required, as explained in Chapter 23. Most of these procedures are not uniform or comprehensive. And they can be incompatible, present major communications barriers, and cause unnecessary misunderstandings, delays, and wasted resources. In turn, most private organizations are not aware of the government procedures in place or of the many resources that can assist them.

Every organization experiences essentially the same threats and has limited response capabilities and resources. Almost every organization is likely to be overwhelmed in a major emergency, and everyone can benefit enormously from outside resources. Yet, each venue tends to use dissimilar models, terms, and procedures that are often unintelligible to the others. In many cases, their models become file-and-forget plans and procedures that no one understands or accepts.

If only as prudent risk-management strategy to avoid liability issues, a single, uniform, and comprehensive standard for processes, procedures, and language is needed. The best way to do this is to adopt the DHS/FEMA methodology, which encompasses all hazards, coordinates all resources efficiently, and is clear, concise, understood, and accepted by everyone involved. Details are given in Chapter 23. This is now required for all federal government departments and agencies, and for state and local jurisdictions as well, if they wish to continue receiving federal grants and assistance.

THREAT ASSESSMENT PROCESS 22 · 9

A useful public/private standard for evaluating disaster/emergency preparedness is published by the Emergency Management Assessment Program⁵ (EMAP), which is a nonprofit professional organization. The EMAP standard covers all the current DHS/FEMA methodology and provides a clear and concise method to determine compliance, as well as a self-auditing process. EMAP also provides an independent, outside security audit of any government organization to ascertain that its emergency management program complies with the federal requirements; it may soon be able to audit private organizations also. The audit is done at the applicant's facility by professional, well-trained volunteers. The whole security-audit process is both rigorous and inexpensive. The EMAP is therefore a suggested alternative to a CPA-provided security audit as well as to the audit procedures of other professional organizations that generally do not include either infrastructure security or compliance with federal procedures.

Security preparedness that is compliant with these standards may become a requisite for the private sector to obtain insurance, to avoid liability and excessive costs, to establish innocence, or to use capital markets for risk management. Knowing these standards will help the threat assessment process and, later, can provide better protection.

22.3 THREAT ASSESSMENT PROCESS. Effective security planning begins with a thorough threat assessment. This process begins by establishing an ad hoc organization, obtaining budget approval and the sponsorship of senior management, and formation of a steering committee that represents all the stakeholders.

The first four tasks are to:

1. Identify all potential threat situations.
2. Determine the likelihood and estimate the direct and indirect costs of each threat.
3. Evaluate and prioritize each threat.
4. Prepare and present a final report that each committee member signs.

This report becomes evidence of due diligence to avoid liability, and becomes the basis for protecting the infrastructure, which is described in Chapter 23.

The way *not* to go through the planning process is for a few people to decide on the threats and risks and then write the security plan. This is the opposite approach to good and effective planning. And it will not work. Very few of the key personnel will accept or even understand the plan, and it will probably be ignored or overlooked when the next emergency arises. Good examples of failed planning are the emergency management plans for New Orleans and Baton Rouge areas that were put in place before Hurricane Katrina in 2005. Apparently, both were good plans. But few officials understood the plans or remembered to use them.

The threat assessment planning process *must* be done thoroughly and completely. For additional perspectives on risk assessment and risk management, see Chapters 62 and 63 of this *Handbook*.

22.3.1 Set Up a Steering Committee. The security planning process is not effective unless all of the stakeholders are represented, and the best way to do this is to establish a steering committee. The committee will help identify and evaluate potential threats, and help develop a comprehensive protection plan.

22 · 10 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

The committee should represent all stakeholders and include as much experience, knowledge, and perspective as possible. Stakeholders should include users, administrators, management, key partners, customers, vendors, and service providers (such as maintenance, repair, and cleaning personnel). A project manager should participate, as should legal, financial, and human resources representatives. Independent experts and facilitators are recommended as well. The best committee chair is usually an outside facilitator who is immune from political or cultural bias, loyalties, or product preferences. It is also wise to seek input from stockholders, lenders, insurers, and community and government officials.

This can be a virtual committee that rarely needs to meet in full session. Communications by email or phone should be sufficient, and an in-house staff can gather data, issue reports, and work with individual committee members. Confidentiality is not an issue at this point. The committee's purpose is to identify threats and to assist in planning. No member need have full knowledge of the actual vulnerabilities or the resulting security systems and protections. Each member, however, should sign off on the committee's final report.

Once its initial mission is completed, it is best for this committee to convene at least annually to review new and changing threats and to assess how well the security systems have performed. This committee not only represents all concerned stakeholders, but it can provide oversight so that management does not neglect to exercise, review, and update the security plan. In effect, the committee provides due diligence that management has fulfilled its fiduciary responsibilities.

22.3.2 Identify All Possible Threats. The first step is to identify all possible threat situations that might affect the information infrastructure. This is not to say that each threat will someday happen but that, conceivably, it might, nearby or even at a distance. The causal connection could be tenuous at best and the event most unlikely, but it could someday happen. It is far better to think objectively about such things and to plan effectively than to simply write them off as events that could "never happen here" or to claim that, if such a disaster does happen here, "there is nothing we can do to mitigate it." Both statements are patently false and potentially very costly.

The threat list can be very long. The format should tabulate each specific threat with a one-line description to clarify what is meant. For example, "flooding" is too general a term to meaningfully assess its risk. Instead, one might list "riverine flooding," which could be the result of heavy rains, snowmelt, or a breached dam or levee; "local flooding," such as a water-main break, severe storm, or a major fire nearby; or "premises flooding," due perhaps to leaking pipes, drains, or leaking windows, roofs, building setbacks, a burst water tank or cooling tower, or fire suppression. The purpose is to have a long list of specific threat situations that can be individually assessed. There can easily be 100 or more threats listed. As the project proceeds, those involved will surely add still more threats and fine-tune the definitions. The threat list should include threats that are direct as well as indirect and threats affecting the general area, the region, or the whole country. It is best to list all threats, including those that are deemed unlikely to occur, as a precaution against inadvertently omitting some that might have more importance than is apparent at first glance.

Unfortunately, threats tend to cascade: One threat situation can create others, and these, in turn, can create more threat situations. For example, severe flooding can close roadways and keep people away from work, impede delivery of supplies necessary to operate (such as food or water or fuel for an emergency generator), disrupt communications, cause mud slides, ignite fires, trigger looting, or cause civil unrest. As

THREAT ASSESSMENT PROCESS 22 · 11

another example, during a fire there will likely be loss of electrical power for equipment cooling, ventilation, or communications, and released water or chemicals that must be contained. Each of these conditions is a separate threat with its own risk. In reality, few hazards occur in isolation, and many of the cascading events may be unexpected and unpredictable, and can themselves trigger still more events.

Therefore, the tabulation of each threat should also include two other columns: (1) events that could trigger the threat, and (2) cascading threats that could result from the threat. This is easily done by numbering each threat and referencing the numbers that may interconnect the various threats. The information is needed to evaluate the impact of each threat. It can only be a general indication for planning purposes of what might happen. The actual cascading will probably be different, but the analysis of the threats will still be valid.

Do not rely on force majeure—for example, a major earthquake—as an excuse that a particular threat may be unavoidable. Force majeure does not limit liability because such events can be anticipated and some steps taken to minimize injury and damage. Nor may acts of war serve as an excuse, for the same reasons.

Finally, the threat list should be divided into major categories such as natural events, man-made events, vandalism, attacks, support system failures, and so on. The remainder of this chapter builds on the threat list and establishes categories best suited to particular needs.

22.3.3 Sources of Information and Assistance. The next step is a compilation of past events: historical records and details as to what happened, when, and how, what injury and damage resulted, was there warning, how fast was the onset? The compilation should include events in both near and distant areas that could possibly cause damage and disruption, either directly or indirectly. Perhaps the 2005 subway and bus bombings in London or the earlier commuter-train bombings in Madrid could affect American workers who commute and become concerned about their own safety.

As the threat list gets bigger, it is necessary for everyone involved to think out of the box and to look for scenarios that past trends suggest could eventually happen. This thinking can be productive when it examines possible scenarios objectively. The author, as a security consultant in New York, worked with tenants in the World Trade Center and tried to suggest disaster planning but was overruled by the landlord who assured everyone that this was “the safest place on earth.”

Information sources to investigate include the weather bureau, libraries, local fire and police departments, state and county officials, utility companies, newspaper files, and local knowledge of past events within the region. Power and communications utility companies can provide outage reports, but their terminology must be clearly understood. For example, an electrical “outage” may only include interruptions lasting more than five minutes. Regional and state regulatory authorities, public utility commissions, industry and professional organizations, and business development groups may also provide useful information, once their perspective and terminology are clarified.

All local, regional, and state emergency management agencies should now have an all-hazard mitigation plan that lists all threats that have occurred throughout their jurisdiction. There should also be dates, descriptions, and map locations in their plan; this information may go back historically for a century or more. And each agency should have a local emergency operations plan that lists all potential threat situations within their jurisdiction. Look particularly at the appendices and annexes that discuss particular types of threats. In addition to what is in their plans, each agency should be of considerable assistance describing possible major threats that are man-made, involve

22 · 12 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

hazardous materials, or are major health risks. The threats they identify are probably your threats as well.

Also, work with each stakeholder for advice on threats they know about or perceive. Ask each person to reach out to their customers, vendors and suppliers, service providers, business neighbors, and consultants and academics as well. There will be much expertise, perspective, and good advice gleaned from these sources. The time seeking it will be well spent.

At the end of this chapter is a list of some reference books that can be helpful. Although some of the procedures are outdated, there is still much useful source data.

22.3.4 Determine the Likelihood of Each Threat. The steering committee should consider the likelihood that each threat may happen. The best way to estimate likelihood is as an annual probability on a relative scale (for example) from 0 (none) to 5 (very likely). A perceived likelihood of once every 100 years (which equals a 1 percent annual probability) could be rated a 1. A 5 percent annual probability might be rated a 2, while an annual probability of 20 percent or more had best be rated as a 5. Threats from severe weather should be rated higher than historical data suggests in order to account for the world's changing weather patterns. Threats with a likelihood of 0 should also be included, as this value can be used later to show the relative vulnerability of each and every threat.

The steering committee can adjust each rating as it deliberates so that the final results are useful and meaningful.

22.3.5 Approximate the Impact Costs. Once all the possible threats are determined and the likelihood of each is estimated, the next step is for the steering committee to deliberate the impact of each threat. Impact considers the potential losses and costs associated with each threat. (Do not consider cascading threats at this point.) Here also, a relative scale of 1 (very low) to 5 (very high) is suggested. A 0 impact rating is probably not relevant, because, in reality, some response and recovery costs will occur.

There may well be individual factors needed to clarify the importance of the overall impact. Such factors might include the likelihood of injury or death, the amount of property damage, the relative response and recovery costs, and, especially for a private organization, the costs of loss of business. The parameters of each factor and the scope of each rating can be adjusted as the planning proceeds, so that the results are realistic and meaningful. The purpose now is to show which threats are potentially the most dangerous or costly if they occur.

Threat assessment planning is best done using a worst-case scenario because this is what tends to happen in reality. Even then, actual costs incurred can be much higher than anticipated. Although the steering committee has only limited knowledge of possible costs, with some discussion the committee should reach consensus as to which rating to apply to each listed threat. Their findings will be reviewed later by others better able to predict potential costs and in a better position to realize that the impact of some threats is much more costly than others.

The impact costs are both direct and consequential. The direct expenses can include costs to locate the trouble(s), stabilize the systems, repair and install replacement infrastructure, reboot, restore databases, and thoroughly test both the data and the systems. Other direct costs that can result from each event include loss of productivity of system users, overtime needed to regain production schedules, temporary contracted services, and interim facilities, or materials and supplies needed immediately.

THREAT ASSESSMENT PROCESS 22 · 13

Add to these potential indirect or consequential costs such as food, sanitation, and lodging; loss of business, customers, or market share, falling stock price, and the costs of public relations efforts to control rumors and adverse news reports. Costs will continue to accumulate during the response phase, throughout the recovery period, and possibly long after. The totals can far exceed expectations.

22.3.6 Costs of Cascading Events. The possibilities and indeed the likelihood of cascading events have been described. Each event adds to the impact costs, but determining how likely is the cascading, and how much the extra cost, is at best an educated guess. If the causative events for each threat are shown on the threat list and the potential cascading events are there also, the impact costs of each individual event should be estimated. However, cascading may invalidate the likelihood factors. If one threat is already occurring, others may be imminent, regardless of the likelihood shown.

The total impact cost will probably be somewhat less than the sum of the parts. There may be some economies of scale during multiple events, as the response activity for one event can take care of another event for little extra cost. Here is where the experience and perspective of the steering committee, security professionals, and management become valuable. However, the combined costs must still be realistic, or the cost-value analysis described in Chapter 23 will be flawed.

22.3.7 Determine the Vulnerability to Each Threat. Deciding which threats are the most important can be subjective and controversial and based on unproven assumptions. The approach suggested here yields quantifiable data that are realistic and will clearly indicate which threats are the most important and why. Individual opinions will vary considerably. The steering committee reports can yield approximations that are statistically valid—if the rules are followed and everything is done carefully. The committee's deliberations should be reviewed by experts and senior management, using calculations to be described, followed by the cost-value analyses described in Chapter 23.

One other common alternative is a matrix showing likelihood vertically and impact horizontally, rating both factors as high, medium, or low. This yields nine levels of vulnerability for each threat ranging from high-high to low-low. Such an arrangement is neither useful nor realistic. Since allowance need to be made for those impact factors that are potentially more costly than others.

Once everyone settles on the likelihood and impact of each threat, the current vulnerability of each threat can be calculated. Fundamentally, vulnerability is a calculation combining the likelihood of an event with how much damage it might cause. The method suggested yields six levels of vulnerability, from 0 (none) to 5 (very high). The next two steps are oversimplified but can be adapted to most needs. Simply multiplying the likelihood times the impact will yield vulnerability values from 0 to 25, which would not reflect the relative impact costs of most threats. Neither does simple addition, as some models suggest.

1. It is best to use several impact cost factors, as described earlier, and then combine each factor using constant multipliers for each, so that the relative importance of each impact cost is preserved. Combination by averaging each cost factor may not yield useful results. Combining the mean value of each factor (or sometimes the highest value) can yield more realistic data.
2. Finally, convert vulnerability calculations to a 0-to-5-point relative scale by choosing a range of results for each scale value.

22 · 14 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

Working through this procedure for all threats will show that the formulas, multipliers, and ranges of scale values may need to be adjusted slightly to yield meaningful results. As already mentioned, simple multiplication and/or addition or otherwise linear formulae will not yield realistic results.

The importance of calculating vulnerability is that it can be done in real time, and the multipliers can also be adjusted in real time as threat levels change. For example, this can be reflected in the National Terrorism Advisory System (NTADS)⁶ by locale or by possible target types. It can also input warnings from local police and intelligence sources. (Some threat information is probably classified, but the multiplier it can provide is not restricted or even sensitive information.) All the data can be displayed on a spreadsheet, even to the point of color-coding each vulnerability level. Everyone involved in security or emergency management can have real-time access to the same screens, which will update continually as conditions change.

It is well to provide a choice of display screens summarized or in detail by regions and threat categories, or prioritized by vulnerability.

Finally, the complete vulnerability data should be restricted and available only to those with a need to know. This information is highly sensitive; it shows potential wrongdoers where the IS infrastructure is well defended and where it is *not*.

22.3.8 Completing the Threat Assessment Report. Once the detailed threat assessment is completed, a general report should be prepared for circulation to all stakeholders. The report is mostly textual and does not suggest weak spots or what the security defenses may be. Each steering committee member should sign off on this report, and so should the senior management and security staff. The purposes of this report are to develop security awareness, for use in training, and as evidence of due diligence that a thorough threat assessment was indeed performed.

More on how to evaluate and manage the vulnerabilities is described in Chapter 23 of this *Handbook*.

22.4 GENERAL THREATS. A wide array of threat situations can affect the IS infrastructure, degrade productivity, and cause anxiety that will reduce performance and morale. The list suggested earlier begins to identify some of these threats. Some other possible situations that are not generally associated with the IS infrastructure are listed in this section. However, the various threats suggested in this chapter are far from a complete list. Therefore, expert advice, local knowledge, and analysis of current events should all be utilized to customize the threat list to the particular needs of each jurisdiction and then to review and revise the list periodically.

As mentioned before, specific threats are usually divided into broad categories that can easily be summarized to provide situational awareness. Generally, threats are categorized as natural hazards, technical and man-made threats, civil unrest, vandalism, and attacks. Other less likely but potentially more serious possibilities include bombs, hazardous materials release, a pandemic, and threats relating to weapons of mass destruction. A final grouping could be local events that are specific to the sites being protected. The groupings themselves are not important as long as all possible scenarios are covered. Some allocations will be arbitrary, as some threats could be put into any of several categories. (But never list any specific threat more than once.) There is as yet no single, uniform, comprehensive threat list format. As the planning process continues, the best threat groupings will likely evolve.

GENERAL THREATS 22 · 15

22.4.1 Natural Hazards. Natural hazard events are becoming ever more frequent, damaging, and widespread, as their patterns seem to be changing. Increasingly recent events include major flooding, severe thunderstorms; hurricanes and tornadoes; blizzards, heavy snowfall, and ice storms; wildfires and heavy smoke; contamination of air, water, buildings, or soil; and the increasing threats of disease. Although earthquakes tend to occur in cycles, these too seem more prevalent currently. Any of these can disrupt business in any number of direct and consequential ways—many unanticipated—depending on the locations of information systems, networks, terminals, data storage, cables, and utilities. Natural hazard events cannot be prevented, but their impact can be mitigated. A closer look at each category follows.

- **Atmospheric hazards.** These can include severe weather events, such as tropical cyclones and hurricanes; severe thunderstorms with hurricane winds, strong lightning, and large hailstones; tornadoes; windstorms; blizzards and heavy snows; ice storms; air pollution and high ozone levels; nuclear fallout (which is always occurring, but at low levels); extreme cold weather; and extreme hot weather.
- **Geologic hazards.** These are mainly landslides and mudslides, land subsidence (sinkholes), and expansive soils due to water.
- **Hydrologic hazards.** These hazards include riverine flooding and flooding of low-lying areas due to heavy or prolonged rains; rapid ice or snowmelt; ice jams or debris obstructions in waterways; coastal damage due to storm surge; erosion of streams; and collapse of roadways, bridges, and buildings. Hydrologic events can cascade to disrupt water and sewage systems, cause food and fuel shortages, and even trigger fires or explosions.

A dam or levee failure is a major hazard that must be carefully considered and then reviewed with experts. Most dams are well maintained and protected from natural hazards but not from vandalism or terrorist attack. Nor are most reservoirs well protected from chemical attack. The consequences of such events can be far more devastating than imagined. The 2005 failure of levees in New Orleans caused by Hurricane Katrina produced catastrophic results over wide areas.

A prolonged drought can be particularly disruptive. There may be many fires and heavy smoke that require evacuation. Potable and cooling water may be in short supply. Both people and equipment must be protected from dust or smoke that cannot be done effectively within many facilities. Hydroelectric power may be rationed and rolling blackouts imposed.

- **Seismic hazards.** Such hazards include earthquakes and tsunamis. Many areas of the United States are at moderate or even high risk of a large earthquake. There were major events in these areas hundreds of years ago, and it is believed the time is nearing for a recurrence. For example, the Ossipee Fault in central New Hampshire produced a “San Francisco–size” earthquake about 300 years ago and may occur yet again, which would cause significant damage to Boston and its suburbs. At the least it would topple equipment cabinets, break cables and wires, and probably disrupt cooling and ventilation to IS systems.

Tsunamis, which are caused by undersea earthquakes, could at some time hit U.S. shores, based on historic events. And like Hurricane Katrina, a major seismic event anywhere in North America could impact businesses throughout a wide area.

22 · 16 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

22.4.2 Other Natural Hazards

- **Major volcanic eruptions.** Volcanic eruptions can spread atmospheric dust worldwide, which can affect weather and cause severe storms, damage cooling systems, disrupt radio and satellite transmissions, and restrict air travel.
- **Wildfire.** Fire can close transportation routes, disrupt utilities and support systems over a wide area, require evacuations, and create heavy smoke that is dangerous to people and equipment.
- **Blight or infestation.** Caused by disease, weather, or insects, blight or infestation can create health problems and disrupt food supplies, which can indirectly impact business.
- **Sunspot activity.** The approximately 11-year cycle of solar magnetic storms that cause sunspots may cause waves of electromagnetic radiation that can disrupt communications and even the electrical power grid.

22.4.3 Health Threats. Of increasing concern is the possibility of health emergencies, such as a SARS or an anthrax outbreak, or West Nile virus in relatively small areas. Such outbreaks can probably be contained with antiviral medicines and vaccines. However, business disruptions will surely result. Of greater concern is the threat of another pandemic that could quickly spread through large populated areas. As yet, there are no effective means to prevent either the onset or the spread of a pandemic, other than isolation and quarantine, which includes closing businesses and schools and perhaps many municipal offices as well. As mentioned earlier as a possible worst-case scenario, personnel refusing or unable to report to work during a major health emergency might total as much as 40 percent of the workforce, and they might remain away from the workplace for as long as 14 months. Consult local health officials how to access detailed information about these threats.

22.4.4 Man-Made Threats. Accidents cause most security problems. Accidents, sloppiness, and blunders are the consequences of bad design, poor quality control, improper installation, failure to update, and poor maintenance. Disruption often occurs during maintenance. Deliberate actions such as snooping, pranks, vandalism, and spying are all increasingly prevalent and sophisticated but are still overshadowed by accidental, unintended events.

Moving furniture or equipment (unless done by trained professionals) can damage wiring, connectors, and other equipment sufficiently to crash systems. And while it rarely happens, substantial “accidental” damage can occur during a labor dispute or when nonunion personnel are brought on site.

Construction work, alterations, repairs, and wiring changes often cause damage to information systems. Crews often drape drop cloths over workstations and equipment to keep off dust and debris. But no one thinks to shut down the equipment first, so it becomes overheated and probably fails either immediately or soon after. Crews also plug in power tools, floor waxing machines, or vacuum cleaners in whatever electrical outlets are handy. If these happen to be dedicated circuits for systems equipment, damage may well result. In like manner, workstation users often plug time stamps, electric staplers, refrigerators, fans, and immersion heaters into outlets intended only for information systems. Such mistakes can cause intermittent problems that are difficult to locate.

GENERAL THREATS 22 · 17

Wiring runs and exposed wire can also create threats. Wiring and cabling are vulnerable in many ways—and are becoming increasingly more so. Data cables tend to be fragile and easily damaged by accident, moving furniture, cleaning and maintenance, improper connections (such as unauthorized equipment), rage, or deliberate attack. Metal wires are also vulnerable to electrical and magnetic interference from nearby light fixtures, motors, transformers, or RF emitters, as will be explained. Lightning is attracted to, and electrical power surges may be induced onto, any metal wires.

Fiber-optic wiring, which is fast replacing metal connections high-speed communications, is especially fragile. Any unusual pressure or a sharp bend in any fiber cable can alter its transmission characteristics and may cause failure. However, fiber is not affected by any electrical or RF inference.

Cables, patch panels, cords, and connectors are often exposed to accidental damage during routine premises maintenance and cleaning. Vacuum cleaners and shampooing and waxing machines can easily damage both signal and power wiring. Carpet shampooing can soak connections and flood underfloor wiring. The casters on chairs and equipment or the moving of furniture often damages cables and connectors.

Intrusions are a major threat. There are many ways to gain access to the infrastructure of information systems. Once an intrusion is successful, repeated disruptions may follow due to accidents, mistakes, snooping, hacking, spying, vandalism, extortion, or deliberate attack. Intrusion is not usually prevented by premises access control. The infrastructure's vulnerability points are different, and early detection is necessary to prevent trouble before it happens.

Intrusions via wiring and cabling are covered in the prior section. Here the emphasis is on preventing access to hardware, distribution and termination panels, patch panels, or any of the utilities that support them. Assessing the possible threats first involves inspecting the existing infrastructure and the defenses already in place. It is also necessary to determine what physical threats and vulnerability points are not sufficiently covered by premises and logical security.

Unauthorized persons should not be allowed access to any type of equipment room, rack, or closet. Everyone granted access, including visitors, should be logged in and out. And as a second layer of protection for critical areas, a gatekeeper should be at hand to observe everyone who enters and leaves, to authenticate each person's identity, and to know why access is needed. Gatekeepers may include guards, receptionists, supervisors, or managers. Except where high security is required, the gatekeeping functions can be performed remotely using surveillance and access control systems, and monitored by motion and proximity detectors.

Workstations located outside of equipment rooms are generally protected logically—using login procedures, tokens such as smart cards, and biometric devices—depending on the sensitivity of the data the system handles. Theft of equipment, components, or removable media is possible, especially during nonbusiness hours. Substitution, rather than theft, is a much more serious threat. The only safe defense is that all data must always be fully encrypted. Should a theft occur, spare units should be immediately available from nearby secure storage. To save time, thieves often cut signal cables rather than disconnect them. It is therefore wise to keep spare cables in a nearby secure storage.

There are various theft-deterrent devices to clamp or tether equipment to nearby partitions or furniture. These devices can be cumbersome and intrusive, may hinder maintenance or repairs, and most are quickly defeated with bolt cutters or a pry bar. Intrusion alarms are more effective but only detect an attempted intrusion rather than deter theft. Intrusion alarms are usually silent and triggered whenever a case or cabinet

22 · 18 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

is opened. Software can trigger alarms also when a cable is disconnected or broken, or when equipment is removed or fails. Good, sturdy equipment and cabinet locks are more effective and can dissuade theft, provided the equipment and cabinets themselves are not easily removed.

22.4.5 Wiretaps. Wiretaps are another means of intrusion whose purpose is to copy data and sometimes also to change it surreptitiously. Most taps can be placed quickly and inconspicuously even in occupied office spaces. And most wiretaps are very hard to detect—if not impossible—except by close visual inspection of all possible entry points, which is time consuming and cumbersome, and must be repeated periodically. Very often, wiretaps are simply not noticed.

Taps are used for surveillance and spying, changing or erasing data, stealing software, theft of service, injecting worms or viruses, or planting decoys to trick responders into thinking they have found the true cause of an incident. Any physical access to the interconnect points, cable runs, servers, clients, or network equipment is a major vulnerability. An experienced person can place a tap within a couple of minutes, unobtrusively, even when under escort and closely watched.

A wiretap is basically a monitoring device. It may be a small box connected directly into a circuit to tap off and record its data. The data can then be retrieved manually or be transmitted by wire or radio signal to a remote location where the information can be permanently stored and analyzed. The monitoring device may be a splitter that can tap off some of the signals traveling through the circuit and reroute them to a removed location or to a small transmitter that is close by and well concealed. The better tapping equipment uses system power, not batteries, so its service life is unlimited. Monitoring the data obtained can occur within the building or outside, via the Internet, from a nearby vehicle, or via a telephone connection to anywhere in the world.

Wiretaps have long been illegal in the United States without a court order. Years ago, an illicit phone tap in New York City, for example, was patched to a leased telephone line to Mexico where the actual surveillance could take place legally. The patch connection was made by someone at a telephone company central office and was not likely to be discovered. But if it eventually was found, the network could be traced and disconnected. This procedure suggests the lengths and expense that a determined adversary will go to, and it also suggests the high value of the information that could be gleaned. This was years ago, long before today's easy access to the Internet, cheap long-distance calls, Wi-Fi, microwave, satellites, and other broadband systems.

Fiber optic circuits are difficult and expensive to tap and sometimes impossible without breaking the circuit or changing their parameters, which can be detected quickly. There is new high-tech gear available for optical taps. But even with this, good results are problematic.

The first method is to create a sharp bend in the cable and detect some of the light rays that may leak from the outside of the bend. The signal is at best very weak. Cable can be constructed to prevent this. Or cables can be run through metal conduit to make access and then creating a sharp bend very difficult, even at junction boxes. Last, a sharp bend will change the optical impedance of the cable run and will slightly degrade its performance, perhaps enough to effectively crash the circuit. The change in optical impedance can be measured, including the approximate location of the anomaly. The cable-bending method may work, but it is hard to accomplish, detectable, and readily observed by inspection. This method can monitor data, but it cannot inject data.

The only direct method to tap a fiber cable is to open a connection, quickly interpose an interface device, and then reconnect the circuit. Alarms should be triggered

GENERAL THREATS 22 · 19

immediately when a circuit breaks or otherwise fails, and these may also give some idea where the problem occurred. But with a good interface device and by working quickly, the tap can be inserted and conditions restored to look normal before anyone can respond. If this method succeeds, the intercept quality will be excellent, and data can be injected into the tapped line as well.

Spying on what is carried by fiber circuits may be possible from within equipment rooms or at other points where the optical data are converted to electronic form or to radio waves. The equipment to do this and its capabilities are highly classified and generally unknown to most security experts.

Wireless transmissions do not need to be tapped, but the process still applies. All that is needed is a very sensitive radio and a high-gain antenna, often situated in an outside vehicle. Chapter 33 in this *Handbook* discusses wireless-network security in detail.

Metal wires are relatively easy to tap into undetectably. Small, inconspicuous induction coils placed next to a wire can easily pick up the signals without penetrating the insulation. Sometimes the taps are done with tiny probes that reach around the circumference of a wire. Induction taps are generally undetectable, except visually. However, induction taps are difficult when multiconductor, twisted-pair cables are used, especially those where each pair is shielded and the entire cable is metal clad, as are some aerial telephone cables to protect them from radio-frequency interference.

When induction is not practical, bridged taps easily placed within an accessible junction box, cross-connect or patch panel, test box, or termination strip. Bridged taps can be employed a physical connection. On long cable runs, usually one or more inline terminal strips interconnect two cables. These are good locations for taps, especially where there may be unused wire pairs within the cable itself that can be used to route the tapped data to a safer location. Bridged wiretaps can also be inserted anywhere along a wire by penetrating the insulation with a needle or cutting away a tiny piece of the cladding and insulation in order to splice in a tap wire. Bridged taps are usually very high impedance, so they do not alter the circuit; therefore, like induction taps, they are undetectable, except visually. However, skill is required to place them to avoid damaging the tapped wires. As with fiber wiretaps, there are also high-tech methods of monitoring what metal wires carry, the details of which are classified.

Today, the biggest threat from wiretaps is that they are usually done by foreign intelligence, organized crime, or terrorist groups that are very well equipped, highly trained and experienced, and have almost unlimited funding to carry out their work. It is highly likely that most of their surveillance is not detected, located, or removed until long afterward—if ever.

22.4.6 High-Energy Radio-Frequency Threats. A controversial area of research involves high-energy radio-frequency (HERF) weapons, which some experts claim pose a serious threat to disrupt, damage, or destroy electrical and electronic systems and equipment over a very wide area. There are several scenarios.

A prototype HERF weapon demonstration at security conferences some years ago showed that a bulky apparatus constructed from easily obtained components can cause personal computers to malfunction while the device continues to operate. It can perhaps cause some system failures as well. There has been no more public information since then as to whether a smaller, portable version could emit a radio-frequency beam of sufficient power to disrupt systems dozens of meters away. Such a device is possible but may not be very practicable, except possibly from a large vehicle.

22 · 20 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

Along with this is another threat called high-intensity radio-frequency flux, which is simply interference from other electronic devices. These can be nearby transmitters, such as high-powered radio, TV, or radar systems, or transmitters in vehicles. Or there may be interference from close-by RF emitters including cell phones, computers, and devices designed for this purpose. Some of these possibilities are design issues, but when portable RF emitters are involved, the situations become infrastructure threats.

In 2004, another approach was reported that is potentially able to destroy unprotected electric and electronic system nationwide by way of an electromagnetic pulse attack (EMP). This report⁷ comes from a commission finding that a low-level nuclear explosion by a missile in the upper atmosphere could cause a massive electromagnetic pulse wave that could wreak catastrophic damage throughout most of the United States. Without extensive shielding, electrical and electronic systems could be destroyed, and power transmission systems as well. However, it appears that such an explosion must be done just right or the damage would be minimal.

The defense against all such RF threats is to shield equipment and wiring effectively. A Faraday cage can envelop and protect everything. This is usually one or more layers of copper screening that are well grounded. There also are Faraday bags to protect and store small items like cell phones or circuit boards.

However, once shielding is installed, it is wise to test that it actually works as expected. The author once toured a large, new data center of a major money-center bank. In touting the huge sums spent on extensive security protections for the new facility, someone mentioned that an impenetrable Faraday shield enveloped the entire room. The only problem was that a transistor radio on top of a cabinet was playing clearly and loudly. Had the room actually been well shielded, no broadcast signal could have penetrated.

Hazardous material incidents can include a release within a building or industrial plant, or from an aircraft, railroad, highway, or waterway accident. Many hazardous substances are routinely transported, stored, and processed throughout the United States. Many of these materials are extremely dangerous when released. They can quickly affect a wide area and may require immediate evacuation. In addition to accidental or deliberate release, these materials may be stolen to become part of an attack elsewhere. And those who attempt to steal, transport, and process the materials into a weapon may endanger others accidentally along the way.

Toxic threats include a buildup of radon within occupied space or in a water supply. Radon is released naturally by many types of rock. It is very prevalent and extremely dangerous if it can enter the body. Radon is an alpha particle that cannot be detected by most traditional measurement devices. Similarly, many toxic substances can accumulate in drinking water or building air. Any of these substances will be disruptive and cause anxiety far beyond the affected areas. Consult local health officials to identify and assess these threats.

Fire and smoke from remote events can be troublesome in several ways. Smoke from events hundreds of miles distant can cause health problems, delay supplies and deliveries, and cause food or fuel shortages. Workers may be unable to get to or from work, and may also be concerned about their families, friends, and homes. Electronic equipment is particularly susceptible to environmental damage from smoke or other airborne particles. And there may be utility problems or loss of cooling caused by a large, distant fire.

Nearby fire or smoke can be particularly troublesome. In addition to the problems just mentioned, an immediate system shutdown and personnel evacuation may be required. Power and data lines may be broken or damaged by heat or water. There may

WORKPLACE VIOLENCE AND TERRORISM 22 · 21

be flooding as well, and hazardous materials released that are potentially injurious to people and equipment. There will also be the need to contain the threats and to clean up the environment before full business can be resumed.

Smoke, dust, or other airborne particles can cause equipment failures, as these block cooling systems, clog filters, and build up inside equipment, constricting ventilation and convection cooling. Although these are mostly maintenance issues, they are also potential threats.

22.5 WORKPLACE VIOLENCE AND TERRORISM. The use of force, harassment, or physical violence is an increasing reality within the workplace. Drug- and alcohol-related incidents are on the upswing, as are workplace crimes necessary to support and conceal such habits. Adding to the threats is the increasing presence of rage within or upon the workplace. Any of these incidents can cause widespread trauma and disrupt business operations for months. Actual violence or terrorism, or simply the threat or the fear of this, can be extremely costly and disruptive.

Any violence situation—whether threatened, imagined, actual, or peripheral—can seriously disrupt information systems, whether the infrastructure is actually in danger or not. Performance, productivity, and morale will all plummet and remain low for a long time following any perceived or actual threat. Full recovery can take many months—assuming that more incidents do not occur in the meantime. Therefore, a safe working environment, good security planning and implementation, and security training exercises are essential so that everyone feels safe.

The likely tools of choice for inflicting widespread injury and damage may soon be weapons of mass destruction (WMDs) rather than the old-fashioned knives, guns, bombs, or arson that inflict only limited damage. WMD devices are far more dangerous, and many are small and easily concealed in a pocket, package, or briefcase.

A small, common-looking object containing biological toxins and powered by a flashlight battery can theoretically kill every person within the largest of office buildings. A vial no larger than a lipstick can contain enough virulent hemolytic viruses to kill every person within a 20- to 50-mile radius, if it is dispersed efficiently. Because few chemical or biological WMD compounds have much odor or color when they are disbursed, occupants, visitors, bystanders, and responders are all likely to be innocent victims. Illness can begin within minutes, hours, or days. Laboratory analysis is needed to identify many of the substances. More time is needed to determine the scope and spread of the damage. Ordinary personal protective gear and breathing apparatus can provide little or no protection.

WMDs can be enormously destructive. A national FEMA, Department of Justice training exercise called TopOff that I conducted in May 2000, simulating a biological attack on Denver, Colorado, resulted in an estimate of 57,000 fatalities. Since then, similar exercises have not released fatality estimates.

WMD threats can be grouped by the mnemonic B-NICE. They include *biological* agents, such as anthrax, cholera, pneumonic plague, tularemia, Q fever, Ebola, smallpox, botulism, ricin, and some others. These are all living bacteria whose incubation periods (the time of onset following exposure) are measured in hours or days.

Nuclear and radiological releases can cause widespread panic but are not likely to cause mass casualties, except to persons nearby. Radiation levels can easily be monitored, and building structures will often serve as an effective shelter.

Incendiary devices are utilized mainly to cause structural fires. A device can be planted surreptitiously and then triggered remotely or by a timer. Rockets and small

22 · 22 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

missiles are increasingly an incendiary threat, as are 9/11-type incidents: airplanes, vehicles, or boats used as bombs.

Explosive devices are similar. These may be stolen or smuggled ordinance or, increasingly, improvised explosive devices made from commonly available materials.

Combined attacks are possible, such as explosive devices used to disburse chemical agents (but not biologicals, which are live bacteria that would be killed by an explosion) or the so-called dirty bombs, which are explosives used to disburse radioactive material.

In a worst-case scenario, all WMD events can be exceedingly dangerous and disruptive, and any of these events may someday happen within the United States. Even the rumor of a WMD event can be damaging and cause widespread hysteria, panic, and create a large army of the walking well. Whether the threat is real or imagined, WMDs are costly to combat.

Most details on WMD and extremist or terrorist threats are classified, but state and federal authorities should be able to provide some insight for an effective threat assessment. The authorities should at least be asked about perceived threats, possible target areas, and regional and local incidents and security concerns.

22.6 OTHER THREAT SITUATIONS

22.6.1 Leaks, Temperature, and Humidity. Threats involving water and other liquids that may be hazardous should be considered also, as well as temperature and humidity conditions where equipment is located. Sprinkler systems in nearby spaces can cause equipment damage, as can liquid leaks from storage tanks, cooling towers, or pipes near equipment areas. Atmospheric conditions near equipment are threat situations as well. Air temperatures that are either too high or too low can cause equipment failure. High humidity can cause condensation within equipment and, worse, a form of galvanic action that degrades connectors that will then fail eventually. Also, low humidity is a threat because this promotes static electrical discharges that can be deadly to electronic gear and is often undetected until the equipment fails without warning. (See also Chapter 23, Section 23.8.7.)

22.6.2 Off-Hour Visitors. Cleaning and maintenance personnel usually work off hours and are often hired by a contractor or landlord who rarely provides much, if any, background checking, supervision, or training. Very few of the personnel are aware of security precautions and most know very little about the IS systems their work can damage. Many are paid poorly, are forced to rush their work, and may understand little English.

One of the classic disruptions can occur when workstations and associated devices are simply plugged into wall sockets or extension cords without labels indicating that they must not be unplugged; cleaning staff may easily unplug such devices in all innocence to run their own cleaning equipment and never even notice that they have abruptly powered down the computing equipment on nearby desks.

Waxing floors and shampooing carpets are usually done off hours by outside services. Moving furniture and changing workstations are often done after hours, as are repairs, alterations in occupied office space, and other major maintenance. Almost always, these people are unescorted, and many are not even logged into or out of the premises or identified in any way. Worse yet, many of these people prop open doors so they can work faster and sometimes so they can take advantage of air-conditioning in adjacent spaces.

OTHER THREAT SITUATIONS 22 · 23

Persons who are unknown often come into the premises without proper authorization. Some work for the landlord or the service organization. Some are delivering food. Some are messengers. And some are snooping, spying, looking to steal, or perhaps bent on violence. Even daytime workers may show up after hours. Therefore, the best security policy is to admit no one after hours until identified positively, logged in, and a need to be there is established. Everyone should be stopped at a reception or delivery desk, with no further access into the workplace until properly cleared. Anyone leaving should be logged out as well, especially if they have been out of sight of where they first entered. Screening visitors at any hour has the added security benefit that outsiders do not see where any IS infrastructure may be located or where money, wallets, or handbags are kept. A night bell or intercom outside an always-locked door should be used to prevent entry to a sensitive facility.

22.6.3 Cleaning and Maintenance Threats. A floor-waxing machine abrades everything it touches and will very quickly destroy unprotected wiring, connectors, dangling cords, or unseen power extension cords (which, by the way, are illegal according to most electrical codes). The waxing machine operator often cannot see any of these items or may not have time to look carefully at what is plainly visible. Carpet shampooing uses a lot of liquid and can flood floor-level outlet boxes and drain into underfloor ducts and conduits. Electrical plugs, receptacles, and unauthorized extension cords for critical equipment are often unlabeled. The cleaning staff does not know they are critical and can unknowingly unplug servers to power their cleaning equipment. Therefore, the threat assessment must first determine whether premises design invites problems, even though cleaning and maintenance personnel do their work carefully.

Users often compound these threats. Many workstations may be logged off but are left running continuously. No one tells users to shut down and cover their equipment before major cleaning or maintenance.

22.6.4 Storage-Room Threats. Rooms used to store computer supplies, paper, or forms are especially dangerous if a fire occurs. For example, stored cartons of paper, forms, or stationery expand when they burn, then burst into a conflagration that spreads quickly and burns at a very high temperature. Such a fire occurred in a high-rise office building in Manhattan. Even though the fire department quickly contained the fire, some steel building columns were so weakened by the heat that the building nearly collapsed. The cause of this fire was determined to be a cigarette butt that fell between stacked cartons of paper. This was assumed to have been accidental but could well have been arson.

Any storage room containing sensitive or inflammable materials must have smoke detectors, sprinklers, or other approved fire suppression systems designed to protect both the room and its contents. There must also be fire extinguishers nearby. Storage rooms should be kept locked, and access should be limited to trusted persons, if only to protect the value of the contents. There should be access control systems with admittance limited to authorized persons only, with open-door alarms provided as well. Delivery persons should always be continuously escorted. Loitering, smoking, drinking, drugs, snoozing, or any social activities must be kept out of and away from all storage areas.

22.6.5 Medical Emergencies. Most medical emergencies will cause business disruptions. People stop work to see what is happening, and if they know the victim

22 · 24 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

they may remain demoralized and unproductive for weeks afterward. Many will fixate on the circumstances and increasingly involve others to compound the problems. Even minor medical emergencies can cause large and long disruptions. In many medical emergencies, the first five minutes can decide life or death, and professional medical assistance is very likely unavailable that quickly. Any death that occurs in the workplace will result in long-term, widespread, and lasting trauma and disruption. (See also Chapter 23, Section 23.9.4.)

There must be first aid supplies, oxygen, and automatic electric defibrillators (AEDs) handy in every workplace. And there must be people nearby who are trained and currently certified in cardiovascular resuscitation and first aid. The costs of a disruption alone are far greater than those for providing ample medical supplies, equipment, and training. In addition, a first aid room and a trained nurse are wise precautions and probably a cost savings as well.

Prompt medical attention is essential for senior executives, visitors, and all workers and visitors. Medical threats will happen and will be very costly if medical assistance is not immediately available.

22.6.6 Illicit Workstation. A convenient method to set up a logical intrusion or attack is to unplug a desktop terminal or workstation that has limited functionality and substitute a full-featured machine well programmed with spyware and analytic utilities. Anyone with a full-featured notebook computer and physical access to the network may be able to connect easily.

Illicit users may then be able to log onto the network, possibly entering the user's own trusted password. They could then search for restricted information and use the full-featured machine to copy network data. Then the illicit user simply connects to a nearby telephone receptacle to export network data via a dial-up connection. (An Internet connection would likely be blocked by the firewall.) The illicit user may use hacking programs to gain supervisory status and spyware to access more sensitive information, crack passwords, steal software, or infect the network with malicious software. The illicit user may install a backdoor entrance into the network that an accomplice can use to spy, monitor network traffic, and modify or destroy data. An experienced user may then erase all evidence of any intrusion.

The intrusion might be done by an in-house employee, contractor, service technician, vendor, or consultant who could breach the network while ostensibly checking a user's machine or LAN connection. It could also be done within a hot zone using a Wi-Fi connection where there may be few, if any, firewall or security protections. Maintenance personnel working off hours also could substitute a terminal inconspicuously. However it is done, intrusion is a serious physical threat.

Security awareness is the best prevention. No one should be swapping equipment or connections unless a manager, supervisor, or nearby workers know the person's identity and what he or she is doing. If there is any possible doubt, the activity should be reported. The second best prevention is good network security, with alarms when any desktop systems is opened, disconnected, or shut down.

Chapter 15 in this *Handbook* includes detailed discussion of penetration techniques.

22.6.7 Other Local Threats. There are many threat situations specific to a group, organization, or community. Usually, many of these situations are identified and assessed as planning progresses. For example, for a community, important local

CONFIDENTIAL THREAT INFORMATION 22 · 25

threats might include vandalism or actual damage to school buildings, schoolyards, or school buses, a building or bridge collapse, interrupted power or energy transmissions systems, damaged fuel storage facilities, or a communications failure. Many possible specific situations will likely be discussed during the planning process, while review and advice by outside experts and government officials probably will lead to more important situations. Some considerations for local threat situations include:

- **Utility disruptions.** Disruptions happen more frequently than reported. Power outages can last a few minutes, hours, or days, caused by storm damage, equipment failure, tree branches hitting power lines, highway accidents that topple utility poles, digging accidents that break transmission lines, and vandalism or worse, such as toppling transmission line towers. There are increasing reports of transmission towers and utility poles deliberately toppled, power outages caused by insulators smashed or wires cut by bullets, and outages that occurred when wires and bus bars were stolen for the value of the metal in them. There can be static or RF interference or spikes on a utility line at any time due to atmospheric conditions, lightning, or the utility switching their feeders.
Communications outages can be even more problematic, and few of these situations are widely known. Communications providers are for-profit enterprises and cannot provide much fail-safe protection, backup, or redundancy when severe weather, vandalism, deliberate attacks, or equipment failure interrupt their services. Although there are not many alternate suppliers to choose from, the only security protections are redundancy and alternate suppliers.
- **Civil, political, and economic disruptions.** Such disruptions may be indicated by an elevated Homeland Security Threat Level or by state or regional alerts. Other possible disruptive events include a demonstration, march, disorderly group, or an unruly crowd. Other economic emergencies include a plant closing, a strike or a lockout, a transportation failure, or a shutdown. There are also the threats of violence of any kind; a hostage incident or kidnapping; sabotage of any infrastructure (e.g., power lines); contamination of food, water, air, or soil; a food or fuel shortage, a spike in energy costs or a shortage; and the repercussions of a major evacuation somewhere in the region.
- **Coordinated attacks.** These attacks are also possible, perhaps even by terrorists. Here, many points of the infrastructure are attacked at once, and many forms of attack may be used. There may be diversions in order to plant surveillance devices, place stronger weapons, or simply to distract response teams and stretch their resources more thinly. And whether spying is intended or not, the goal of a coordinated attack is to inflict maximum damage and disruption.
- **Solar storms.** Solar activity can also cause large problems. Sunspots affect electrical distribution systems as well as electronic systems, and they can severely disrupt satellite, microwave, and emergency communications, and radio, TV, and radar transmissions.

22.7 CONFIDENTIAL THREAT INFORMATION. Many other threat scenarios should not be described publicly because they are easily done by anyone with a grievance, and they tend to incite copycats. Other threats that will not be described utilize simple tools or devices that are readily obtained, are inconspicuous while in use, and can be safely hidden after a crime. Avoiding such threats is difficult, and

22 · 26 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

apprehension is unlikely. Nonetheless, mitigation is possible, and sometimes even deterrence, once these threats are known and understood.

Most vendors, installers, and consultants have long lists of such unmentionable threats and ways they know of to disrupt, snoop into, or destroy specific types of information systems. No threat assessment can be complete without asking all these information sources for whatever attack methods they can suggest as well as for methods of preventing or mitigating each threat.

Many useful sources of information and guidance are not generally public—many sources, indeed, that Internet search engines have not discovered. Some of this material is simply not publicly available, and other sources may be derived from classified documents. But briefings or redacted information may be made available to those who need it. More information may be available from local, state, and federal authorities, regulatory bodies, peer groups, business or vendor associations, professional groups, and security experts and consultants. Each of these sources may be willing to share information not available to the general public.⁸

The FBI is the ultimate source of threat information. But this information is mostly classified, and without the proper security clearance and a need to know, the FBI will not divulge much. A nationwide group of people involved in IS security is sponsored by the FBI and does provide some useful (sensitive but unclassified) information, once each member is vetted and approved. This group is called InfraGard, and there are chapters in most states. Visit www.infragard.net for a list of local chapters and contact them about membership.

22.8 SUMMARY. A vast array of possible physical threat situations can disrupt the infrastructure of information systems and thereby also disrupt business productivity and performance. The list of specific threats may identify several hundred situations, and each one should be included and considered during a threat assessment process that is well done and thorough.

Some of the threat situations are obvious and well known, while many others are much less so. Many are newly emerging, suggested by recent events worldwide and by reexamination of historic data. There also are issues of apathy, denial, and ignorance, where some persons feel that such things will never affect them or that nothing can be done to prevent a disaster. Both assertions are patently false and potentially very costly if pursued. Most threat situations can and will happen somewhere, someday. But whether they do or not is a matter of statistics and not conjecture.

Every possible threat situation can be mitigated to some extent by careful security planning and concerned management to minimize injury and damage. In fact, good physical security can be affordable, effective, and efficient. The first step is a thorough threat assessment to identify and consider all possible threat situations. The next step is to determine the likelihood that each threat may occur, the potential impact if it does, and the vulnerability of the organization within the context of its current security protections. Each step can be calculated statistically to determine the best possible options.

This chapter begins a threat assessment process that takes all these factors into account. Chapter 23 then suggests some ways to protect the information infrastructure and completes the threat assessment process with a cost-benefit approximation to determine the vulnerability of each threat.

These chapters are based on a uniform and comprehensive methodology recommended by the federal government and now required for all federal, state, and local

NOTES 22 · 27

agencies and departments. Compliance is strongly recommended for the private sector, if only as an effective means of risk management.

The value of this approach is that good security is a wise investment. Anything less than good security is a waste of time and money.

22.9 FURTHER READING

- Alibek, K., T. Dashiell, A. Dwyer, and S. Lane. *Chem-Bio Handbook*, 3rd ed. Jane's Information Group (ASIN B008DM4KT4), 2005.
- Baker, P. R., and D. J. Benny. *The Complete Guide to Physical Security*. Auerbach (ISBN 978-1420099638), 2012, 360 pp.
- BSI. *IT-Grundschutz Catalogs*. German Federal Office for Information Security (Bundesamt für Sicherheit in der Informationstechnik). English version, 2012. https://www.bsi.bund.de/EN/Topics/ITGrundschutz/ITGrundschutzHome/itgrundschatzhome_node.html
- DoT. *Emergency Response Guidebook* (2012). U. S. Department of Transportation, Pipeline and Hazardous Materials Safety Administration. A useful quick reference to the characteristic of hazardous materials, including some WMDs, with a table of isolation and initial-response guides. Free copies may be downloaded from www.phmsa.dot.gov/hazmat/library/erg
- Fennelly, L. *Effective Physical Security*, 4th ed. Butterworth-Heinemann (ISBN 978-0124158924), 2012, 384 pp. FEMA. *Multi-Hazard Identification & Risk Assessment (MHIRA)*. U.S. Federal Emergency Management Agency. Publication 9-0350 released in 1997. This is an excellent reference manual, 355 pages. All types of natural hazards are described in detail, and maps show the likelihood of each threat throughout the country. See www.fema.gov/library/viewRecord.do?id=2214
- FEMA. *State and Local Mitigation Planning, How-to Guide*. U.S. Federal Emergency Management Agency. There are nine manuals, FEMA Publications 386-1 through 386-9. www.fema.gov/hazard-mitigation-planning-resources#1

22.10 NOTES

1. See, for example, Duncan Chappell and Vittorio di Martino, *Violence at Work*, 3rd ed. (Geneva: International Labour Organization, 2006). Review available from www.workplaceviolence911.com/docs/20060624.htm
2. The International Labour Organization of the United Nations provides a good view of the rise of workplace violence worldwide. For example, look at the paper at www.ilo.org/public/english/protection/safework/violence/violwk/violwk.htm
3. www.cesnur.org/testi/FBI_004.htm. Every U.S. police department received a copy of the Project Megiddo report and may provide access.
4. The National Institute of Standards and Technology (NIST) within the Department of Commerce publishes many standards and procedures that relate to computer security. The current list, as of November, 2013, at csrc.nist.gov/publications/index.html
5. EMAP, P.O. Box 11910, Lexington, KY 40578, Tel: 859-244-8222; www.emaponline.org
6. See www.dhs.gov/national-terrorism-advisory-system

22 · 28 PHYSICAL THREATS TO THE INFORMATION INFRASTRUCTURE

7. Report of the Commission to Assess the Threat to the United States from Electromagnetic Pulse Attack. Available at www.globalsecurity.org/wmd/library/congress/2004_r/04-07-22emp.pdf
8. Multi-State Information Sharing and Analysis Centers (www.msisac.org/about). Regional Information Sharing System Automated Trusted Information Exchange (www.riss.net). DHS Daily Open Source Infrastructure Report (www.dhs.gov/dhs-daily-open-source-infrastructure-report)

Computer Security Handbook, Sixth Edition, Volume I
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

INTRODUCTION TO PART III

PREVENTION: TECHNICAL DEFENSES

The threats and vulnerabilities described in Part II can be met in part by effective use of technical countermeasures.

The chapter titles and topics in this part include:

23. **Protecting the Physical Information Infrastructure.** Facilities security and emergency management
24. **Operating System Security.** Fundamentals of operating-systems security, including security kernels, privilege levels, access control lists, and memory partitions
25. **Local Area Networks.** Security for local area networks, including principles and platform-specific tools
26. **Gateway Security Devices.** Effective recommendations for implementing firewalls and proxy servers
27. **Intrusion Detection and Intrusion Prevention Devices.** Critical elements of security management for measuring attack frequencies outside and inside the perimeter and for reducing successful penetrations
28. **Identification and Authentication.** What one knows, what one has, what one is, and what one does
29. **Biometric Authentication.** Special focus on who one is and what one does as markers of identity
30. **E-Commerce and Web Server Safeguards.** Technological and legal measures underlying secure e-commerce and a systematic approach to developing and implementing security services
31. **Web Monitoring and Content Filtering.** Tools for security management within the perimeter
32. **Virtual Private Networks and Secure Remote Access.** Encrypted channels (virtual private networks) for secure communication, and approaches for safe remote access

III · 2 PREVENTION: TECHNICAL DEFENSES

33. **802.11 Wireless LAN Security.** Protecting increasingly pervasive wireless networks as a majority of the world's Internet users access information through their mobile phones and tablets
34. **Securing VoIP.** Security measures for Voice over IP telephony, a cost-effective method for increasing long-distance collaboration
35. **Securing P2P, IM, SMS, and Collaboration Tools.** Securing collaboration tools such as peer-to-peer networks, instant messaging, text messaging services, and other mechanisms to reduce physical travel, and to facilitate communications
36. **Securing Stored Data.** Managing encryption and efficient storage of stored data, including cloud-storage services
37. **PKI and Certificate Authorities.** Concepts, terminology, and applications of the Public Key Infrastructure for asymmetric encryption
38. **Writing Secure Code.** Guidelines for writing robust program code that includes few bugs and that can successfully resist deliberate attacks
39. **Software Development and Quality Assurance.** Using quality assurance and testing to underpin security in the development phase of programs
40. **Managing Software Patches and Vulnerabilities.** Rational deployment of software patches
41. **Antivirus Technology.** Methods for fighting malicious code
42. **Protecting Digital Rights: Technical Approaches.** Methods for safeguarding intellectual property, such as programs, music, and video, that must by its nature be shared to be useful

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 23

PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

Franklin Platt

23.1	INTRODUCTION	23·2	23.5.3	Overt, Covert, and Deceptive Protections	23·10
23.2	SECURITY PLANNING AND MANAGEMENT	23·2	23.6	ACCESS CONTROL	23·12
23.3	STRATEGIC PLANNING PROCESS	23·3	23.6.1	Locks and Hardware	23·12
23.3.1	Attractive Targets	23·4	23.6.2	Card Entry Systems	23·13
23.3.2	Defensive Strategies	23·4	23.6.3	Proximity and Touch Cards	23·13
23.3.3	Who Is Responsible?	23·5	23.6.4	Authentication	23·14
23.3.4	One Process, One Language	23·5	23.6.5	Integrated Access Systems	23·15
			23.6.6	Portal Machines	23·16
			23.6.7	Bypass Mechanism	23·16
23.4	ELEMENTS OF GOOD PROTECTION	23·5	23.6.8	Intrusion Alarms	23·17
23.4.1	Segmented Secrets	23·5	23.6.9	Other Important Alarms	23·18
23.4.2	Confidential Design Details	23·6	23.7	SURVEILLANCE SYSTEMS	23·18
23.4.3	Difficulties in Protecting the Infrastructure	23·7	23.7.1	Surveillance Cameras	23·19
23.4.4	Appearance of Good Security	23·7	23.7.2	Camera Locations and Mounts	23·20
23.4.5	Proper Labeling	23·8	23.7.3	Recording Systems	23·20
23.4.6	Reliability and Redundancy	23·8	23.7.4	Camera Control Systems	23·20
23.4.7	Proper Installation and Maintenance	23·9	23.7.5	Broadband Connections	23·21
23.5	OTHER CONSIDERATIONS	23·10			
23.5.1	Equipment Cabinets	23·10	23.8	FACILITIES DESIGN	23·21
23.5.2	Good Housekeeping Practices	23·10	23.8.1	Choosing Safe Sites	23·22
			23.8.2	Physical Access	23·23

23 · 2 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

23.8.3	Protective Construction	23·24	23.9.4	Medical Emergencies	23·43
23.8.4	Using Existing Premises Alarms	23·27			
23.8.5	Clean Electrical Power	23·28	23.10	COMPLETING THE SECURITY PLANNING PROCESS	23·43
23.8.6	Emergency Power	23·30	23.10.1	All-Hazard Mitigation Plan	23·43
23.8.7	Environmental Control	23·36	23.10.2	Cost-Benefit Analysis	23·44
23.8.8	Smoke and Fire Protection	23·38	23.10.3	Security Response Plan	23·45
23.9	MITIGATING SPECIFIC THREATS	23·41	23.10.4	Implementation, Accountability, and Follow-Up	23·46
23.9.1	Preventing Wiretaps and Bugs	23·41			
23.9.2	Remote Spying Devices	23·42	23.11	FURTHER READING	23·46
23.9.3	Bombs, Threats, Violence, and Attacks	23·42	23.12	NOTES	23·47

23.1 INTRODUCTION. Chapter 22 in this *Handbook* describes a systematic approach to identifying threats to the physical information infrastructure. This chapter describes practical recommendations for implementing good physical-security protection for information systems (IS).

There was a time when insurance was enough to cover most threat situations. However, today disaster, workplace violence, and terrorism insurance coverage may not be available or affordable without an independent, outside security audit to ascertain a good security program is in place.

There is one nearly ironclad legal defense against negligence allegations, and this is to demonstrate that the organization is following the commonly accepted federal security planning and management procedures. Usually, such compliance alone is sufficient to prevent negligence suits from even being filed; any remaining allegations are usually minor and easily covered by insurance. Conformity to the commonly accepted federal procedures can quickly be ascertained by a recent security audit.

The security audit will examine compliance with all documented facets of the planning, preparation, policies, training and exercises, management and oversight, and response and recovery plans and procedures. These are all the necessary components for protecting the information infrastructure.

Chapter 54 in this *Handbook* discusses audits, and Chapter 60 discusses insurance for information systems.

23.2 SECURITY PLANNING AND MANAGEMENT. In the United States, there are many laws, regulations, and federal directives that should be incorporated into the security planning and management of any organization. Even when compliance is not yet mandated, it is highly recommended simply to avoid unnecessary and potentially huge costs of defending against allegations of negligence, and establishing in court that the security procedures in place are sufficient.

Some of the new requirements that may affect almost any organization in the United States are briefly summarized in the following list.

National Incident Management System (NIMS) Compliance: The National Incident Management System (NIMS) defined by the Federal Emergency

STRATEGIC PLANNING PROCESS 23 · 3

Management Agency (FEMA) requires all departments and agencies of the federal, state, and local governments to implement its standards. NIMS is an emergency-response template which supplements and unifies the incident command systems long in use by most response agencies.¹

National Response Framework: The National Response Framework (NRF) provides an emergency operations plan which all state and local governments have long been required to have. The NRF now unifies the content and format for all such plans. Compliance also requires certification of senior officials.²

National Infrastructure Protection Plan: The National Infrastructure Protection Plan (NIPP) defines recommendations for the protection of critical infrastructure and key resources. It proposes a unifying structure that includes partnership initiatives, a risk-management framework, and information sharing.³

Robert T. Stafford Disaster Relief and Emergency Assistance Act of 1988 (Public Law 100-707): In combination with the Disaster Mitigation Act of 2000 (Public Law 106-390), every state and community must have a current and FEMA-approved *All-Hazard Mitigation Plan* and a current and approved *Emergency Operations Plan* in place.⁴

The Sarbanes-Oxley Act of 2002 (SOX): Public corporations must report all of their internal financial controls, which may include those relating to security management—including future events that can materially affect future earnings. Most of the threats identified in Chapter 22 can do this, so that the possible mitigated costs of each threat as determined at the end of this chapter may need to be included in the SOX reporting.⁵

The Health Insurance Portability and Accountability Act of 2002 (HIPAA): All organizations that receive and hold personally identifiable individual health information must keep those data.⁶

The Gramm-Leach-Bliley Act of 1999: Financial institutions, lenders, advisors, accountants, and businesses that process or receive financial information must protect those data against unauthorized disclosure and modification.⁷

Emergency Management Accreditation Program (EMAP): EMAP, “an independent nonprofit organization, is a standard-based voluntary assessment and peer review accreditation process for government programs responsible for coordinating prevention, mitigation, preparedness, response, and recovery activities for natural and human-caused disasters. Accreditation is based on compliance with collaboratively developed national standards, the *Emergency Management Standard* by EMAP. Accreditation is open to all U.S. states, territories, and local government emergency management programs. Anyone can subscribe to receive standards and guidance materials.”⁸

23.3 STRATEGIC PLANNING PROCESS. The planning process must think ahead strategically to determine the best possible security options to add maximum value and contribute to profits, enhance productivity and performance, and avoid intrusiveness. The process of identifying all possible threat situations and a statistical method to predict the likelihood of injury or damage and the current vulnerability of the enterprise to each individual threat is described in Chapter 22. This chapter builds on the planning process, suggests ways to better protect the infrastructure and thereby reduce its vulnerability, and describes how to implement the security planning and management process.

23 · 4 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

Some facets of good physical protection to consider in the planning process are outlined next. These suggest some critical elements and how to protect them, and some possible weak points in the information infrastructure. Section 23.11 completes the planning and mitigation process.

23.3.1 Attractive Targets. Much of this chapter deals with what amounts to physical hacking. As with its electronic counterpart, physical attacks cannot be predicted as to when, where, or how they will occur. Electronic hacking from within the organization, over its networks, or from anywhere in the world via the Internet is still the best way to break into most information systems. However, physical intrusion can be harder to detect and locate, and often is more damaging and costly. In addition, threats from employees or contractors who are bad actors with *authorization* to be unsupervised in their work with or near the physical infrastructure of information systems are a serious threat.

Chapters 12 and 13 in this *Handbook* discuss the psychology of computer criminals and dangerous insiders.

23.3.2 Defensive Strategies. There are many effective defensive strategies to consider. The planning process must evaluate each approach and how best to combine them strategically to maximize effectiveness and minimize costs. These defensive strategies are common:

- *Prevention* so that specific threats do not affect the enterprise.
- *Deterrence* so that specific threats are not likely to occur.
- *Mitigation* to reduce each threat to tolerable consequences.
- *Redundancies* so there are no critical links in the infrastructure that cannot be bypassed. There are many methods of redundancy, such as multiple data paths; bidirectional data loops; parallel or distributed processing; alternative support systems and utilities; and many more.
- *Early warning* to detect impending trouble and delay onset, so that fast response can prevent or minimize any disruption.
- *Layers of security*, which are like the concentric layers of an onion, so that several layers of security must be penetrated before a target can be reached. This adds reliability, because a failure or breach of one layer does not compromise the other concentric layers.
- *Insurance* that can reimburse some of the recovery costs but usually few of the response costs. Insurance coverage often excludes many threats and can be costly. Insurance may not cover (or even be available for) gross negligence, some acts of God, flooding, terrorism, or acts of war.
- *Capital markets*, which are less costly than insurance and better able to lay off larger and broader risks.
- *Self-insurance* to establish retentions (which are funds accumulated for the purpose) in the hope nothing serious ever happens.
- *Contract security services* that are performed in-house, which basically transfer risks, but do not necessarily mitigate threats.
- *Outsourcing*, which is another option that introduces still other and often unrealized threats and vulnerabilities.

ELEMENTS OF GOOD PROTECTION 23 · 5

23.3.3 Who Is Responsible? Effective security requires both governance from the boardroom and oversight by senior management. Top management must also actively sponsor it and insist that everyone involved understands, supports, and respects the security plans and procedures. Accountability and oversight are essential, as is insistence on periodic security exercises, review, and updating.

If possible, the infrastructure-security manager should *not* also be in charge of IS or corporate or premises security. Although managers of these areas face similar threats and their groups must work closely together, their priorities and levels of response to any incident are different. It is best therefore that two or more specialized security groups report to one senior executive officer with clear lines of authority. Even if the role cannot be full time, defining a formal position for an infrastructure-security manager is a wise investment.

23.3.4 One Process, One Language. As mentioned previously, the security planning and management process must be uniform and comprehensive, and applied effectively and efficiently. There should also be a current audit done by an outside, independent group to ascertain both security preparedness and compliance with the pertinent laws, codes, and regulations. The least costly and most effective approach is to follow federal guidelines. The Department of Homeland Security (DHS) and FEMA procedures provide a uniform and comprehensive methodology and generally accepted standard practices that are uniform, comprehensive, and most likely to deter liability as well. This process is recommended for any organization, public or private.

FEMA makes the distinction between an emergency, which is a situation that an organization can handle with its own resources, and a disaster, which is when internal resources are likely to be overwhelmed and outside help is needed. The distinction is not always clear-cut. “Crisis” is not a defined term. FEMA also recommends planning and preparation for worse-case situations, which often turn out to be what happens. Many threats can become disasters to a business, even when outside resources respond quickly. For example, the consequences of an equipment room fire can be disastrous when the local fire department responds with axes and water hoses, cuts off all electrical power to the building, and smashes out windows to vent smoke. Even a major incident that occurs away from the immediate premises can necessitate an evacuation, cutting off building power, and creating serious disruption.

23.4 ELEMENTS OF GOOD PROTECTION. Protecting the infrastructure generally requires different and stronger defenses than premises security or IS logical security can provide. The infrastructure protection must be effective, efficient, and affordable. Yet it must also be nonintrusive and user friendly. Too much protection is unnecessarily costly and often counterproductive. Too little can be even more costly and endanger productivity, morale, and goodwill.

Some of the requisite elements for just the right amount of protection are provided in the sections that follow.

23.4.1 Segmented Secrets. Although *security by obscurity* is rightly derided if it is the primary focus of protection, to maintain good security, no *one* person should know or control *all* the details or inner workings of the security systems and procedures. If total understanding of the security systems is segmented into several parts and the knowledge distributed among collaborators, there is much less likelihood of misuse, fraud, or error, and less dependency on a few key persons. The shared knowledge contributes to separation of duties and pervasive vigilance.

23 · 6 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

However, the more people with knowledge of the security systems and procedures, the more these systems become vulnerable. It is easier to compromise security through coercion or extortion, or by an unwise remark inadvertently disclosing information. The list of those with inside knowledge may include managers, administrators, maintenance personnel, users, partners, suppliers, customers, vendors, and consultants. Although many individuals must know at least some of the security protections, no one needs to know all of the details.

Beyond segmenting, there is another precaution needed for good protection. If members of the same group must share sensitive knowledge, a “two-person” rule should apply—an example of separation of duties. This principle is especially important when anyone is able to modify the security infrastructure, alarms, or event logs. The two-person rule says that any modification to the security system requires two authorized persons working together and that there will always be an audit trail showing who did what and when. This procedure also is used when two or more groups share responsibility for common elements of a security system.

23.4.2 Confidential Design Details. It is often quite easy for others to locate and identify critical infrastructure components. Many times the information is clearly shown on public documents. The signage on infrastructure locations often clearly describes what is inside, and critical components are often put in spaces easily accessible to outsiders. These are all unsafe design practices.

Many types of documents clearly indicate IS areas, support systems, and other infrastructure. The documents may show the locations of equipment, wiring, utilities, and cable runs. Documents that are likely to reveal sensitive information include building plans that show floor and office layouts, furniture, wiring, and equipment locations. Architectural and engineering documents, drawings, and specifications often show sensitive information. These documents often list the function of each area, and even an occupant’s name or title. Other types of documents breach security as well. Examples may be as-built plans, alteration and construction plans, electrical and communications wiring diagrams, patch-panel and cross-connect setups, as well as contract drawings or proposals, shop drawings, installation plans, maintenance diagrams, and, especially, documents filed with code-compliance and regulatory agencies. Good protection requires that all such documents be controlled and kept securely stored. Better yet, sensitive information should be removed and alphanumeric designations that are cross-referenced to sheets than can easily be kept classified should be shown.

All of the listed documents are routinely distributed to a wide range of sources, such as interested contractors and bidders, vendors, suppliers, and maintenance providers. Building managers, landlords, and often real estate offices are likely to keep copies on file, and many other persons can readily access the documents and copy them. Most such documents are publicly available or obtainable via the sunshine laws of each state, by court order, or simply by deception. Moreover, many legitimate persons obtaining or receiving the documents have no internal document control or security provisions.

All room and area designations should always be alphanumerical. Descriptive or functional names should not be used for any area. This caveat pertains to the entire premises, all public areas of the facility, and all building mechanical and core areas. Functional designations or terms such as “treasurer,” “marketing director,” “security desk,” “computer room,” “network closet,” or “telephone room” should never appear. Nor should the names of any occupants, departments or functional groups, or individual tenants’ spaces be given. Use only alphanumeric designations.

ELEMENTS OF GOOD PROTECTION 23 · 7

Lists that correlate area and room numbers, functional areas, or any descriptive names should be kept secure in a locked file, as should equipment room drawings, patch panel connections, and wiring plans. All these must be readily available to system administrators in the event of a system failure or when doing maintenance or upgrades, but only on a need-to-know basis. Ideally, even managers do not have access to this information, except when accompanied by security personnel. As well as controlling access to such documents, it is equally vital that every document be kept current.

Finally, security personnel should review all plans, drawings, and documents for construction, alterations, equipment moves, or any other physical changes before the information is issued. Security personnel should also review the invitations to bid and all the drawings and specifications, review again the quotations received, and continue to review every as-built document produced throughout a project. Internal security should review all these as a matter of policy. If the security personnel do not have sufficient time or expertise to review everything quickly and thoroughly, outside independent experts can be valuable. When properly chosen, these experts can provide broad experience, perspective, and evidence of due diligence, so that management cannot later be accused of negligence in protecting information, people, or property. Finally, the security leader should sign-off that no project documents violate any internal security policies.

23.4.3 Difficulties in Protecting the Infrastructure. Decades ago, most IS equipment was housed inside a single computer room, and most of the external terminals and peripheral equipment were located nearby. Cable runs were short and generally within secure areas, and access controls, alarms, and surveillance could easily cover the critical areas. There were often security guards as well. But in those days, there were fewer threat situations and less cost if trouble did come. Many organizations then were apathetic and security was lax. They were lucky, but a few incidents did occur and often at great cost. There were also likely many undiscovered security breaches and still more that were never reported.

Today, the IS infrastructure is much larger and more complex, and the potential costs of trouble are far greater. Today's infrastructure is much more interdependent, and it now includes many more equipment and network rooms. It extends to many more telephone and utility closets, and interconnects widespread and diverse desktop, peripheral, and remote nodes. Today's infrastructure is increasingly harder to protect, and the future outlook suggests many more threat situations with the potential to cause major business disruption.

To further complicate security, there are now many more and diverse IS interfaces and a complex infrastructure to protect. Interfaces now include direct and switched wiring, wireless topologies, and infrared coupling. Access to the Internet, local area networks, and wide area networks may now utilize combinations of metal wiring and fiber optics, wireless, satellite, TV cable, and microwave links. Some of these interfaces will be dedicated, others switched, and still others temporarily patched. It will be difficult to even locate all the interfaces, yet each must be protected, as must be their cable runs and the utilities that support them. Today, the early warnings from the security alarms and defenses must be so effective that most trouble can be prevented before it happens.

23.4.4 Appearance of Good Security. The appearance of an armed fortress is usually counterproductive. This usually intimidates and obstructs both visitors and staff more than it protects them. The same is often true for too many guards and

23 · 8 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

receptionists stationed to block entry points to internal areas (unless this is deemed absolutely necessary). Such barriers tend to enrage anyone who may be already anxious, and can provoke attacks. Barriers tend to be ineffective; for example, even if the glass is bulletproof and spray-proof, the pass-through holes likely are not. Most barriers can be breached as well, and the person(s) inside injured or at least traumatized. For example, some years ago most banks rushed to put their tellers behind thick glass panels. However, they soon discovered that the glass provoked trouble rather than deterring it. The banks quickly removed both the glass and the pass-throughs, which many businesses then bought and installed, with similar results. Anger and potential violence are best avoided by good facilities design, security systems and access controls, and training in security awareness and violence prevention.

There are also considerations whether security devices should be covert or appear in plain sight for all to see. Defenses that are in plain sight can serve as deterrents and promote a feeling of safety, but they can also be vulnerable themselves. Cameras can be spray-painted, shot out, or knocked aside with a club. Any exposed wiring (which is a security no-no in itself) can be cut or shot through. The alternative is to use concealed devices and also dummy devices or “honeypots” that are obviously positioned. (Both are discussed later.)

23.4.5 Proper Labeling. Good security requires quickly locating trouble spots, with certainty that the nomenclature of every component, cable, and connector is clear and current. There must be proper and consistent labeling of all data and power cables, both inside and beyond the secure areas. No label or tagging should reveal confidential information. Many vendors, installers, contract personnel, and in-house staff tend to use their own labeling and tagging systems, some of which are clear and understandable, while others are not. The labeling and tagging must match the documentation, plans, and drawings as well. Sloppy wiring management is often the norm as changes are made but not marked or documented, or wiring abandoned and not removed. The status of in-house personnel, installers, and maintainers is likely to change frequently, creating all the more potential for misinformation and confusion. The inevitable result is poor protection and slow response to a security incident. There are generally accepted wire-management procedures, and a single method should be utilized throughout the information infrastructure. It is especially important that all authorized personnel understand and accept the labeling system.

The Telecommunications Industry Association/Electronic Industry Alliance (TIA/EIA) publishes the generally accepted authority, TIA/EIA Standard 606-B, *Administration Standard for the Telecommunications Infrastructure*,⁹ which describes labeling of cables, connectors, patch panels, cross-connections, and equipment. This standard also requires labeling firestops, grounds, special-ground circuits, and neutral wires. The *National Electric Code* (NFPA 70),¹⁰ published by the National Fire Protection Association, also includes standards for labeling cables and conduits. Local codes may augment the national codes. Generally, any substantial alterations, moves, or changes, and usually all new construction will require full compliance with the current standards throughout the premises—which is a wise security investment as well.

23.4.6 Reliability and Redundancy. The first requisite of reliable system performance is reliable equipment, systems, and infrastructure, properly installed and maintained. Unfortunately, some components may be poorly designed, subject to erratic quality control, damaged in transit, or applied wrong. Much equipment now includes a fail-over mode to maintain full performance when a failure occurs.

ELEMENTS OF GOOD PROTECTION 23 · 9

Another component of high service levels is *redundancy*: concurrently accessible equipment and parallel paths to take the load should one component falter. Effective redundancy for business continuity also requires alternate sites that are off premises to process and store information. It is best that there be multiple alternate sites, each within a safe environment and well distant so that problems at the primary site will not affect any alternates.

Inside each computer room, storage systems should be redundant arrays of independent disks (RAID) allowing important data to be fully replicated. For critical data, the RAID storage systems should themselves be redundant. Servers should also be redundant. Multiple parallel servers with load balancing are a wise investment, so that one server will automatically take over another's load if it falters or must go off-line for any reason. The load-balancing features also provide potentially higher throughput for the overall system by allowing parallel computations and input/output. All critical equipment should have redundant power supplies, fans, and hard drives that can be diagnosed quickly and hot-swapped easily.

Good system manageability is another vital requirement. This includes hardware that can detect trouble before it happens and, if possible, pinpoint the exact location. It also includes good management software with warning and alert capabilities and good logging systems. Remote management capabilities must also be private and secure to preclude penetration or denial-of-service attacks. Good security management also requires that any changes or disabling of alarm parameters should require two authorized persons to be physically present inside the equipment room and simultaneously logged on. Good system management adds some cost but is a wise investment in effective security and oversight.

23.4.7 Proper Installation and Maintenance. Good protection requires that all information systems, equipment, and wiring be installed properly, according to the manufacturer's instructions and the intended usage. All the wiring must conform to, or exceed, local code requirements. Data wiring should be independently tested and certified that it meets specifications, current standards, and, if possible, anticipated future needs.

Out-of-the-box equipment hookups and installations are common and the cause of many system failures. Most security features are disabled when components and software are shipped. Proper installation requires careful setup, customization, and performance testing, for which adequate time and resources must be allocated. Promptly installing the latest modifications, service packs, updates, and security patches is also vital to maintaining performance. Delays of days, weeks, and even months often intervene, while new threats emerge or threat levels escalate until the new defenses are in place. Once installed, the information systems and infrastructure must be periodically reviewed, tested, and kept up to date. Frequencies of such reviews and testing must reflect the known or estimated rates of failure and the consequences of such failures.

Administrators, installation, and maintenance personnel must be properly trained, experienced, and, in many cases, certified or licensed as well. Each person's credentials should be checked before being permitted on site. Given the limited staff, time, and budgets available, there is often more lip service than actual certainty in the process of reliability assurance. Management must understand that proper installation, upkeep, and maintenance together constitute cheap and effective assurance that IS performance is never compromised.

For more details of data backup methodologies, see Chapter 57 in this *Handbook*; for business-continuity planning and disaster recovery, see Chapters 58 and 59.

23 · 10 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

23.5 OTHER CONSIDERATIONS. There are many other factors that can significantly establish and improve physical security. A few of these are discussed in the next sections.

23.5.1 Equipment Cabinets. Most IS equipment is now open-rack mounted to save floor space, for more reliable performance, and for easier access. Although enclosed equipment cabinets cost more than racks, they offer much better protection. Equipment mounted within a closed cabinet can be better ventilated and cooled, kept free of contaminants, and may also escape damage from external particulates, water, smoke, or liquids. Should there ever be overheating or smoke generated within an enclosed cabinet, the condition can be detected quickly and usually resolved before equipment or wiring are damaged. Locks keep out those unauthorized, or at least delay access and leave evidence of trouble. Open-door alarms provide another strong layer of protection, and wiring conduits to cabinets can also be protected. Cabinets with redundant fans can better monitor and maintain ventilation and cooling, which, in turn, facilitates more equipment mounted in less floor space. In all, closed cabinets can provide an additional layer of protection against accidental or deliberate damage to the infrastructure.

23.5.2 Good Housekeeping Practices. All food and drink must be kept out of equipment rooms, since they can cause considerable damage if spilled on equipment, a monitor, keyboard, connectors, and wiring and cable harnesses. Food also attracts insects and rodents (which are found in many buildings), many of which also like to eat wiring. In addition, grease and other contaminants from food can easily be spread by operators onto critical surfaces of equipment such as input output devices and cause read/write errors or even complete shutdown of the devices. For example, a single greasy DVD can ruin a DVD unit or high-capacity backup-tape unit. Worse, a contaminated read/write head or entry slot can in turn spread contamination to other media and then on to other devices. Space for food and drink should be provided outside the equipment room, where routine maintenance personnel can keep the food area clean. A convenient dispenser for alcohol-based hand wipes can encourage cleanliness for everyone entering the equipment room.

Loose papers, books, supplies, newspapers, and trash are fire hazards and also must be banned from every equipment room.

23.5.3 Overt, Covert, and Deceptive Protections. Effective protection of the IS infrastructure requires many hidden elements—such as concealed surveillance cameras, sensors, and detectors—and all of the wiring that supports them. But good protection also needs some clearly visible elements. It is important to consider which devices are best hidden to protect them and which should be visible as deterrents.

Overt devices are ones that are evident to workers and visitors, or whose presence is implied by other visible objects, such as warning notices. These visible devices, which suggest that some sort of security exists, are intended to deter troublemakers, so that all but the most determined attackers will go elsewhere. Examples are surveillance cameras, access controls, visible alarm boxes, and visible sensors. Although most overt devices are active and recording data, some may be inexpensive dummy devices that only look real, perhaps with slowly blinking indicator lamps to heighten the effect. Covert protection, however, must not be noticeable to either visitors or insiders. There

OTHER CONSIDERATIONS 23 · 11

must be no indication that these protections exist, what they are, how they might function, or where they are located. Most effective security systems operate covertly; examples include stealth and silent alarms, concealed early-warning systems, perimeter and proximity sensors, access monitors, and many other surveillance devices that are not readily seen.

It is important also to conceal the wiring that interconnects all protective systems and the utilities that support them. Whether any part of a system is visible or not, the wiring that connects it should not be. Although overt devices may themselves be vulnerable, they will generally advertise that there is good security here and everyone within the premises can feel safe. However, visible devices can sometimes be covered or spray-painted, knocked aside with a club, or shot to disable them. An expert thief can also defeat many hidden systems, if he or she knows what they are, where they are located, or how they are connected.

Another approach to protection involves deception. Dummy devices that look like surveillance cameras, access control devices, and alarm sensors can be placed to attract troublemakers, who may think they can physically damage, disable, or circumvent the system. These visible devices are intended to distract potential troublemakers and divert them away from vulnerable areas. Some devices are deceptive in that they are not what they appear to be but are actually alarm sensors to measure motion, proximity, sound, or anything that disturbs the device. Deceptive devices often are used to divert troublemakers away from vulnerable people and infrastructure, by offering them a “honeypot”: an attractive target to distract them, but often a target equipped with an alarm device, surveillance cameras, or other means of identifying a perpetrator and gathering evidence.

There is a gray area between what management can legitimately do to protect its information systems and what may be unethical or illegal actions. Management has a legal and fiduciary responsibility to protect people and property, and those who support deception say that these techniques are increasingly necessary to protect an organization. Others insist that this amounts to entrapment or violates privacy rights. State and local regulations and interpretations vary widely and are continually changing. It is necessary to check carefully with local officials, legal advisors, and insurers to determine what is acceptable and how to manage such risks. Management must then decide to what extent these techniques may be effective and whether less contentious approaches will suffice.

Whether the protection devices themselves are overt, covert, or deceptive, the security systems behind them must not be obvious. No one seeing or knowing about the elements of a security system should be able to deduce the details of the system, the functionality, or where and how it is monitored. An observer may notice a particular device or product or the suggestion of a vendor's standard security solution, but the particulars of the protection systems must remain obscure, and all the wiring that supports them must be hidden or disguised as well.

Everyone involved must be aware of the security policies and procedures. Conspicuous signs should advise that anyone entering the premises may be monitored, as may all communications. All of the security policies and procedures should be understood and accepted by everyone involved. Employees and other on-site personnel should receive periodic security-awareness training and briefings. And there should be periodic security exercises and drills to test the procedures and reinforce the training.

Finally, protection must not be intrusive. Security cannot limit productivity or IS performance in any way. Instead, the protection must contribute to a feeling of safety and security within the workplace and thereby enhance productivity.

23 · 12 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

23.6 ACCESS CONTROL. Access control systems are but one layer of good infrastructure protection. They are usually used in conjunction with surveillance and perimeter control systems in order to provide sufficient early warning to head off trouble before it happens. Effective access control requires three tiers of support that are described next. The strength of each tier and its integration with other security layers determines the security effectiveness.

- 1. Privileges.** This tier determines whether a person seeking entry is authorized. It is the initial entry-request process that may use an access or proximity card, radio-frequency identification (RFID) key, or a combination lock. Since many of these devices can be lost, borrowed, stolen, or copied, and many can be quickly defeated, there is usually not much effort to ascertain just who is seeking entry. Therefore, privileges alone are not strong security.
- 2. Authentication.** It is usually necessary to identify a person seeking entry with some degree of certainty. To do this, the person must possess or know something unique that others cannot readily duplicate. Examples include personal identification numbers (PINs), electronic keys, entry cards, and biometric devices. PINs and passwords may be used, provided they are strong and well implemented. Some of these approaches merely strengthen the privileges process but can still be copied or defeated.
- 3. Audit trail.** A log is required for each entry attempt to show the date and time, the identification of the person, and the action taken by the access control system. Access-denied and unable-to-identify events should trigger immediate alarms. Logs must be analyzed in a timely manner for anomalies or unusual patterns. Where better access control is needed, each person's exit also must be authenticated and logged.

See Chapter 28 for more details about identification and authentication in access controls.

23.6.1 Locks and Hardware. Strong protection begins with high-quality locks, door hardware, and access control systems that are nonintrusive yet strong enough to deter most unauthorized entry. Lock types should be hard to pick and should use keys that are hard to duplicate. Examples are Medeco® locks and keys with dimpled sides. Ace® locks with circular keyways require special tools to pick. Many types of keys can be created from lock numbers, so keep these numbers stored securely.

No lock is completely safe. Someone with equipment, experience, and time can open any lock, often very quickly and without causing attention. Where key mastering is used, an experienced person with a key to any single door can open the lock cylinder and copy the mastering system. Therefore, additional layers of protection are needed.

Interior areas accessed only by a few people usually can be secured with a strong push-button combination lock. Key locks are not appropriate, because the keying cannot be changed periodically or quickly when a key is lost or someone leaves. And misplaced or lost keys may not be reported for days or weeks. However, key locks tend to be stronger and less vulnerable to vandalism, so keys may be the best alternative for outside areas or for doors that remain open during business hours. Wherever keys access critical areas, there should be spare lock cylinders and a new set of keys stored on site that can be utilized quickly when a change is needed. Once a lock cylinder is changed, the old cylinder should be rekeyed and a new set of keys produced.

ACCESS CONTROL 23 · 13

Locks that use an electronic key are particularly effective. Electronically activated cylinders can replace existing mechanical cylinders, and many do not require any wiring, so the hardware and operating costs are minimal. Most electronic keys have a small cylindrical tip that is touched to the lock for access. Both the lock and the key can log each event, identify the specific lock or key used, the date and time, and whether access was granted. And conveniently, electronic keys are not much larger than mechanical keys. However, both can be defeated with the proper equipment and skill. Often, two independent entry locks are utilized to provide stronger, relatively inexpensive protection.

RFID devices can provide good access control, provided the activating device is not lost or stolen. RFID can be an improvement over card entry systems in that it can activate a lock from several feet away. RFID access cards can include photo ID, name, and perhaps an optical stripe for encrypted identity information. As with other systems, RFID can be breached with the right equipment and skills.

Another inexpensive upgrade of key locks is the card-access lock similar to the ones used by hotels. Many do not require wiring and are battery operated, so the keying of each door remains unchanged and an old key remains valid. Therefore, without central wiring, most card-access locks offer limited security and cannot trigger an alarm although some do at least log events.

To prevent intruders from bypassing the lock mechanisms by using credit cards or thin metal sheets to push the bolts inward, the area over the actual latch should be protected with an *astragal*—a flap of metal bolted to the door. In addition, latches should include anti-tampering features that block external movement of the moveable elements into the door when the locks are active. The hinges for doors should not be accessible from outside the secured area when the doors are closed: It is a simple matter to break and remove the hinge bolts and remove a locked door if the hinges are on the outside.

No matter how strong the access control systems, doors have their limitations. Absent a vault-type door and hardware, a determined attacker with pistol and silencer can gain access readily. A small water cannon that is transported in what looks like a toolbox can breach a standard metal door, or the partition surrounding it, with one shot and very little noise.

23.6.2 Card Entry Systems. One of the best means of physical access control is a card entry system, especially the newer systems that are increasingly capable and less costly as well. A central card entry system often controls the entire premises: all entrances, exits, elevators, rest rooms, and many other interior doors. Access cards are usually similar to a credit card and can be carried concealed or worn as identification badges. Access cards are usually imprinted with a full-face photo, the individual's name, the organization's logo, and often a printed stripe with biometric information. Access cards often show number codes to indicate authorized areas of access, and are usually color-coded to indicate status and whether the person must be escorted.

For more details of identification and authentication systems and card entry systems, see Chapter 28 in this *Handbook*.

23.6.3 Proximity and Touch Cards. Some cards can be held near a sensor (*proximity cards*) and some can be held in contact with a sensor (*touch cards*). Both types of tokens can reduce wear on physical tokens that have to be put directly in contact with the sensors.

23 · 14 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

For more details of identification and authentication systems and proximity and touch cards, see Chapter 28 in this *Handbook*.

23.6.4 Authentication. Anyone can use an access control token that may be borrowed, lost, copied, altered, stolen, or taken by force. Therefore, authentication is another layer of security that is needed to establish the identity of the person seeking entry with some degree of certainty. Authentication devices commonly include a biometric scanner, a numeric keypad, and visual or voice identification by a computer or by another person. All such devices come in varying security strengths and each can be used in combination.

When an access token is read, the system must verify that the *token* has the requisite privileges for that place, date, and time. If it does, it then becomes necessary to authenticate the identity of the person using the token. For this purpose, a numeric keypad was once the most common device whereby the token user entered a personal identification on a number keypad. Now most use a touch screen where the numbers will appear in random order. If the system validated the PIN, it activated an electric door strike to momentarily unlock the door. This system can be slow, cumbersome, and easily compromised. In some systems, everyone uses the same PIN and is supposed to keep it private. However, PINs can be forgotten, lost, or discovered by others.

Visual authentication is a better approach. This can be done by a computer or by a security guard or receptionist who can see the entrance or monitor it via a surveillance camera or a video intercom. Emergency-type video intercoms work well because they provide the remote authentication with a visual image and can monitor sound continuously, so that the authenticator can speak with and challenge whoever approaches the door. When identify is verified, the system or other person activates the electric door strike to unlock the door. This system offers stronger security and is faster than a keypad. It also facilitates recording and logging all entrances and exits, especially during off-hours. Breaches can happen, if the person seeking entry can deceive the guard or receptionist.

Biometric scanners offer the strongest security. They offer faster, more positive identification of every individual and do away with the need to remember a PIN. Biometric scanners can read any of these personal attributes, which are listed somewhat in order of their current popularity: fingerprint patterns, facial characteristics, voiceprint, handprint, signature, retinal details, or height and weight.

Most biometric systems can be adjusted to be highly sensitive (which is slower and may require repeated entry attempts, but is very hard to breach) or less sensitive (which is still fast but may result in some false authentications). Before choosing a biometric system, it is necessary to determine that the users understand and will accept it. For example, some people balk at having to use a retinal scanner, while others may feel that any biometric device invades their privacy. The latter reason is invalid because most biometric systems cannot be used to identify any person not already known to the system. Increasingly, especially in critical public places, these systems can check every individual against criminal and terrorist databases and quickly alert law enforcement. It is difficult to steal an identity using biometrics because one must also know the encoding algorithms used.

Fingerprint scanners are the most common and are becoming increasingly powerful. Initially, these utilized optical scanning, which could be fooled by photographs, wax impressions, or by a severed digit. The newer capacitive scanners use electronics rather than optics and can provide nearly certain identification. Usually, three of the user's fingers are "enrolled" in the system in case some fingers are later bandaged or dirty. If

ACCESS CONTROL 23 · 15

all the enrolled fingers are incapacitated, single-use passwords can be used to bypass the system.

Most fingerprint scanners cannot identify an individual by name, but only that a person seeking entry matches the person whose biometric identity has been enrolled. Most scanners do not conform to the uniform, automatic fingerprint identification standards used by law enforcement. Instead, they scan a small area of the finger and apply proprietary algorithms and encryption. A template from one system is usually meaningless to another.

Accuracy of fingerprint scanning is affected by the angle and position of the finger and by the pressure applied. Most systems allow sensitivity adjustment to optimize enrollment and verification times and success rates as well as to minimize delays and false negatives that require repeated access attempts. While well suited to most applications, fingerprint scanners may not be appropriate where the user could have dirty or thickly callused hands or must wear gloves (such as healthcare workers).

Facial recognition is the basis of another popular scanning system. It uses graphics technologies and any surveillance camera to measure the size and relationship of prominent facial elements. Most systems are proprietary and cannot be used to identify an individual by name, but some others are compatible with law enforcement standards. Sensitivity and accuracy are dependent on the distance, position, and angle of the head as well as on the background lighting. Cameras at all entry points must be positioned to photograph the subject at the same angle that their faces were enrolled. Not all facial recognition systems offer strong security; some can be fooled by a face mask or a photograph. Others are best left at their highest sensitivity settings to avoid spoofing, which may require increased false positives and multiple attempts at entry.

Voice-print scanning can be used, but mostly for access to a terminal or workstation. The better systems display random words on the monitor so a prerecorded response cannot be used. Most use the workstation's microphone or, increasingly, a cell phone. Voice scanners can be affected by hoarseness, so there should be a one-time password access provision. These systems can be useful for remote login, especially while traveling, although the low bandwidth of dial-up circuits may lessen the system's usefulness.

Retinal or iris scanners are considered the best security, but authentication can take a few seconds. For access control, the user must generally look closely into an eyepiece, which is traumatic to some people. But for access to a terminal or workstation, a Web camera is generally placed on top of the monitor, 17 to 19 inches from the eye. The user's head must be positioned to align a white dot within a colored circle. The user must hold still, without blinking, while the scan proceeds automatically. The camera can be used for surveillance to see who is near the workstation.

For details of biometric authentication, see Chapter 29 in this *Handbook*.

23.6.5 Integrated Access Systems. Biometric scanners are used increasingly by other applications, and most can readily coordinate with an access control system. Applications include network user authorization, and access to terminals and workstations and to software applications and data. Biometric readers can be built into laptops, keyboards, mice, or peripheral readers. They can be used with access cards, badges, and proximity or RFID devices for authentication with some degree of certainty. Scanners that are mostly proprietary are currently integrated with encryption systems (e.g., virtual public networks, public key infrastructure, and smart cards) to authenticate transactions including credit cards, financial, banking, and automatic teller machines. Biometrics are increasingly used to identify hospital patients, welfare

23 · 16 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

recipients, people who frequently enter the United States, and similar applications that involve identifying a diverse or widely dispersed group.

Infrastructure security protection can be independent of, or integrated into, a comprehensive premises security system. Either cards or badges can provide many other functions beyond basic access control. Applications include off-hour access to the building, elevators, and premises; control of building entrances, restrooms, and parking areas (especially at night); purchases at a cafeteria or company store, or for charging stationery or materials picked up from a supply room. They can also be used to receive classified documents.

For greater protection and efficiency, an integrated enterprise-wide, Web-based system can control access to various premises, locked doors, infrastructure components, networks, workstations, applications, secure data, and Web connections. This arrangement offers comprehensive security by logging every event. Centralized logs can yield much more meaningful security information, because an integrated system provides better early warnings to head off trouble before it happens.

23.6.6 Portal Machines. Airportlike security checkpoints are increasingly being installed outside airports. A portal machine is the archway one must walk through during the airport security process. It detects concealed metal, such as guns and knives, and tools that might be used to cause trouble. A portal can be used to detect IS components (such as storage media) being smuggled into or out of the premises. For example, the Texas Capitol building had such a system installed in May 2010; however, Texans carrying concealed handguns need merely show their license to be able to circumvent the scans entirely.¹¹

Newer trace-portal machines can detect many explosives as well. The person entering is asked to stand still for a few seconds while the machine releases several puffs of air and captures samples that are then analyzed for a number of hazardous or explosive substances. When a guard is not present, a computerized voice instructs the person to proceed, or a turnstile can hold the person in place.

Such devices might be appropriate for all who enter premises requiring ultra-high security and at entrances to critical areas or equipment rooms as well.

In January 2013, the Transportation Security Administration (TSA) announced that 200 specific full-body scanners would be removed from U.S. airports because of concerns over the resolution of the nude images produced by the Rapiscan™ machines.

In December 2012, the DHS announced that it would contract with National Academy of Sciences for a thorough study of the intensity of electromagnetic radiation in security scanners.¹²

23.6.7 Bypass Mechanism. Whatever the systems used, there should be one bypass key that can open every door that uses electronic or electrical controls. The bypass key is for emergency use only. It is a unique key that is not on any mastering system and is available only to a few highly trusted people. The cylinders should be heavy duty and very hard to defeat, with the keys nearly impossible to copy. Careful control and protection of each key is essential. The loss of a key may not be discovered quickly, and the time and costs of rekeying every lock will be substantial. Meanwhile, every door on the access control system is vulnerable.

Bypass passwords for individual users also may be needed. These passwords should trigger a silent alarm whenever used, so that security can respond to the site or verify identity by telephone, surveillance, or intercom. One-time passwords provide the best security.

ACCESS CONTROL 23 · 17

23.6.8 Intrusion Alarms. Intrusion alarms are necessary to provide perimeter and early-warning alarms and are usually needed as extra layers of security. There are several methods of intrusion detection. Digital surveillance cameras with motion detection are best because they can monitor visually what is happening and record what they see. Other methods include proximity and pressure sensors mounted within the perimeter walls or floors or inconspicuously within the room. Most of these sensors can detect intrusion and forced entry and can pinpoint the location of trouble, but they provide no details, monitoring, or evidence-recording capabilities. Proximity and pressure sensors can protect long-distance perimeters, cable runs, and utilities inexpensively. Concentric layers of such devices are necessary for sufficient early warning to prevent trouble from happening.

The best motion detectors use digital closed circuit television (CCTV) surveillance cameras that can sense movement while observing and recording the event. Miniature cameras that are inconspicuous or concealed are particularly effective and increasingly inexpensive. Several cameras can record pictures from many angles and often can identify an intruder positively. In a larger area, cameras often use swivel-tilt mounts and zoom lenses, which can be controlled automatically and remotely as well. Color cameras are preferable, as are cameras that automatically adjust to light conditions, including near darkness.

Some CCTV cameras include a military-type night-scope mode, which is relatively inexpensive and functions well in near-total darkness. These cameras also work well in normal light and can switch automatically to night-scope mode when necessary to see clearly and to record evidence. Other types of cameras, such as infrared, work well where there is no (humanly) visible light.

Other intrusion detectors use radar technology or measure changes in capacitance or inductance to sense intrusion. Most cameras and detectors can trigger an alarm as soon as they are disabled or lose power. Detectors may be wall- or ceiling-mounted devices as an overt means of deterrence, but these may then be vulnerable to spray paint, wire cutters, a gun, or a club. Therefore, intrusion detection sensors are usually concealed, or at least inconspicuous.

Perimeter alarms are especially important in building core areas, or public areas within or outside a building, to provide ample early warning of an intrusion that might soon affect a secure area. Digital video cameras are best for this but may be ineffective over large areas or long distances. Therefore, many proximity devices are used to monitor intrusion. Most utilize long sensor wires that are surface-mounted inconspicuously or hidden within partitions, ceilings, and sometimes inside of conduit. These systems detect the presence of a standing adult or a large animal that may come within a few feet of the sensor wire. The sensor wires can be very long, so zoning is often necessary to pinpoint an incident at least within a general location.

A better and cheaper alternative can be fiber optic perimeter alarms. Developed for the military and national security installations, the fiber optic systems are very sensitive and can monitor, evaluate, and record events. The sensor wires can also be embedded inside drywall or masonry partitions, ceilings, floors, or conduit and will detect both pressure changes and ambient sound. Because they do not measure proximity and can monitor and evaluate events, false alarms are less likely. They can warn of an impending accident or efforts at forced entry, and may soon be able to locate the event as well. These systems use software that can discriminate between recognized events and situations that are unknown or potentially dangerous. All the events are recorded and can be replayed and reviewed by a remote operator at any time.

23 · 18 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

Increasingly, inexpensive monitors are available that connect wirelessly via Wi-Fi, mesh networks, cell phone, satellite, a radio channel, or sometimes via microwave. Such devices save the high costs of direct wiring to remote areas and the need to protect and conceal the wiring. But there is also a trade-off with reliability and security, as any radio-frequency device can be monitored or jammed and may also be spoofed.

In addition to intrusion alarms, environmental alarms should be provided to measure temperature, humidity, and smoke, fire, and flood situations within all critical areas and equipment cabinets as well, as described in Section 23.8.1.

23.6.9 Other Important Alarms. A duress alarm system is recommended within most critical areas. This is usually a large red button or an emergency intercom conspicuously mounted near the exit. It is used if someone is injured or ill, when trouble occurs that other alarms will not immediately sense, or if violence is threatened. The emergency intercom types are best because each party can talk and listen. Those with a CCTV camera are described here, but many inexpensive emergency intercoms that provide audio only can be useful. Security personnel can constantly monitor all sounds within a secure area, and anyone inside can readily talk with security personnel. Duress alarms are usually silent. And activation devices can be concealed or located inconspicuously in case a potentially violent situation erupts. Duress alarm activators inside of cash drawers, in the knee wells of desks, or under counters have been common for years.

Another form of duress alarm can be integrated into card-access systems or password systems, for example, swiping the card backward may allow access but also silently trigger alarms at security stations. Proximity cards could be programmed to recognize peculiar methods of waving the cards which might pass inspection of those applying duress yet initiate appropriate defenses and alarms. Entering a keycode with an additional sequence could do the same. All such duress-alarm signals should be part of every employee's training for use of the security systems.

Beyond access control and authentication, it is also important to know whenever a locked door is not fully closed. There should be a sensor that warns whenever an access-controlled door is open or ajar. The sensor is normally built into the door buck (frame) to provide a silent alarm. An open-door alarm system delays for a few seconds in order for one authorized person to enter or exit normally. A short delay prevents leaving a door ajar for someone else to push open, piggybacking when another person enters or leaves, and tends to prevent anything large being taken in or out. The open-door alarm also prevents a door from being propped open at any time. The door-ajar sensors should be concealed at all times so users are unaware of their existence. Otherwise, they may be taped, jammed, or otherwise defeated.

A final security protection is to prevent *double-dipping* whereby an authorized person requests multiple entries in order to admit unauthorized persons. The access control system can prevent this; once persons enter a space, they must be logged out before they can try again to enter.

23.7 SURVEILLANCE SYSTEMS. Today, surveillance systems are designed and laid out to document fully every event, to facilitate positive personal identification, and to provide legal evidence when needed. Today's digital cameras, controllers, and recorders can do all this and more. The old analog or film cameras and recorders are inadequate and should be replaced, especially since the new surveillance systems are much less expensive.

SURVEILLANCE SYSTEMS 23 · 19

If the protection is designed well, cameras can provide an undisputable, accurate, historical record that is available instantly. Cameras never sleep and are not distracted by repetitive tasks or boredom. More important, cameras can provide early warning, and can document events from many perspectives, concurrently and synchronously. Cameras increasingly incorporate microphones and interface with emergency audio intercoms, the better to assess a situation, assist people on the scene, or challenge suspicious persons. The very presence of visible but inaccessible cameras usually deters most troublemakers. Protecting the infrastructure requires both early warning and identifying the nature of trouble before it can happen. The new surveillance systems can do this effectively. And they can be integrated with other alarms and with premises security systems for strong, seamless protection.

23.7.1 Surveillance Cameras. Surveillance cameras are far more effective and much less expensive than guards, watchmen, or extra receptionists. However, for strong, redundant, and flexible security protection at important locations, both cameras and people will be needed.

Surveillance capabilities are becoming increasingly effective protective devices. They can now see more, detect some dangerous hidden materials, and, increasingly, decide what is important to monitor.

Equipment rooms and other sensitive areas once relied on motion detectors for security, but now digital cameras with motion detection work better. Digital cameras also function well indoors and outside, and under most light conditions, from sunlight shining into the lens to near, if not total, darkness. When used outdoors, digital cameras are immune to all but very severe weather conditions. They are, however, affected by heavy smoke, snow, ice, and rain. Wind-blown objects and sometimes birds or small animals can distract their motion-detection systems.

Today's digital systems are far more capable and reliable than the earlier analog or film cameras. Cameras were once the weak link because the details were not clear or too small. Many of the old cameras were needed to cover an area, and even then, necessary details were often out of focus. Now cameras are much smaller, less expensive, and take advantage of auto-focus, tilt, and zoom capabilities to yield much better pictures. Almost all cameras now use color as well for better clarity. Some also incorporate an inexpensive, monochrome, night-scope mode that works well in total darkness. Digital cameras automatically correct electronically and mechanically for varying ambient light conditions. Most can correct for background lighting conditions that would otherwise cause subjects to be underexposed, for sunlight glaring directly into the lens, and for unusual brightness that would otherwise wash out a useful image. All of the images can be viewed in real time, and one or more persons can simultaneously control the views and scroll back to freeze and enlarge frames. All viewers can adjust the brightness, contrast, colors, and zoom, and can apply filters to bring out foggy details. Any viewer can save anything viewed.

All but the smallest indoor cameras now include a microphone so that both video and audio are recorded. Outside the facilities, some installations interface with an emergency intercom so background sound is continuously monitored and dialog recorded. Emergency intercoms are particularly good deterrents because security personnel can confront possible troublemakers and advise them that they are being recorded. Most will back off or flee before trouble occurs. Although many camera control systems provide electronic zoom, the primary zoom function should be mechanical. Electronic zoom systems are inexpensive, but they lose resolution and images can quickly become unidentifiable.

23 · 20 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

Miniature cameras can be particularly useful and save money as well. Cameras less than one inch in diameter can provide sharp detail. Some units with zoom lenses and pan and tilt mounts can be hidden within any common object.

Opinions vary as to the ethics of using concealed cameras, and some state laws limit their use as invasions of privacy. Generally, the rules of good surveillance practice are the same as monitoring telephone conversations, email, and Internet use. In order to protect itself, the organization has a right to see and hear whatever is going on in and around its premises. But signs should be posted so that workers and the public are advised that all persons entering the premises may be monitored for safety and security purposes. Company policies should explain this, and why surveillance is necessary.

23.7.2 Camera Locations and Mounts. To serve as a deterrent, some cameras should be readily visible, but this may present some problems. Overt cameras placed too high are not good for identifying people, while cameras placed too low can be spray-painted, covered, or smashed. Any visible camera can be shot out (perhaps unnoticed using a silencer). Older cameras are good for overt use because they are large and may provide backup surveillance. Another option is dummy cameras, some of which have visible indicator lights or automatically pan to attract attention. Some dummy cameras are in fact alarms that will trigger if the unit is disturbed or if its wires are cut.

Surveillance of a wide area requires an elevated camera that can see greater distances, zoom in on details (especially faces and vehicle license plates), and closely follow an event. Wide-area views are also necessary to spot multiple troublemakers. Several camera angles will likely be needed to gather good evidence.

Outdoor cameras are usually attached to a swivel, tilt mount, and inside of weather-proof domes or larger enclosures when large zoom lenses are needed. There is usually no attempt at concealment because the cameras are too high to reach, and they provide good deterrence as well. Many may be inconspicuous, however, or even disguised. When large cameras and enclosures are used, the direction the camera is looking in may be obvious. Often, troublemakers will create a distraction to lure the cameras away from the real problems.

23.7.3 Recording Systems. Once VCRs were the norm, but now most recording is done with hard drives that can store hundreds of hours of audio and video taken from multiple sources. Increasingly, solid state drives are also used.

Once, only one person could review the tapes and then, usually, only after the recording was complete. Now one or more people may access and analyze the information simultaneously. Once, 31 VCR tapes were recommended for each system: one tape for each day of the month, after which the tape was inspected and reused. Now one hard drive can record a month's data, which can then be inexpensively archived on DVDs or at remote storage locations. The VCRs often recorded in a lapse-time mode to save tape until an event occurred to trigger real-time recording. Now everything can be recorded continuously in real time to examine fully the happenings before, during, and after an event.

23.7.4 Camera Control Systems. Camera control systems commonly can direct the swivel (pan), tilt, mechanical zoom, aperture, and background lighting of each camera. Some control systems can automatically pan cameras back and forth across large viewing areas and also provide motion detection. Some control systems

FACILITIES DESIGN 23 · 21

can also stop panning and zoom in on any unusual event. This is a valuable feature that must be carefully programmed so the system is not distracted by spurious events or deliberate diversions. Usually, each camera can also be controlled manually and often from remote locations.

A major advantage of a digital system is that each camera provides continuous images, usually at about 30 frames per second, which fiber optic and other broadband connections bring to the recording system in real time and high resolution. The control systems usually allow one or more persons to roll back the images without affecting the real-time recording. Each viewer can scroll backward and forward, freeze frames, zoom, crop, or enhance the image electronically, and save any material to another medium or system as evidence.

23.7.5 Broadband Connections. The advantage of broadband camera connections is that the information is available in real time and with high resolution, and that more camera control functions are possible. Fiber optic cabling does this best, but it should be dedicated and well protected. Remote cameras can also be connected via a LAN, wireless, or broadband Internet connection using Transmission Control Protocol/Internet Protocol (TCP/IP), if done carefully, so that the data are secure and other traffic is not impeded.

Digital multiplexing is another advantage of broadband connections. Multiplexing over metal wiring once necessitated delays, lower resolution, and fewer frames transmitted per second, even though each camera was providing continuous, high-quality images. Now digital multiplexing over fiber connections allows all data from all cameras to be transmitted simultaneously. Multiplexing is accomplished remotely as signals from several cameras are combined into a single broadband connection to the control system. Connections among control systems are also multiplexed to minimize line charges. The signal interfaces between each camera and the multiplexer can be fiber optic, microwave, a network, or the Internet. This way, hundreds of cameras can be networked economically.

Another benefit of today's TCP/IP connections for cameras is that they can be Internet-enabled, allowing monitoring stations to receive data without having to retrofit special-purpose cables throughout the structure. A disadvantage of such systems is that they are therefore susceptible to the same attack methods that are developed to hack any other Internet-based systems. Proper use of encryption (e.g., virtual private networks) for camera output and controls can help protect the surveillance systems against electronic tampering. Encryption can also interfere with the *playback attacks* so popular in movies, where criminals insert devices to cause an infinite loop of playback and fool observers into believing that all is well in the compromised surveillance zones.

For additional information on penetration techniques, see Chapter 15 in this *Handbook*.

Radio or microwave multiplexing is another option, but it is not very fast and is subject to interference or disruption during severe weather conditions.

23.8 FACILITIES DESIGN. Good protection begins with good facility design that will ensure the safety of the information infrastructure and the people who use it. Protective systems become expensive and inefficient when the premises themselves do not facilitate good security. Proper interior design, facilities layout, engineering, and construction can maximize security effectiveness, minimize costs, enhance productivity

23 · 22 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

and morale, and generally boost profits. The starting point is an inspection of all sites and review of all as-built plans and construction documents.

Premises security inspection and review are best augmented by using independent, outside security experts. Comprehensive and objective architectural, design, engineering, IS infrastructure, and premises security experience are all needed and may not all be available internally or from vendors. The inspection and review process always must be threat-specific and must relate to a thorough threat assessment (which is described in Chapter 22). The premises inspection can then serve to validate the vulnerabilities identified by the threat assessment. Once some of the vulnerabilities are corrected, the premises inspection should be repeated.

Effective infrastructure protection can prevent trouble from happening. To do this, there must be ample early warning, which, in turn, requires good facilities planning and design, effective premises and infrastructure security in place, and the awareness and vigilance of everyone in the area. All must work together seamlessly, efficiently, and proactively.

Good facilities design can be efficient, nonintrusive, cost effective, and inexpensive—even within existing facilities. Here are some guidelines and suggestions.

23.8.1 Choosing Safe Sites. If possible, the building site for any facility that will house IS infrastructure—especially for highly secure applications—should be far from obvious manmade threats such as fuel depots, railway lines, suspension bridges, high-speed highway overpasses, and airport flyways. Sites should avoid obvious natural threats such as floodplains, areas subject to forest fires, active volcanoes, and old underground mines subject to subsidence of the overlying areas.

Access to the site should include provision for emergency evacuation of employees and visitors and access for emergency vehicles to reach the facility even if the main public access is blocked.

In cities, the building should not be situated in high-crime neighborhoods or in industrial parks with inadequate security supervision or histories of theft, vandalism, or sabotage. Decisions should include attention to public records indicating sufficiently competent police and fire-response teams; interviews with other users of industrial or commercial buildings in the area may be helpful in making decisions about the final choice of location.

Sites for equipment rooms and utility closets should be protected from possible threats. Infrastructure sites should be located far away from all piping, tanks, or equipment that uses any liquid that could possibly leak or spill. Most plumbing tends to leak unexpectedly at some time or can be intentionally breached, so it is well to assume that any pipe, connection, container, or pump will eventually burst or leak. Placing sensitive sites at a distance from such threats is safer and cheaper than any other form of protection. The danger zone where leaks may spread must include an ample vertical distance and horizontal area. Begin the danger zone with a pyramid from any leak source, downward several floors, and outward horizontally well beyond the infrastructure areas.

Most buildings require fire-suppression sprinklers in all areas. Any may activate by accident, vandalism, or an actual fire, and may cause considerable flooding. Also, infrastructure sites should be located away from windows, exterior walls, roofs, skylights, building setbacks, and storm drains, which are all potential sources of flooding. Treated water used in heating, air conditioning, and solar systems presents a worse problem, in that the chemicals can quickly destroy electronic equipment, connectors, and wiring.

FACILITIES DESIGN 23 · 23

If all of the infrastructure cannot be positioned at a safe distance from all liquids and hazardous materials, there must be special protections and alarms. Protections include sealed, waterproof construction, drains, berms or other diversion devices, and proper materials close at hand to control any spill. There should be environmental alarm sensors near where leaks or spills could occur. Floor drains must also protect the equipment areas, especially those that use sprinklers; otherwise, cleanup will be difficult.

Infrastructure sites should not be visible or accessible from any public area. Infrastructure wiring or cables should not run through the building core, public, or mechanical areas, including rest rooms, slop sinks, janitorial closets, and stairwells. Avoid placing equipment or cables where any persons might loiter. All equipment room entrances should be clearly visible from workplaces where employees can readily observe all comings and goings. Choosing inherently safe sites and entrances greatly reduces both risk and costs because less security is needed.

For effective security control, there should only be one access point to each critical area, used for both entry and exiting. However, if local fire or building codes require a secondary means of egress, a delayed-access panic bar is usually acceptable. Such a system delays releasing the exit lock for a few seconds, while an alarm sounds and a surveillance camera is triggered. There should also be surveillance cameras with motion detection throughout all secured areas. All locked doors should look as alike as possible from the outside and be identified only by a room number that looks similar to that of any other premises door. No functional name, person's name or title, or any other means to identify what is inside a locked area should be apparent. No signage or directory should include a functional, personal, or departmental name, but only area or room numbers and directional arrows if needed. Only floor, room, or suite numbers should appear on premises signs, on floors.

23.8.2 Physical Access. Access to the facilities site should be controlled. Guard stations at the entrance won't stop terrorists, but they can provide early warning of trouble. All such stations should be under constant surveillance and have panic buttons installed for emergency notification of the security teams that there's a problem at the entrance to the site.

Reinforced concrete bollards can prevent physical attack on a building by interfering with even high-speed ramming by ordinary vehicles such as cars, vans, and light trucks. They will not necessarily prevent attacks using commandeered large vehicles such as loaded dump trucks.

The entrances to buildings should not include large glass walls, which could be demolished easily to gain access to entrance halls for detonation of bombs.

Physical access to all parts of the information infrastructure should be restricted. All information system and network equipment must be inside equipment rooms that are always locked and accessible only to those who need to be there. All utility and wiring closets, power distribution panels, patch panels, wiring blocks, and terminations should be located inside equipment rooms that are always locked. If possible, do not allow unrelated systems, equipment, or utilities inside a restricted area, so that a technician working on an unrelated system cannot access the information infrastructure. If this cannot be avoided, IS personnel should always escort others entering these areas. In high-security areas, all persons entering should be escorted, or use two-trusted-person teams where each person observes the other. Guards and facilities-security personnel are good premises-security escorts but probably do not know the infrastructure and are not the best choice here. IS personnel are better escorts, even though they may

23 · 24 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

not know all the security details. In any event, all persons entering must be positively identified and each visit logged.

It is wise to put critical electrical distribution panels inside an equipment room, so they are quickly accessible and also protected. This is a safety issue as well; someone working on equipment or wiring can readily see that the circuit breaker is off. Otherwise, electrical distribution panels should be located inside a locked area that is unmarked from the outside other than by a coded location number. Panels are often located in public areas in many buildings; they must then be securely locked and alarmed as well. Whenever an electrical panel controls anything critical, access to the area should be restricted and the room alarmed. These precautions reduce any loss of power to critical systems by accident or intentionally.

One of the cases used in training courses about facilities security is the unlocked patch panel discovered by a physical-security expert in the public hallway of a hospital; the director of the intensive care unit (ICU) was horrified to realize that the breakers controlling the life-support equipment for the patients in the ICU were open to any child or malicious person.

A mantrap can best control access to critical equipment rooms. A mantrap is a two-door arrangement with an inner door, outer door, and a short corridor in between. Both doors are interlocked; one must be closed before the other can open. The corridor usually is constructed with a shatterproof, full-height glass partition (often one-way glass) on one or both sides for surveillance. Both the doors are usually windowless. And each door must have a strong access control system. But for safety, the outer door can usually be opened from inside the corridor using an alarmed panic bar. Usually, one or more surveillance cameras are positioned to identify anyone entering or leaving the mantrap.

Emergency intercoms within the mantrap corridor and at the entrance and exit points are strongly suggested. Conspicuous duress/assistance buttons should activate silently, so security personnel can monitor the area and speak with or challenge anyone who cannot pass through properly. The alarms are best silent so that anyone under duress is not further threatened. Emergency intercoms often include small cameras that are inconspicuous or concealed, so that security personnel can see what is happening, assist if someone is ill or somehow becomes trapped, and avoid public-safety issues. A well-designed mantrap can be valuable for all but heavily traveled entrances.

Mantraps can tighten security in many ways and are usually not intrusive. They can detect *piggybacking*, when one or more extra persons closely follow someone who is authorized to enter. Mantraps in themselves do not preclude propped-open doors, but they make removal of objects from the room difficult and risky. The access control log and surveillance recordings can identify troublemakers and provide strong evidence to convict those who might otherwise be suspects or persons unknown.

23.8.3 Protective Construction. The design of a new building can improve physical security.

- Avoid vulnerable construction methods such as using structural columns in large, open spaces on the ground floor; destroying such columns using explosives could cause serious damage to the structural integrity of the building.
- Don't allow architects to include decorative indentations on the outside of the building; such vertical patterns can provide the means for rock-climbing attackers to climb the walls easily.
- Place all IS in a hardened, isolated area.

FACILITIES DESIGN 23 · 25

Equipment rooms require sturdy partitions for good security and to support the considerable weight of wiring and equipment. Moreover, the partitions and walls must remain safe and stable during any seismic activity, such as heavy road traffic or a sonic boom, explosion, or earth tremor. Floors may also need to be reinforced. Sturdy partitions deter forced entry, which may be otherwise accomplished with little more than a pocketknife and a fist. Existing walls, partitions, and floors should be inspected and any subsequent alterations approved by a structural engineer. Consider too what might be needed long into the future; changes later may be very costly or simply impractical.

By definition, windows are not sturdy partitions. Avoid allowing nonsecurity personnel to pressure architects into including large windows, allowing visibility and possible intrusion into secured areas. Allowing uncontrolled visibility into operations centers may allow malefactors to learn more than they should about personnel, procedures, and schedules simply through observation.

Security doors should be sturdy also. They should be metal, fire-rated, and relatively heavy and use heavy-duty hardware. Try not to call attention to controlled doors by any distinctive external appearance. If many occupied areas use wood doors, the security doors should not stand out. Use wood-faced metal doors and hardware that looks similar to all the other doors. Sometimes secure-looking dummy doors that lead to nothing important are used for deception or as an alarmed honeypot to draw troublemakers away from secure areas.

Well-constructed partitions and ceilings will also seal out smoke and contain smoke in the event of an interior fire. Weatherstripping around the perimeter of each door is recommended to keep out dust, contaminants, and humidity and to trap any smoke.

If possible, wiring that must run inside a door should be routed within the hinges so that no wires are visible at any time. Exposed wires can be damaged by accident by cutting, or compromised in many ways. The major hinge manufacturers can supply special hinges to conceal most wiring and match the appearance of other premises hardware.

Masonry partitions are usually unnecessary unless there are special structural or very high security requirements. Drywall partitions with metal studs are usually sufficient but should be extra sturdy. Type-X fire-rated drywall panels at least three-quarters of an inch thick are recommended. Better yet, use double half-inch- or five-eighths-inch-thick panels. Existing drywall partitions can easily be double-paneled for added security and strength. Masonry partitions, especially, and often drywall as well are usually faced with sturdy fire-rated plywood for attaching equipment supports and wiring. Whether the plywood is mounted on stand-offs so that wiring can be run behind it is a matter of preference. Usually, it is more efficient to run all wiring exposed inside a secured area.

Do not use a suspended (hung) ceiling in any equipment room. Suspended ceilings just add cost, inconvenience, and diminish the volume of the room. The plenum space above the suspended ceiling is a fire hazard, and everything inside it must be fire rated, including any exposed cables, which adds unnecessary costs. Most building codes require separate fire detection zones and suppression heads within every plenum, which is a major cost. Remove any suspended ceilings and get more useful space for cables, better air circulation, and easier maintenance.

Avoid raised floors for the same reasons. They are very costly, functionally unnecessary for most installations, and take up extra floor space for ramps. Like suspended ceilings, raised floors create a plenum space that needs separate fire detection zones and suppression heads. Raised floors usually restrict the ceiling height because the space was not designed for this, so the room must be enlarged to accommodate everything,

23 · 26 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

and often with little expansion provision for the future. Raised floors soon become dust traps and, in time, usually a clutter of new, old, and abandoned cables that no one can figure out. Many equipment failures have been caused by overheating due to airflow restricted by too shallow a raised floor or by too many obstructing cables. A raised-floor plenum is rarely needed to supply cooling air. Surface-mounted ducts and registers usually can do the job better and much cheaper, and are easily cleaned and modified. All of the wiring can usually be routed efficiently above and between the equipment or by using inexpensive cable troughs mounted on walls and ceilings. If needed, floor outlets can access trench ducts in the floor, which may already exist, unused. Conduit can be installed through the floor slab and along the ceiling below for special needs.

It is important for all wiring outside of equipment rooms to be protected inside of metal conduit. This conduit should not be thin-wall or plastic but rugged, heavy-duty metal. Thick-walled metal conduit is strong, harder to damage or breach, and provides good magnetic and electronic shielding as well. Metal conduits can easily be, and should be, well grounded. Obtain expert advice on where and how to connect the grounds to avoid interference. Metal conduit also may serve as a proximity sensor to warn when something gets too close. Alarm wire concealed within the conduit offers early warning of trouble and often can pinpoint where it is located. Sometimes an inert gas or gel is used to pressurize conduits to protect the wires; a pressure drop indicates a leak or a breach but does not indicate where it may be.

All conduits should look alike, whether they carry power or data or control security systems. Although the diameter of the conduit must vary, the general appearance should be the same. Do not label or mark the outside of any conduit except with an alphanumeric code. Cables inside of a conduit should not be marked either, except alphanumeric. Generally, any wires within conduit should emerge only inside of a secure closet, equipment room, or junction box, where the wires should be labeled with appropriate codes that can be deciphered using up-to-date electronic summaries available only with authorization.

Data cables and wiring are fragile, whether copper or fiber optic. Any undue pressure, bending, stretching, vibration, or crimping can alter transmissions or cause failure. Fiber optic cables are especially fragile; metal conduit is usually the best and least expensive protection. This avoids special sheathings that are often cumbersome and costly.

Any wiring, whether metal or optical fiber, can be improperly specified, installed, or terminated. Substandard wiring or installations may function well temporarily, but future failure is likely, possibly hard to locate quickly, and will be costly to fix. Therefore, it is important that all cables be acceptance-tested and certified by an independent expert before the installer is fully paid.

Critical cable runs should be equipped with alarm facilities from end to end, whether the conduit carries power, data, or security. There are several ways to alarm a cable run. Outdoor conduits are often pressurized with nitrogen to keep out humidity, or with a special gel to keep out oxygen and humidity and to stabilize the wires inside. Interior conduits can be pressurized and alarmed in the same way. Monitoring the pressure provides an alarm when trouble starts (including failed seals that must be fixed quickly). The system is effective and provides early warning, but breaches cannot be pinpointed, and any future wiring changes may be difficult.

Proximity and pressure-sensitive sensors also can alarm the entire length of critical cables. A monitored run of conduit can be very long and may continue through areas that are difficult to protect or offer no concealment. While surveillance and intrusion detectors can protect most vital areas, there is often much infrastructure that can be

FACILITIES DESIGN 23 · 27

protected only by sensor wires running the full length of the conduit. Mechanical pressure sensors will detect unusual vibration, while proximity sensors indicate a change in the magnetic or electronic environment surrounding the cable or conduit. Newer systems utilize fiber-optic sensors that monitor sound pressure. Some of these systems are smart enough to distinguish routine, harmless events from possible trouble, and many can roughly pinpoint the location as well. Sometimes the conduit itself is the sensor, or an external wire is attached to the conduit, but these approaches are often ineffective.

23.8.4 Using Existing Premises Alarms. Various codes require workplaces to have visible and audible fire alarms. And most workplaces have voice paging, emergency intercom, surveillance, and premises security systems as well. All of this equipment can be utilized effectively to augment and support information infrastructure security. Audible alarms are used when persons at the scene must take immediate action, such as to lock down or to evacuate. Conversely, silent alarms provide early warning and allow security to monitor the scene discreetly, gather evidence, and respond and assist as needed. All these alarms can be integrated into infrastructure protection systems to provide better early warning and extra layers of protection.

All alarms and alerts should be transmitted to a central security desk or guard station. The purpose is to document and manage incident response, summon assistance quickly, monitor the scene, accumulate evidence, and support all of the response resources. Central management is especially necessary when threats cascade or multiple events occur, as they often do. Security managers, IS managers, the infrastructure security manager, and key department heads also should be notified immediately. Some of these people may be offsite or at remote locations but will need to communicate effectively with at least the operations center. Notifications and the subsequent communications should be quick, private, secure, and logged to document the events. One or more online backup security control and operations centers provide redundancy, support and assistance, and strong security.

An effective method of premises-wide alert uses voice codes broadcast over a paging system. These are usually scripted and often prerecorded so that alerts can be initiated automatically, remotely or manually. Hospitals do this effectively with their color-named codes, equivalent to silent alerts, which do not seem unusual to the public. An effective system of alert codes in a large organization also uses the names of fictitious persons. In a smaller setting, such as a school, where the names of all personnel are known, there can be alert messages to an individual to take an innocuous action, which is understood to be an alert code. Additional alphanumeric information in a message can identify the general location of an incident. As in a hospital, it is well to add similar codes for other routine purposes, so the public will generally tune out all the paging. This system is particularly useful when violence is threatened or has erupted.

Although most security personnel now have portable radios, there may be many areas of no reception, and few use earpieces so others cannot readily hear what is happening. However, all security personnel must know immediately when trouble looms, and they must be alerted in a way so as not to excite others. Also, everyone inside the premises needs to know when an emergency threatens. Indeed, everyone has a legal right to know and to promptly receive instructions for their own safety. Anything short of this will result in considerable liability. Therefore, effective procedures, clear simple instructions, good preparation, and periodic training can protect everyone and provide strong security.

23 · 28 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

The widespread availability of mobile phones has added another channel for emergency notification. Employees can be notified of emergencies using the Short Message Service (SMS) for text messages. Using computer-generated messages or even manual messages, security personnel can quickly and easily reach many of the people involved in emergency response. If phones are configured to use a characteristic tone for such messages, and if the phones are not generally used for casual messaging, such arrangements may be useful adjuncts to other notification methods.

23.8.5 Clean Electrical Power. Protecting electronic equipment requires a source of electrical power that is consistently *clean*. Clean power usually means alternating current with tightly limited variations of amplitude and waveform. Power outages and brownouts (low-voltage states) can cause obvious trouble, but numerous other disturbances can disrupt or damage the information infrastructure. Some of these include dips and sags, spikes, transients, and magnetic or radio-frequency interference. Most of these are intermittent and not necessarily present when the circuit is tested. Understanding each term is not as important as knowing that a wide variety of problems commonly occur in power lines, randomly and without warning.

Brownouts are particularly harmful. These are voltage reductions by the electric utility that can cause air-conditioning equipment, cooling water pumps, and ventilating systems to malfunction or to shut down. The associated equipment may be damaged unless it is quickly shut down or switched over to an uninterruptible power supply or a motor generator.

Few power disturbances will destroy circuits or crash systems immediately, but most can cause cumulative damage. Each incident can weaken electronic circuits that will eventually fail for no apparent reason. Poor equipment quality is often blamed, because the cumulative effect was not recognized. Worse, replacement equipment will probably soon fail also.

Power disturbances can be measured to determine whether a particular circuit seems to be clean. Usually, a recording device is left in place continually for at least a week to measure and log the details of every event that occurs. (Avoid any test instrument that merely logs an unnamed event but provides no details.) An independent engineer who is not a vendor may best provide testing that is objective, comprehensive, and covers the entire facility.

Some circuits are likely to show intermittent disturbances caused by something nearby or within the building and sometimes by faulty wiring. Knowing what power problems arrive via the main service helps to determine whether the utility is at fault and to isolate where in the building other disturbances may originate. There may indeed be numerous causes of power disturbances, and all of them may be intermittent, which is why continuous seven-day, 24-hour monitoring is the minimum recommended.

Another major electrical problem is improper grounding that can damage sensitive equipment and cause interference in cables. Electricians must comply with national and local electrical codes, but they do not necessarily understand or provide the special grounding necessary for sensitive IS platforms. Most heavy equipment manufacturers require that each unit have an isolated ground connection with a dedicated wire all the way back to the central building ground bus. A few manufacturers do not provide installation specifications unless asked, and some installers disregard them to remain price competitive.

Opinions vary as to the best building ground configuration for information systems, and as local conditions can vary significantly, no one approach is best. It is wise to consult an independent engineer to inspect the grounding configurations, and to

FACILITIES DESIGN 23 · 29

recommend and certify local code compliance. It is also important to provide separate circuits for all IS equipment, and where the equipment plugs into a receptacle, there be no more than a single receptacle. Separate circuit breakers are usually required for all equipment that can draw high current, especially if the load may cycle on and off. This is required by code for large motors, such as pumps, air conditioning, and elevators, whose cycling can cause dirty power on other branch circuits. But copiers and large laser printers (especially older ones) can also create electrical disturbances when starting and when the fuser-heaters cycle. All types of lights can cause a dip or a surge when many fixtures are switched on or off at once. (The newer fluorescent light fixtures with electronic ballasts conserve power and cause much less interference.) A separate circuit connection somewhat isolates the hot and neutral wires from other circuits, but interference may be generated through the grounding connections that often are daisy-chained with many other circuits to cut costs.

Do not share a dedicated circuit with any of this equipment, which can readily disrupt and damage electronic equipment: time stamps, electric staplers, coffeepots, refrigerators, heaters, fans, or any other device with a motor or solenoid. Even if inaccessible, a dedicated outlet should be a single receptacle, not a duplex. It is all too easy for vacuum cleaners, power tools, or maintenance equipment to use the other half of the duplex outlet. This will cause a severe disturbance and may trip the circuit breaker. There must be plenty of convenience outlets that are readily accessible for all noncritical needs.

Yet another cause of power problems is excessive solar activity. These events can be measured only when they occur, which is randomly during an unpredictable interval of several years that peaks about every 11 years. Solar disruptions occur only during daylight hours. High solar activity occurred in 1988, causing major power outages and radio interference in Montreal, Canada. Daily solar activity reports, forecasts, pictures, and historical records are available at www.dxl.com/solar and other sites. (See also Section 22.6.7.)

There are several remedial options when power disturbances are suspected, encountered, or even possible.

1. Eliminate the problem at the power distribution panel. Better grounding, more separate circuits, suppressors, filters, and isolation transformers may help. But this type of remediation can be difficult, costly, or unreliable.
2. Use a surge suppressor near the equipment it protects. This is inexpensive but useless against brownouts, outages, and severe disturbances.
3. Employ an uninterruptible power supply (UPS), which provides battery backup power, for each piece of critical infrastructure equipment. One type of UPS is activated only when the incoming power fails, although its battery is always being charged. A better type is always online, acting to filter the incoming power, suppress most surges, compensate for minor brownouts, and maintain full system power for five to 10 minutes following an outage—enough time to allow an orderly system shutdown. A third and best type of UPS always powers its load from its batteries, thus isolating the load from the power line and providing optimum protection.
4. For systems that draw high wattage, a motor-generator (MG) set will eliminate most power problems. An MG set is an electric motor driven by the utility power. The motor is coupled to a generator that will always supply clean power to the load. The generator is usually voltage-regulated automatically, although the

23 · 30 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

frequency can vary, which may disrupt some timing circuits. Usually, if there is a power outage, mechanical momentum of the unit will provide sufficient power for an orderly equipment shutdown. Motor generators are still used for the ultimate in filtration and regulation, but there must be an electrical bypass to facilitate maintenance.

UPS units should power and protect the servers, network and telephone equipment, computers, and critical monitors. Most UPS units provide outlets with no backup power but with noise suppression for printers, transformer “bricks,” fax machines, and peripherals that do not need to be kept running. Most of these devices are somewhat expendable and quickly replaced by spare units if one is damaged.

Individual UPS units, placed near the equipment they protect, cost less and can be powered off by the operator to better protect equipment that may be vulnerable, even when the equipment is already shut down. This can provide an extra layer of protection where lightning might strike. Larger UPS units are used in equipment rooms where they can also monitor and log all power events and trigger remote alarm indications. Most good UPS units can initiate an automatic equipment shutdown when the power fails, the UPS batteries are low, or someone intervenes manually. There can be an issue, however, when some of the protected equipment cannot be restarted until the UPS batteries are fully recharged. Some UPS units are also network devices that can report their condition to a remote location.

Many UPS units also provide telephone and Ethernet line filtering and suppression, and these features should be used if possible. Lightning transmitted over communications wires can readily damage telephone instruments and modems. Power and communication line spikes can occur asymmetrically, and can devastate equipment that one disturbance alone would not damage. A good UPS unit with communications line suppression is best able to stop both types of spikes.

Another benefit of UPS units is that when an emergency generator is on line, the electrical power is usually much dirtier than normal. The extra filtration and stabilization provided by the UPS units may be the difference between having IS equipment operating or crashing.

23.8.6 Emergency Power. Most of the critical systems and infrastructure must remain fully operational during any electrical power problem. Filtration and suppressors cannot compensate for power outages, handle most brownouts, or cope with major electrical interference. Disturbances might come from lightning (even when it is too distant to be seen or heard), severe solar activity, or radio-frequency interference corrupting the utility power. UPS units can deal with some of these conditions, but only briefly. Therefore, a backup emergency generator may be the only way to continue operations during sustained power problems.

Although backup generators are often the only alternative, they are not a panacea. Generators are expensive and complex to install, require at least monthly exercise, and are not always reliable. Their voltage regulation is marginal; during sudden load changes, the output voltage and frequency may fluctuate as well. As the load increases, more current is drawn, and, if the generator is overloaded, the voltage will drop and the frequency may drop below 60 Hertz (which can disrupt IS timing circuits). Because the load current increases to meet the power demand, the amount of heat generated by the equipment being powered will increase as the square of the current. Much of the IS equipment powered is inductive, and there can be a large starting power surge when it is turned on or restarted. Generators, therefore, must have ample reserve capacity

FACILITIES DESIGN 23 · 31

beyond the anticipated equipment loads. And given their cost, generators should have ample reserve capacity for the future as well.

Another issue is whether a particular generator can provide sufficiently “clean” electrical power to operate IS equipment as well as power for the other emergency needs of the facility. Be sure, therefore, that the generator specified has ample capacity and that it is intended for use with electronic equipment. Even then, interference from large motors or lighting systems can affect the electronic equipment.

Because backup generators are expensive and complex, planning is often shortsighted, and many installations are not well designed or adequately tested. The inevitable result is that many generators do not perform as expected. Here are a few examples of what can be overlooked.

After considerable discussion at a major money-center bank in Manhattan about what seemed to be the excessive cost of a backup generator, the project was begrudgingly approved. The generator was to power two identical computer systems running in tandem to support a critical securities trading operation. Because the generator would actually cost more than the two computers, cost was an issue until the bank realized that the generator would pay for itself within one day of use. Soon after completion, a sudden citywide blackout erupted and the outage lasted for three days. Despite much inconvenience carrying fuel from the elevator up a flight of stairs to the rooftop generator, the unit performed flawlessly—one of the few systems in Manhattan that did.

Many other generators did not start or cut over properly, despite warm, clear weather conditions. And others did not support the necessary infrastructure. Some installations did not think to include the power requirements of HVAC, so the computers had to be shut down within a few minutes to avoid overheating, even though there was ample electrical power for them. Generator power for other necessary support functions was neglected. These included network components and communications systems, lighting for users, an elevator for access by key people and to carry fuel to the generator, security and access control systems, and at least basic electrical provisions for food and rest for those keeping the vital systems running. Very few businesses thought to include all of the necessary support functions on their emergency power systems. This incident happened some years ago when power outages were considered very unlikely in Manhattan. Generators then were somewhat of a status symbol. But today, sudden blackouts anywhere are far more common.

A related example of shortsightedness occurred in a large suburban luxury hotel operated by a prestigious hotel chain. Following a severe thunderstorm, the power utility advised the hotel that they must lose power for several hours to repair a damaged substation. Given ample notice, the hotel set out hundreds of candles in the corridors, dining, lounge, reception, and pool areas, and started their emergency generator, which then cut over automatically as soon as the blackout occurred. Emergency exit signs and emergency lights in the corridors and stairs all worked properly. As expected, their batteries soon died but the candles functioned long and well. The generator also powered one elevator, the computer and telephone systems, the access control system, and all the room locks. The generator performed as expected, but the emergency response process did not.

Even with ample warning, no one thought to shut down the other elevators or to post signs to use the stairs. Two very frightened people were trapped in the dark, between floors, proving that a generator can be a liability and not a benefit unless operating procedures are carefully planned, well implemented, and periodically reviewed. There should have been a security checklist used whenever the generator started.

23 · 32 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

Another recent example involved a state's emergency operations command center, designed to remain safe and fully operational no matter what events might occur. A large generator powered all the critical systems. Everything had been tested many times and had operated smoothly as expected. But then trouble came during a heavy thunderstorm in the vicinity. Electrical power for most of the city flickered several times and then returned to normal. However, the generator tried to start at the first sign of trouble and then faltered as the power returned. A few seconds later when the power again flickered, the generator system had been damaged and was unable to start. The state was lucky that the generator was not needed then, but it was out for several days for repairs.

Most power failures begin with flickering and momentary outages, which can incapacitate a generator system that is not set up properly. Most mission-critical generators are set up to start the engine automatically, and many transfer power automatically as soon as the generator comes up to speed. Manual start-up and transfer are more reliable and cheaper, if trained personnel are always available. The best way to sequence automatic operation follows.

1. After the first start-up signal, the start-up sequence must continue until the engine starts, a failed-start timeout occurs, or the sequence is terminated manually.
2. Power does not transfer until the generator is fully up to speed, at a reliable operating temperature, and the utility power is unusable. All three conditions should occur before transfer, and there can be manual overrides as well.
3. All transfers back to utility power and the generator shutdown should be done manually. It is best also to be able to transfer each circuit individually to utility power.

There are countless examples of critical backup generators failing to operate as expected. Here are some suggestions to determine whether a generator is necessary for protecting information systems and how to utilize a generator efficiently and economically.

- Investigate the outage history of the utility feeders that serve the premises. The electric utility can usually provide this data; if not, the state's Public Utilities Commission usually can. Be sure to ask how the terms are defined, because an "outage" may only include interruptions that continue for more than several minutes. Also ask whether more reliable feeders are available. Loop feeders that are powered from both ends are more reliable and often serve critical equipment. Ask whether the distribution transformers isolate and filter out power disturbances, whether they can also regulate the incoming voltage, and, if so, the specifications.
- Find out which other customers share the same feeders, and visit them to discuss their experiences and to determine if they use heavy machinery. Although some safeguards are possible and may be at little or no cost to the utility customer, past history is not always a reliable guide to the future. The distribution grid changes as more heavy loads are added. Today, the threat of extended power problems is far greater than in the past and is increasing rapidly. UPS units, motor-generator sets, and backup generators may all be a necessity in mission-critical applications.
- Determine which of the IS infrastructure components need backup power from an emergency generator. Most critical information systems, equipment, networks, and infrastructure must be at peak performance at all times. And so must all the

FACILITIES DESIGN 23 · 33

office areas, support systems, utilities, and personnel needed to operate them. Outages can drag on for days or weeks with key people isolated and living inside the facility to keep the systems running. The generator power must serve all of these needs.

- Consider these support systems that may require emergency power:
 - All the IS security systems, protection and monitoring devices; perimeter surveillance, and access control systems, the security stations and consoles.
 - Fire stairs (which may become the primary means of entry and egress), emergency exit doors, fire alarms, and intercoms whose batteries will quickly discharge. Also the need for these batteries to begin recharging immediately as soon as backup power is available.
 - Heating, ventilation, air-conditioning (HVAC), and process-cooling systems, including all the controls, fans, pumps, and valves needed to operate the critical and support systems. In addition to equipment cooling, it is best to provide room comfort for users, operators, and administrators. Area air conditioning may not be possible, but at least supplementary heating in winter and adequate ventilation will be needed.
 - Sufficient lighting for key personnel, equipment rooms, utility closets, corridors, rest rooms, and food service. Many individual light switches can conserve power and generator fuel. Battery-powered lights are suitable only for immediate emergency egress and cannot provide area lighting.
 - Enough live convenience outlets for test equipment, work lights, and any accessories that must be used. Live receptacles may also be needed for portable fans.
 - Sufficient food service equipment and refrigeration, running water and sanitary facilities, and a sleeping area for 24/7 operations that may have to continue for several days.
 - An elevator for critical access to the site, for medical emergencies, delivery of food and supplies, and to carry fuel for the generator.
- Compile a list of all the items a generator must power. Then total the rated power of each item to determine the size of the generator and the number of circuits needed. Power ratings for equipment usually are shown on a nameplate near the power connection and listed in the instructions. Ratings may be given in watts, amperes, or volt-amps. Generally, watts and volt-amps are assumed as equivalent. The latter value is the product of multiplying the rated voltage (e.g., 120 volts) by the rated amperes, while the former multiplies that number by the equipment's power factor. Large generators are rated in kilowatts (1,000 watts) of power. Units intended for short duty cycles cost less and may fail during prolonged, continuous duty. An experienced engineer should review this process.
- Consider the costs per hour of lost productivity, business, and goodwill if any information systems are disrupted. Add to this the recovery costs. In the example of the bank given earlier, the first day that the generator was needed saved the entire cost of the backup power system. The second and third days of that particular outage were sheer delight to the bank as most of their competitors faltered.

Electric codes may require, and good practice dictates, that a generator be sized to handle the sum of the individual start-up loads. This may seem wasteful, because not

23 · 34 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

all loads start up at once and an average operating load will be somewhat less than the sum of the parts. It is nonetheless a wise practice to provide for the maximum rated load, with additional spare capacity for improved reliability and future additions. There are several reasons for oversizing the generator. When power is first transferred to the generator, the sum of the initial surges can far exceed the anticipated load. All of the UPS units and other battery-operated devices will begin to recharge, and all equipment motors may concurrently draw their maximum surge currents. Extra generator capacity ensures a smoother transfer with better voltage regulation and enhances the system's reliability.

Most large generators produce three-phase power. And each of the three outputs should be balanced so that each "leg" draws about the same power. To do this, heavy motors, multiple light fixtures controlled by one switch, and other surge-producing equipment may have to be divided among the three legs. Existing wiring at the distribution panels probably will need changing to balance the legs. It is desirable, but not always possible, to reserve one leg for clean power to critical single-phase electronic systems. Balancing each leg is a tricky business best done in consultation with an independent engineer.

As many electronic systems as possible that are powered by a generator should also be protected by UPS units as discussed earlier, even though these add to the generator's load. There will be large voltage surges, dips, sags, and over-voltage conditions as the generator loads are switched and constantly change. Power disturbances will be much greater because electrically noisy motors and lighting cannot be isolated. The UPS units should include noise suppression and voltage regulation as well. And even with all of this, IS equipment will be stressed and may fail.

Locating the generator is the next challenge. The choices are on a roof or setback, inside the premises, or outdoors. Each site has advantages and obstacles. Outdoor generators can be the easiest and cheapest to install but also more expensive to operate. Outdoor generators are noisy, often unsightly, subject to vandalism, and local ordinances may restrict them. When located outside, weatherproof housings are needed to protect the engine, generator, and fuel tank. Most engines used outdoors need to be kept heated, which can become a high overhead expense. Noise is another problem, and persons nearby may object. It is important to use good mufflers and to get written permission from nearby property owners and other tenants. Outdoor units should be fenced with plenty of room for maintenance and fueling. A generator shed is best, if possible, but this does not reduce the need for heating and a good muffler. The whole installation should be securely locked and protected by an open-door alarm and motion detectors, and be in the view of surveillance cameras. Floodlights may deter vandalism and will assist refueling.

Generators on roofs or building setbacks present other problems, and these installations too may be restricted by local codes. The first problem is weight. Structural reinforcement probably will be needed. The next problem is getting the unit in place, which may require a crane or a licensed rigger. Very few building elevators come up to the roof level, and they may not be able to handle even a disassembled generator's parts. All the generator components may have to be rigged up outside of the building or manhandled up the fire stairs.

Installations on top of building setbacks will need a special access door, and moving heavy equipment across a finished floor requires heavy planks and floor protection (e.g., sheets of Masonite or plywood) under the casters to avoid considerable floor damage. There must be sufficient space on the roof or setback to fuel and service the generator safely. Noise will usually be a problem and vibration as well.

FACILITIES DESIGN 23 · 35

Indoor installations offer both advantages and challenges. An indoor location that is sometimes feasible is a heated garage-type ground-floor room with a standard garage door to open when the generator operates. This arrangement is good because it is inconspicuous, fireproof, easily protected, and convenient for fueling and maintenance. And, should a generator fail, a trailer-mounted unit can be hooked up easily.

Inside generators may be prohibited by building or fire codes. Large rooms are needed to facilitate fueling and maintenance, and large ventilation systems to dissipate the considerable engine heat. The engine exhaust can be well muffled and piped outside, while engine-intake air is ducted in from outside. Heating and ventilating the room must be designed correctly for both very hot and very cold weather. The room must be fireproof and soundproof with fire alarms and a suppression system that uses chemicals or dry-head sprinklers that cannot freeze. The floor may need reinforcement and vibration isolators. A floor drain is advisable and must be environmentally approved.

There are advantages to indoor installations. The generator and its fuel can be kept warm easily. Starting is easier and more reliable. Fueling is easier without having to brave the elements. There is less chance of water in the fuel, which can be fatal to diesel engines and maintenance is much easier.

Problems with building installations include building codes that allow only small day tanks for fuel. Every few hours, a lot of fuel must be carried in to keep the generator running. Fuel cannot be stored inside most buildings, and an elevator may not be running or available to help bring in fuel cans.

There are many possible fuels for emergency generators. Diesel fuel is the most efficient, and diesel engines can operate continuously for days but are hard to start, especially in cold weather, and cannot be hand cranked. Home heating oil is basically the same as diesel fuel and can be substituted at any time that diesel fuel is not available, but this requires extra fuel filtering.

If liquid fuel is used, the fuel tank should be full at all times to avoid condensation. Fuel additives can prevent gumming and assist starting. Make sure all diesel fuel is treated for use in a cold climate. Refiners normally do not use this process except in winter, but untreated diesel fuel turns to a gel near freezing temperatures and the fuel will not flow. Never let a dealer "cut" diesel fuel with kerosene, which is corrosive. Diesel fuel also requires additives to avoid bacteria buildup that will clog fuel lines. There should be OSHA-approved cleanup materials ready for any future spills or leaks.

Natural gas or propane are the most convenient fuels. Either one eliminates the day tank and refueling. These engines are the least polluting, and they start much easier, require no preheating, and can be hand cranked. But most are not designed for prolonged continuous duty. Gasoline engines are prohibited by many building codes and are rarely used except for small, portable generators. Gasoline is far more dangerous to handle and store, and gasoline engines do not hold up well under heavy loads.

Continually monitor the engine oil level and be ready to add oil as soon as it is needed. Most generators automatically shut down when the oil level is low. Some also shut down when overheated. Any unexpected generator shutdown will be catastrophic, so monitor closely for early warning signs of trouble.

Once the desired size and type of generator is decided, there are other considerations:

- Automatic engine controls and load transfer switches can be unreliable and may cause damage. Avoid these if possible. However, generators can be monitored and controlled remotely, as well as on site.

23 · 36 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

- Automatic starting can be unreliable. If the engine does not start quickly, the battery will quickly discharge, especially diesel engines in cold weather, which require glo-plug heaters. If at all possible, someone should be present during the starting process, using a checklist to verify proper operations and then transferring the load manually when the generator is ready. Switches that automatically transfer the load are expensive and sometimes fail. Always transfer back to utility power manually, and do this only after sensitive systems are put into a standby mode. Automatic transfer can cause major damage if the utility power flickers and goes out again or if the voltage or the frequency fluctuates during transfer, as it often does. Do not shut off the engine automatically. This is best done manually, and not until utility power is flowing smoothly.
- The best transfer switches allow each of the major circuits to be transferred individually to minimize the inevitable fluctuations likely to occur when everything is switched over simultaneously.
- An emergency generator must be exercised regularly. The manufacturer will specify when and how the units should be exercised. Usually, this must be done monthly and at medium to heavy load. When critical systems are involved, good security practice is to exercise the generator weekly. There should be a written, initialed log entry for each event, including each exercise, inspection, maintenance, oil check, and refueling. Always log operating hours.

Despite the cost and complexity, there is a great feeling of contentment in having a good emergency generator system that functions smoothly, especially when other organizations may be floundering. Once the generator performs well during a real emergency, even skeptics realize the value added.

23.8.7 Environmental Control. Even though today's information systems do not need as much cooling or the precise environmental controls that legacy systems once demanded, good control of temperature and humidity, good ventilation, and clean air are still important. Information systems can function reliably only when proper environments are maintained in both equipment rooms and user workplaces. But each area requires a different approach.

Air conditioning is basically intended to cool people; equipment should be cooled by a functionally different system, which is best called process cooling. The systems should not be intermixed, nor should either one substitute for the other. Building codes require HVAC within all occupied spaces, where people may congregate, or where there are workstations. Building codes also set minimum ventilation requirements for occupied space, including a minimum percentage of makeup (outside) air to be constantly brought into each occupied space so the inside air does not become stale. Most codes do not consider the needs of electronic equipment.

Electronic equipment has many special needs, and many are incompatible with the people comforts required by the codes. Most electronic equipment operates continually, whereas air conditioning operates mostly during business hours. Air-conditioning cooling systems may be shut down for maintenance, during a power brownout, off hours, or in cool weather. By contrast, process cooling must operate continuously and every day, so parallel and redundant systems are often used. The same air should be well filtered and recirculated with no makeup air added to introduce dust or contaminants. This also reduces the cooling capacity needed, so process-cooling equipment can be of smaller capacity and cheaper to operate.

FACILITIES DESIGN 23 · 37

Electrical equipment and wiring also need good humidity control, which process-cooling systems are designed to provide. These systems are designed to be easier, and faster to clean and maintain. Often many components are redundant and hot-swappable. Increasingly, the cooling unit is on the floor or ceiling of the equipment room, so that few ducts, dampers, or registers are needed.

All IS processing, storage, and network equipment should be inside dedicated equipment rooms, which also should be designated as unoccupied spaces to avoid the code-imposed air-conditioning requirements. Avoid using terms such as "computer room" or "data center," which are usually construed to be occupied spaces.

Both the process cooling in equipment rooms and the air conditioning in work areas must provide humidity control. It is important that relative humidity be controlled between 40 and 60 percent at all times, regardless of the climate or season.

When the relative humidity falls below 40 percent, which can easily happen in cold weather, static electricity is generated as people move about. Static charges can quickly accumulate to become many thousand volts, and a spark will jump to any object a person touches that is differently charged. Even though such a spark may not be felt, several thousand volts can annihilate electronic circuits. For example, a static charge jumping from a fingertip to a keyboard can cause serious damage to storage media and circuits. Much of the damage may not be readily apparent. Actual failure may be delayed, so the cause is not identified. Grounded strips can be installed on workstations, and service personnel should wear grounded wrist straps, although these do not completely stop the problem. The only effective solution is always to keep the relatively humidity above 40 percent.

Relative humidity above 60 percent also causes problems that will eventually destroy equipment and wiring. Above 60 percent, condensation and mold will begin to damage some components. Above roughly 80 percent, galvanic action occurs and will eventually cause serious trouble. The process is often called silver migration because most electronic connections are silver-plated. The phenomenon is similar to electrolysis (electroplating), but here the two metals are immersed in high humidity rather than a liquid and there is no external current needed for galvanic action to occur. Molecules of one conductor begin to physically move toward and attach themselves to another less active metal. Even though both surfaces may be gold or silver or copper plated, it is likely that they differ slightly in composition. Therefore, galvanic action will occur whenever the humidity is too high. Connector pins and sockets can disintegrate, fuse together, or fail electrically due to pitting. Printed circuits can also fail. Although this galvanic action happens slowly, it accumulates and is irreversible. The failures are usually without warning, and almost always, poor quality is blamed, rather than high humidity.

The only protection is to control humidity in both equipment rooms and work areas. Process-cooling and air-conditioning systems commonly do this by several methods. Both systems dehumidify naturally when cooling and can use a reheat coil to warm output air if the humidity is too low. Also when the humidity is too low, water is added using a spray, atomizer, or a wet screen through which the supply air is pumped.

There are additional protections, which are wise to install and maintain. In a cold climate, all areas and workplaces with electronic equipment should have low-static floor surfaces. This can be low-static carpeting or floor tile made for this purpose. Do not rely on sprays to control static electricity; they soon dissipate. Be sure that equipment room walls and doors are well sealed so that humidity, dust, and contaminants cannot migrate. Be sure the walls are well sealed from slab to slab and that the slabs themselves are impervious to dust and humidity.

23 · 38 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

23.8.8 Smoke and Fire Protection. Smoke and fire must be prevented within any equipment room because otherwise, considerable damage and disruption will occur quickly. No matter how small the incident, the effects of either smoke or heat are cumulative. Systems will eventually fail and usually without warning. Obviously, smoking tobacco and other substances must be prohibited—but equipment rooms can be tempting for untrained staff because these rooms may seem to be a safe, cool place to sneak a smoke—and they may even be largely unattended.

Equipment room doors should never be propped open by cleaning, delivery, or maintenance personnel, or others working outside when the air conditioning is off.

The first level of prevention is to keep everything combustible outside of equipment rooms. Paper and supplies not in actual use should be stored outside, never within an equipment room. Although most documentation is now made available electronically, reference materials or documents for emergency response when systems are down may be kept within the room; these should be stored inside fire-resistant files or cabinets for secure and rapid access in an emergency. There should be no trash receptacles within an equipment room, and shredders should be outside, under strict control. There should be a clear and firm policy that nothing combustible can remain inside an equipment room, and frequent inspections should be held to verify compliance.

Equipment and suppliers storage contents rooms should be designed to protect and not to accommodate people. There should be no unessential furniture within equipment rooms, especially desks that can become cluttered and that are not rated as fire resistant. A metal table with one small drawer and one or two metal chairs with fire-resistant upholstery are usually sufficient.

Inadequate firestops are a major threat that is often overlooked. A firestop is usually a sleeve and a special material to prevent smoke or heat from penetrating an opening in a partition, floor, or ceiling. It also stops the spread of flame. Many firestops are needed throughout the premises, including the building core and the mechanical, utility, and equipment areas—even within a one-story building. Firestops are rigorously required by most codes, but compliance is often inadequate and the devices often are breached by subsequent alterations or wiring changes.

Partitions, building walls, floors, and ceilings must all be fire-rated in accordance with the national and local building codes and other regulations. Proper construction usually is specified by an architect or engineer, and compliance is inspected or certified as soon as the construction is complete. Inspection often occurs before all of the mechanical systems and wiring installations are finished. Subsequently, if any penetration or opening is made through a wall, floor, or ceiling, its fire rating is thereby invalidated and should be recertified. The walls for all equipment-rooms must be slab-to-slab construction with no access for air transfer below or above the walls.

The Underwriters Laboratory (UL) or similar recognized authority rates and approves commercial firestops before they can be sold legally. Each manufacturer then specifies the approved applications, installation, and maintenance procedures necessary for compliance. It is therefore wise to utilize specialized vendors with extensive training and experience installing and inspecting firestops and to have them conduct periodic premises inspections and certifications.

Inadequate firestops are particularly common in the core areas of older buildings or where tenants occupy multiple floors. While proper firestops may have been provided during construction, installation of piping, cables, conduit, and subsequent wiring changes often breach them. Wires often poke through large, gaping holes hidden by a hung (suspended) ceiling or behind equipment racks or wire troughs. Proper firestops must be installed and inspected whenever changes occur. Many installers do

FACILITIES DESIGN 23 · 39

not understand this, or cut corners hoping no one will notice, or assume others will take care of it.

Without proper firestops throughout, fire and smoke can, and probably will, spread surprisingly quickly. And so will dust. There can be substantial liabilities if any people are harmed or equipment is damaged because of improper firestops or inadequate fire-rated construction. The costs, time lost, and reputational damage will be huge. Periodic and thorough fire inspections by an independent and qualified expert will quickly discover building and firestop violations.

Smoke is far more dangerous than flame. And all smoke is toxic! It contains high levels of carbon monoxide, which is invisible, odorless, and quickly fatal. Smoke is the product of the combustion of many materials, and most of these are dangerous to breathe. Some are immediately fatal. Even a little smoke can do considerable harm to humans and much harm to electronic equipment. Smoke is deceptive; even when there does not seem to be very much smoke or heat, and visibility looks good, people within or passing through the area quickly become disabled and some may soon die.

The first priority is the safety of people. Get everyone away from any smoke immediately, and keep everyone away. Only trained responders with proper protective clothing, equipment, and self-contained breathing apparatus should enter any smoky area. Generally, respirators are not enough protection and may leak as well. There must be no heroics; crawling through smoke on the floor or breathing through a wet rag are desperate measures that should be attempted only when unavoidable to escape the area. Everyone should wait in a safe place until firefighters arrive and then follow their instructions.

The best way to prevent an equipment room fire is to keep anything combustible outside the room. Documents, manuals, and supplies should be stored outside the room in closed metal cabinets. Inside furniture should be limited to a metal table and a chair or two. All waste receptacles should be outside. Once combustible materials are eliminated, the only smoke that develops will be from electrical overheating. Electrical fires rarely occur in an equipment room, and those that do occur are likely to be very small, brief, and cease as soon as electrical power is removed. (Note that most computing components now operate on five volts or less, so that a short circuit is no more dangerous than, for example, a shorted flashlight, which presents no smoke hazard.) While sometimes noticeably acrid, there is usually little visible smoke. Therefore, sensitive fire and smoke detectors and an effective means of fire suppression are needed and required by most building and fire codes. Good detectors can provide enough early warning to ward off trouble and injury.

Enough smoke or heat to cause actual equipment damage requires an electrical current higher than most components can draw. Circuit breakers and fuses usually will open before there is much smoke or damage. Perhaps the greatest risk is smoke from the ballasts in low-quality fluorescent light fixtures, which can put out considerable black smoke. Any smoke is corrosive and may condense on connectors and printed circuits, which may then eventually fail.

There should be smoke detectors in every equipment room that are connected to a central alarm system. There should be enough detectors to cover the entire volume of each room. Each detector should include an electric eye to look for haze or smoke, ionization sensors to detect products of combustion well before any are noticed by humans, and rapid-rise-in-temperature detectors in case there is enough heat buildup to cause damage. Even though detected, nothing will stop smoke generated by overheated wiring or components until the electrical power is cut off or other heat source is removed.

23 · 40 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

There must be a fire suppression system in every equipment room. Both code compliance and good security practice requires this. Fire suppression is best accomplished with sprinkler heads that spray water mist, even though some unprotected equipment may be damaged if the water is not effective quickly. Special waterproof protective covers are often kept near equipment in case of accidents such as ceiling leaks or a damaged sprinkler head. But if an area is already smoky, no one should attempt to place the covers.

Wiring, connections, and most components will dry themselves, even when soaked. The process may be hastened with lint-free towels and careful use of hair dryers. Keyboards, monitors, UPS units, power supplies, some disk or tape drives, and especially printers may be damaged and should be replaced until they can be inspected. Hard drives are usually hermetically sealed and unaffected. A few other components could be damaged by excessive heat, although water mist is very effective in quenching heat sources. Plenty of replacement items should be safely stored nearby. Handling damaged low-voltage components (such as most circuit boards) presents little risk to people—provided there is not too much water and the persons know what they are doing and how to avoid damaging the components. If in doubt, shut down the components temporarily.

Enclosed equipment cabinets offer the best protection regardless of the room's fire suppression system. Enclosed cabinets can monitor temperature and humidity, detect and contain smoke, sound alarms, and often contain systems to suppress a fire before trouble occurs.

Halon 1301 fire suppressant was once widely used in critical areas. But Halon is a fluorocarbon whose manufacture has been banned for many years. Today's chemical systems are designed differently; one example uses the FM200 Suppression Agent made by Siemens. The claimed advantage of the chemical suppressants is that humans can breathe the agent, at least while they are exiting the area. Another fire suppression system uses carbon dioxide, which is effective and less expensive, but can extinguish people as well as fires. The problem with all chemical agents, including carbon dioxide, is that they quickly mix with smoke and become very toxic. The agent itself may be safe to breathe, but the smoke mixed with it is not. These systems are also very expensive.

Regardless of the suppression system, there should be controls and a shutoff near the room's exit, but not accessible to a perpetrator. Generally, an audible, continuous alarm indicates that the suppression system is about to activate. There should be postpone buttons on the control panel, and perhaps remotely as well, that will delay activation for about two minutes while someone intervenes. The postpone mode generally pulses the audible alarm. A silent alarm indication should remain activated whenever a fire suppression system is disabled or the alarms are silenced.

The next level of protection utilizes several fire extinguishers. These are the most useful protections because the suppressant can be aimed where it is needed and not throughout the room. Carbon dioxide is best because it does not leave a residue. Chemical, powder, and foam extinguishers also work well but are hard to clean up. ABC-type extinguishers are best because they are effective for combustible materials, flammable liquids, and electrical fires, respectively. Several handheld extinguishers are better than a few large, heavy units. All fire extinguishers should be conspicuously wall-mounted or placed immediately inside and outside of entrances. An OSHA-approved red patch placed on the wall nearest to every extinguisher highlights its location. Also, check other OSHA, local code, and insurance requirements that may apply.

Supply air from the process-cooling equipment should be shut down quickly and automatically to avoid recirculating the smoke. The IS equipment may have to be shut

MITIGATING SPECIFIC THREATS 23 · 41

down soon thereafter before it overheats. It is best to shut down everything promptly and automatically in an orderly sequence—cooling, IS equipment, electrical power, and lighting—and then evacuate. Shut down the lighting, in case it is part of the problem. Shut-down should occur automatically with manual intervention from controls inside the room or remotely. Battery-powered exit and emergency room lighting are advisable so responders do not need flashlights.

A so-called crash cart is a good investment. This is used during a smoke condition, a water leak, and, it is hoped, before a fire suppression system activates. A crash cart is kept outside or nearby major equipment rooms and rolled to where it is needed. The cart usually contains covers to keep smoke and water out of racks and off equipment, large fire extinguishers, and sometimes respirators or self-contained breathing apparatus. The crash cart should include quick-reference procedures, and a checklist for protecting and shutting down the room, as well as safety and notification procedures—usually printed on plastic. The crash cart should be inspected and the procedures reviewed monthly, and there should be periodic training and exercises to practice using the equipment. Before the smoke and water covers are used, be sure the equipment is first powered off. Crash carts were important for yesterday's computer rooms but are increasingly unnecessary in a well-designed equipment room.

Finally, be sure to have smoke-exhaust systems available to quickly purge the areas of smoke. Most fire departments have portable purge fans with long fabric and wire hoses to reach outside. Do not allow anyone to use a respirator or breathing apparatus unless it is approved for this purpose and has been properly fitted to a trained person.

23.9 MITIGATING SPECIFIC THREATS. Several other threats should be considered before good infrastructure protection is possible. Some of these situations are unlikely but potentially very costly if they should ever occur.

23.9.1 Preventing Wiretaps and Bugs. Most wiretaps are placed at wiring junction points. Vulnerable spots are within equipment rooms, wiring closets, junction boxes, wiring blocks, or data receptacles. The tap wire can be fiber or coax or utilize a pair of unused conductors inside an existing cable. It is likely to be a small wire that is hardly noticeable, running to an inconspicuous place where monitoring and recording can occur. Once removed to a safe place, the data can be extracted by phone, wireless, Internet, or manually. Tapped data may even be encoded and stored on the victim's own network. Video and/or audio bugs used for spying are similar to wiretaps in that once the data is monitored, it must then be sent elsewhere for retrieval.

Unless all system data are encrypted—including all data, voice, and video traffic—wiretap protection must be strong because detection is difficult at best. First, determine which cables are critical and inspect the entire cable run. All cables should be inside of metal conduit. Data and power conduits should look similar and with no markings or labels except alphanumeric codes. Keep critical conduits as inconspicuous as possible, and away from places the public might access. There must be strong access controls, intrusion alarms, motion detectors, or surveillance where terminations, connectors, or wires can be accessed. Critical cable runs must be protected over their entire length. See Section 23.8.3 for ways to protect conduit and exposed cables.

Data cables between the desktop and wall or floor outlets are potential wiretap sites. Cables, harnesses, and connectors within office furniture systems may also be compromised. Reasonably good protection is possible with careful design, with devices that harden the data cabling against the possibility of a wiretap, and that detect disconnecting or tampering with any data wires.

23 · 42 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

For continued protection against wiretaps and bugs, even when all data are fully encrypted, there must be periodic and thorough visual inspections, sweeps for any unusual radio-frequency transmissions, and careful cable testing to determine any anomalies. Everything done must be logged and quickly analyzed. Unfortunately, most spying is never detected and can continue undetected at the will of the perpetrators.

23.9.2 Remote Spying Devices. There are very sensitive radio receivers that can monitor information system data through walls or from outside the building without the use of an inside bug or wiretap. These devices can simply listen to the data from afar. Such equipment is not available publicly and is well beyond the means of all but the best-financed spies. However, there are many such systems in use today, and many more will be available as prices drop. Any organization whose data are very valuable is a potential target. The best protection is good shielding around equipment rooms and thick-wall metal conduit for data cables, and everything must be properly grounded. There are also interference transmitters that may help; these broadcast white noise that can overwhelm signals radiated from the IS infrastructure.

23.9.3 Bombs, Threats, Violence, and Attacks. Violent events are unpredictable and potentially devastating. These are not accidents, but deliberate attacks, intended to disrupt, cause damage, and spread fear. The tragic attacks of 9/11 and their aftermath have proven the vulnerability of people and of their infrastructures. The vulnerabilities remain today, and the risks are even greater.

Protection against violence must be threat-specific, and all possible threats must be addressed as described in Chapter 22. Effective deterrence and mitigation then become a matter of strengthening the protections described throughout this chapter, which need not be very costly considering the response and recovery costs that could otherwise result. Premises or corporate security must deal with most threats of violence, but the infrastructure needs special protections to avoid disruption and to mitigate the downtime and cost consequences of any such event.

Reasonably good protection and mitigation measures can be simple and inexpensive. State and regional bomb squads or explosives units can advise and assist in many ways, including current briefings. Weapons of mass destruction (WMDs), other than nuclear weapons, are fast becoming a real threat, especially because many such devices are small and easily concealed. WMDs include chemical and biological agents and incendiary devices, while even small amounts of radioactive materials disbursed by an explosive *dirty bomb* can spread panic. The government considers these devices very serious threats, with businesses and their infrastructures as likely targets.

New federal office space must now be certified as bomb resistant, so that an explosion or the impact of a truck bomb cannot collapse the building. Officials can usually provide a current threat briefing and suggest protective measures. Although small areas may be destroyed, the structure will not collapse.

Just like some possible threats, there are also effective protection devices that cannot be mentioned. These are not marketed publicly and therefore are unknown to dealers, resellers, or distributors. The costs can be reasonable because they are only sold direct. Developers may supply classified systems for military or government use and offer declassified versions to other selected users. In this way, developers can restrict knowledge of their products to as few people as possible, so that others cannot discover how to recognize or circumvent them.

Many consultants who have worked with financial, regulated, or very large private companies know some of these specialty vendors. Usually, a consultant will approach

COMPLETING THE SECURITY PLANNING PROCESS 23 · 43

the vendor and discuss what is needed; the vendor may then contact the customer directly.

Finally, given today's environment of violence, get to know key local, state, and federal law enforcement and investigative officials, and ask their suggestions how best to protect an organization.

23.9.4 Medical Emergencies. Mitigating medical emergencies requires a first aid room on the premises, first aid and some medical supplies, a registered nurse if possible, and many workers trained in first aid and cardiopulmonary resuscitation (CPR). All security personnel and guards should be certified in first aid and CPR.

An automated external defibrillator (AED) on site will save lives and can be operated by anybody in an emergency. A portable AED currently costs about \$1,000, and the suggested training is inexpensive. An AED is now required in all federally managed buildings. Many shopping malls, places of public assembly, and commercial aircraft are now equipped with one or more units.

Oxygen is often necessary to save lives and prevent permanent impairment. Most sites equipped with an AED also have oxygen units. Good portable units cost \$800 or less.

23.10 COMPLETING THE SECURITY PLANNING PROCESS. The last step necessary to protect the information infrastructure has four components. Absent any of these components completed thoroughly, good security is not possible. They include:

1. Develop an all-hazard mitigation plan.
2. Develop all of the mitigation options for each identified threat and perform a cost-benefit analysis to determine which options are best.
3. Develop an overall security response plan to show who is responsible for what.
4. Complete the necessary implementation, accountability, and follow-up procedures.

23.10.1 All-Hazard Mitigation Plan. FEMA provides extensive resources for *Multi-Hazard Mitigation Planning* on a wide range of topics and including email updates.¹³ *State and Local Mitigation Planning How-to Guides* can be helpful in this final step.¹⁴

- FEMA Publication 386-1, *Getting Started: Building Support for Mitigation Planning*, establishes the hazard-mitigation process.¹⁵
- Publication 386-2, *Understanding Your Risks: Identifying Hazards and Estimating Losses*, shows a method of cost-estimating potential losses due to flooding using tables rather than by calculation. These tables quickly show that losses can be far greater than expected.¹⁶
- Publication 386-3, *Developing the Mitigation Plan: Identifying Mitigation Actions and Implementation Strategies*, provides guidance on developing the mitigation strategy.¹⁷
- Publication 386-4, *Bringing the Plan to Life: Implementing the Hazard Mitigation Plan*, discusses guidance on implementation.¹⁸

23 · 44 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

- Publication 386-5, *Using Benefit-Cost Review in Mitigation Planning*, provides heuristics for selecting among mitigation options to maximize cost effectiveness of expenditures.¹⁹
- Publication 386-6, *Integrating Historic Property and Cultural Resource Considerations into Hazard Mitigation Planning*, discusses how to work with local planning authorities to protect cultural capital.²⁰
- Publication 386-7, *Integrating Manmade Hazards Into Mitigation Planning*, explores how to plan for terrorist or criminal attacks in emergency planning.²¹
- Publication 386-8, *Multi-Jurisdictional Mitigation Planning*, emphasizes the importance of collaborative response planning (Federal Emergency Management Agency 2006).²²

Once all possible threat situations have been identified and assessed as described in Chapter 22, the next step is to develop as many options as possible to prevent, deter, or mitigate disruption, injury, or damage from each threat. Although some threats cannot be prevented, there are always ways to prepare for and mitigate their impact. Usually, there are many mitigation options, so the objective is to determine which options are the most practical and affordable. The only objective way to do this is with cost-benefit analysis (described in Section 23.11.2).

In actuality, the options to protect against many different threats will be similar, but each option should be retained until the best mitigation strategy for each threat is determined. Here also there is likely to be one common mitigation strategy that covers many different threats. All credible threats should be listed in the mitigation plan, but the mitigation projects laid out will be far fewer in number.

When the next step is finished, the All-Hazard Mitigation Plan can be completed. The FEMA 386 how-to manuals described above provide the suggested format and content. FEMA's *How-To Guide #9* (FEMA 386-9), entitled *Using the Hazard Mitigation Plan to Prepare Successful Mitigation Projects*, is an extensive guide to using all the other resources mentioned above.²³

The complete mitigation plan should be for official use only, and not released except to those with a need to know this information. The complete plan would be very helpful to a potential troublemaker because it shows where the organization is vulnerable, and to which threats. If disclosure of the complete plan is not well controlled, its contents could be leaked by, extorted from, or sold by an insider. However, the executive summary of the plan and abbreviated findings should be circulated widely, so that all stakeholders know that much is being done to protect them.

23.10.2 Cost-Benefit Analysis. Security is pointless unless it is cost effective and also adds value: that is, the cost of mitigating each threat must be less than the potential benefits and savings of the event not occurring, because if the protection is effective. The costs of every option can easily be determined; these are the initial and ongoing costs. Some future benefits, though, will be intangible, and all will have to be approximated. The long-term benefits of something not happening must reflect the approximated costs of:

- Disruptions that would reduce the productivity of the business and the performance of its information systems
- Morale and performance that could plummet, and remain very low because people feel unsafe and unprotected

COMPLETING THE SECURITY PLANNING PROCESS 23 · 45

- Loss of business or customers until operations could be restored to normal
- Response and recovery costs including extra time and overtime, expenses including lodging and meals, temporary facilities, public relations, and legal defense costs that are all likely to be incurred
- Legal, public relations, and other services and expenses to repair reputational damage, and fallen stock price, and to restore goodwill

Not all of these costs will follow every threat. But then again, there may also be additional, unexpected costs as well. In general, the response and recovery costs of any major security event tend to be far greater than expected. Nonetheless, each situation can be studied and some costs determined in order to facilitate a statistically valid cost-benefit analysis.

Cost and benefit information have no meaning unless each is associated with a common time frame. The likelihood of each threat should be assessed on an annualized basis (see Section 22.3.4), so that both the mitigation costs and the potential benefits can be amortized over the same life cycle.

There are many methods for cost-value analysis, and most are beyond the scope of this chapter. However, for those who are not financially trained, the federal government has a good system, freely available, that is widely required within the government. The U.S. government calls this system *BCA*, which stands for benefit-cost analysis. As part of this system, FEMA has developed a *Benefit-Cost Analysis Tool*, which assists grant applicants with financial analyses, such as net present value. The same site includes hazard data.²⁴

One particular advantage of the BCA system is that the OMB publishes current cost-of-funds data needed to project any costs. Some of the private models tend to use wildly optimistic (or grossly out of date) future interest rates, which invalidates meaningful results. Again, it is wise to use a federal government model simply as a matter of risk avoidance. The BCA model is widely used and required for many grant applications.

23.10.3 Security Response Plan. A hazard mitigation plan is needed to document the threat assessment done in Chapter 22 and to list the mitigation options and the predicted costs and benefits associated with each; a security response plan is also needed to direct each stakeholder according to the type of threat experienced. The purpose of the security response plan is to define clearly who is in charge and who will do what, when and how, when any threat occurs.

The new NRF format revises and reestablishes the Emergency Support Functions (ESF) concept for each support activity that each organization may need.²⁵ There are now 16 numbered ESFs, beginning with ESF-1, “Transportation,” ESF-2, “Communications & Alerting,” and on through ESF-16, “Animal Health.” It is recommended that the standard ESF titles be retained, even though many may not be applicable to a nongovernment organization, if only to maintain a uniformity and language that everyone understands. Additional ESFs will be needed to mount an effective response, but number these as 17 and upward.

One useful quick reference in the new response plan format is an Emergency Support Function Assignment Matrix (which is often Figure 1), a one-page graphic that shows which agency or department has primary, secondary, or support responsibility for each ESF. This is a handy quick-reference guide for management and staff when trouble comes.

23 · 46 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

Once a response plan is ready, standard operating procedures should be written by each affected department to outline the procedures it will use to respond. The term is in quotes because these are usually issued as guidelines so that there can be some flexibility to adjust to the actual conditions. The requirement for rigid adherence to a procedure is an invitation that invites litigation.

The complete security response plan should also be for official use only, and released only to those with a need to know. After all, this plan shows troublemakers just how the organization will respond. However, a press release summarizing the plan should be widely circulated so that all stakeholders know the organization is trying to protect them.

23.10.4 Implementation, Accountability, and Follow-Up. Once all the plans are completed and signed off by management and key officials, the job of implementation begins. Of the How-To Guides mentioned earlier, FEMA 386-4, *Bringing the Plan to Life*, will be helpful here.

This first implementation step is the most critical to the security-protection process. It begins with training so that everyone involved understands and accepts the procedures. There must then be periodic exercises and drills to test the response plan and to validate that the training has been effective. Plans that are not periodically tested are soon forgotten.

Every exercise and drill must be reviewed to determine what went right and, more important, what did not, and how to do better in the future. So too should every emergency response be reviewed. Documentation before, during, and after the event is important. There will be some lessons learned from each event, and these lessons should be used to update and improve the security systems to work better in the future.

It is also critical to establish accountability for the infrastructure security. There should be only one person in charge of each function. Responsibility cannot be spread among management or departments, nor can responsibility be worn as a second hat for someone with many other duties. The senior responsible authority must set up schedules for periodic review and update of the plans and procedures, training, and exercises to make sure that the security program remains current and effective.

Good security management must also include oversight. Good planning must begin in the boardroom, and the directors must also provide continuing oversight to ascertain good security. Outside, independent, auditors who are directed from the boardroom are best able to validate the current condition of the security program. The auditor's written opinion is evidence of whether the organization is fully compliant or not. The best procedure for an infrastructure security audit is suggested in Section 23.2.7.

For management's own peace of mind (and possibly as a requirement for maintaining insurance and obtaining credit), there may be periodic security inspections, testing of defenses, and some penetration tests, including deceptions to gain access done by independent professionals.

23.11 FURTHER READING

The further readings suggested in Chapter 22 will also be helpful in this chapter.

The Congressional Research Service report to Congress dated June 2006, entitled

Federal Emergency Management and Homeland Security Organization: Historical Development and Legislative Options, is interesting reading for anyone burdened with security planning and management. This outlines the federal government's struggles since 1947 to get the system right. It also lists legislation

NOTES 23 · 47

before the 109th Congress to fix the problems. The report can be downloaded free at www.fas.org/sgp/crs/homesec/RL33369.pdf

For an interesting discussion of the inner workings of federal response systems and how better security planning and management can be implemented, see Christopher Cooper and Robert Block, *Disaster: Hurricane Katrina and the Failure of Homeland Security*, Holt Paperbacks (ISBN 978-0805086508), 2007: 352 pp.

23.12 NOTES

1. Federal Emergency Management Agency (FEMA), “National Incident Management System,” FEMA Website, January 30, 2013, www.fema.gov/national-incident-management-system (Federal Emergency Management Agency 2013).
2. FEMA, “National Response Framework,” FEMA Website, August 20, 2012, www.fema.gov/national-response-framework (Federal Emergency Management Agency 2012).
3. FEMA, “IS-860 National Infrastructure Protection Plan (NIPP) Course,” FEMA Emergency Management Institute, May 24, 2011, www.training.fema.gov/EMIWeb/IS/is860.asp (Federal Emergency Management Agency 2011).
4. FEMA, “Hazard Mitigation Planning Resources,” FEMA Website, June 15, 2012, www.fema.gov/hazard-mitigation-planning-resources
5. “Welcome to Sarbanes Oxley 101,” Sarbanes-Oxley-101.com Website, January, 2013, www.sarbanes-oxley-101.com
6. Department of Health & Human Services, “Health Information Privacy,” [hhs.gov](http://hhs.gov/ocr/privacy/hipaa/understanding/index.html) Website, January 2013, www.hhs.gov/ocr/privacy/hipaa/understanding/index.html
7. Federal Trade Commission, “How To Comply with the Privacy of Consumer Financial Information Rule of the Gramm-Leach-Bliley Act.” Bureau of Consumer Protection Business Center, July 2002, <http://business.ftc.gov/documents/bus67-how-comply-privacy-consumer-financial-information-rule-gramm-leach-bliley-act>
8. Emergency Management Accreditation Program, “EMAP Home.” EMAP Website, 2010, www.emaponline.org/index.php
9. Maguire, Valerie, “TIA Publishes New Standard for Telecommunication Administration,” Network Infrastructure Blog, July 20, 2012, <http://blog.siemon.com/infrastructure/tia-publishes-new-standard-for-telecommunication-administration>
10. National Fire Protection Association, “NFPA 70: National Electrical Code—Current Edition 2011,” NFPA Codes & Standards, 2011, www.nfpa.org/aboutthecodes/AboutTheCodes.asp?DocNum=70
11. Ward, Mike, “Start of airport-style security at Texas Capitol: Some visitors complain that licensed gun holders are sped through,” Austin American Statesman, May 21, 2010, www.statesman.com/news/news/state-regional-govt-politics/start-of-airport-style-security-at-texas-capitol/nRs3C
12. Department of Homeland Security, “Airport Passenger Screening: Backscatter X-Ray Machines: Solicitation Number RSEN-13-00004,” Federal Business Opportunities Website, December 13, 2012, www.fbo.gov/index?s=opportunity&mode=form&id=0af3059baf5a4b278a75936152e93253&tab=core&_cview=0
13. FEMA, “Multi-Hazard Mitigation Planning,” FEMA Website, September 12, 2012, www.fema.gov/multi-hazard-mitigation-planning
14. FEMA, “State and Local Mitigation Planning Fact Sheet,” FEMA Website, May 2011, www.fema.gov/pdf/media/factsheets/2011/mit_state_local_plan.pdf

23 · 48 PROTECTING THE PHYSICAL INFORMATION INFRASTRUCTURE

15. FEMA, "Mitigation Planning How-To Guide #1: Getting Started: Building Support for Mitigation Planning (FEMA 386-1)," FEMA Website, September 2002, www.fema.gov/library/viewRecord.do?id=1867
16. FEMA, "Mitigation Planning How-To Guide #2: Understanding Your Risks: Identifying Hazards and Estimating Losses (FEMA 386-2)," FEMA Website, August 2001, www.fema.gov/library/viewRecord.do?id=1880
17. FEMA, "Mitigation Planning How-To Guide #3: Developing the Mitigation Plan: Identifying Mitigation Actions and Implementation Strategies (FEMA 386-3)," FEMA Website, April 2003, www.fema.gov/library/viewRecord.do?id=1886
18. FEMA, "Mitigation Planning How-To Guide #4: Bringing the Plan to Life: Implementing the Hazard Mitigation Plan (FEMA 386-4)," FEMA Website, August 2003, www.fema.gov/library/viewRecord.do?id=1887
19. FEMA, "Mitigation Planning How-To Guide #5: Using Benefit-Cost Review in Mitigation Planning (FEMA 386-5)," FEMA Website, May 2007, www.fema.gov/library/viewRecord.do?id=2680
20. FEMA, "Integrating Historic Property and Cultural Resource Considerations into Hazard Mitigation Planning (FEMA 386-6)," FEMA Website, May 2005, www.fema.gov/library/viewRecord.do?id=1892
21. FEMA, "Integrating Manmade Hazards Into Mitigation Planning (FEMA 386-7)," FEMA Website, September 2003, www.fema.gov/library/viewRecord.do?id=1915
22. FEMA, "Multi-Jurisdictional Mitigation Planning (FEMA 386-8)," FEMA Website, August 2006, www.fema.gov/library/viewRecord.do?id=1905
23. FEMA, "Mitigation Planning How-To #9: Using the Hazard Mitigation Plan to Prepare Successful Mitigation Projects (FEMA 386-9)," FEMA Website, August 2008, www.fema.gov/library/viewRecord.do?id=3388
24. FEMA, "Benefit-Cost Analysis," FEMA Website, February 1, 2013, www.fema.gov/benefit-cost-analysis
25. FEMA, "National Response Network," FEMA Website, August 20, 2012, www.fema.gov/national-response-framework

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER **24**

OPERATING SYSTEM SECURITY

William Stallings

24.1 INFORMATION PROTECTION AND SECURITY	24·1	24.3.5 Protection Based on Virtualization	24·10
24.2 REQUIREMENTS FOR OPERATING SYSTEM SECURITY	24·2	24.4 FILE SHARING	24·11
24.2.1 Requirements	24·2	24.4.1 Access Rights	24·11
24.2.2 Computer System Assets	24·3	24.4.2 Simultaneous Access	24·12
24.2.3 Design Principles	24·4	24.5 TRUSTED SYSTEMS	24·12
		24.5.1 Trojan Horse Defense	24·14
24.3 PROTECTION MECHANISMS	24·4	24.6 WINDOWS 2000 SECURITY	24·14
24.3.1 Protection of Memory	24·5	24.6.1 Access-Control Scheme	24·16
24.3.2 User-Oriented Access Control	24·6	24.6.2 Access Token	24·16
24.3.3 Data-Oriented Access Control	24·7	24.6.3 Security Descriptors	24·17
24.3.4 Protection Based on an Operating System Mode	24·9	24.7 FURTHER READING	24·20
		24.8 NOTES	24·21

24.1 INFORMATION PROTECTION AND SECURITY. This chapter reviews the principles of security in operating systems. Some general-purpose tools can be built into computers and operating systems (OSs) that support a variety of protection and security mechanisms. In general, the concern is with the problem of controlling access to computer systems and the information stored in them.

One of the core concepts in all such discussions is the *process*, which is defined as the execution of a specific piece of code by a particular user at a particular time on a particular processor.

Four types of overall protection policies, of increasing order of difficulty, have been identified:

- 1. No sharing.** In this case, processes are completely isolated from each other, and each process has exclusive control over the resources statically or dynamically assigned to it. With this policy, processes often “share” a program or data file by making a copy of it and transferring the copy into their own virtual memory.

24 · 2 OPERATING SYSTEM SECURITY

2. **Sharing originals of program or data files.** With the use of reentrant code, a single physical realization of a program can appear in multiple virtual address spaces, as can read-only data files. Special locking mechanisms are required for the sharing of writable data files, to prevent simultaneous users from interfering with each other.
3. **Confined, or memoryless, subsystems.** In this case, processes are grouped into subsystems to enforce a particular protection policy. For example, a “client” process calls a “server” process to perform some task on data. The server is to be protected against the client discovering the algorithm by which it performs the task, while the client is to be protected against the server’s retaining any information about the task being performed.
4. **Controlled information dissemination.** In some systems, security classes are defined to enforce a particular dissemination policy. Users and applications are given security clearances of a certain level, while data and other resources (e.g., input/output [I/O] devices) are given security classifications. The security policy enforces restrictions concerning which users have access to which classifications. This model is useful not only in the military context but in commercial applications as well.¹

Much of the work in security and protection as it relates to OSs can be roughly grouped into three categories.

1. **Access control.** Concerned with regulating user access to the total system, subsystems, and data, and regulating process access to various resources and objects within the system.
2. **Information flow control.** Regulates the flow of data within the system and its delivery to users.
3. **Certification.** Relates to proving that access and flow control mechanisms perform according to their specifications and that they enforce desired protection and security policies.

This chapter looks at some of the key mechanisms for providing OS security and then examines Windows 2000 as a case study.

24.2 REQUIREMENTS FOR OPERATING SYSTEM SECURITY

24.2.1 Requirements. Understanding the types of threats to OS security requires a definition of security requirements. OS security addresses four requirements:

1. **Confidentiality.** Requires that the information in a computer system be accessible only for reading by authorized parties. This type of access includes printing, displaying, and other forms of disclosure, including simply revealing the existence of an object.
2. **Integrity.** Requires that only authorized parties be able to modify computer system assets. Modification includes writing, changing, changing status, deleting, and creating.

REQUIREMENTS FOR OPERATING SYSTEM SECURITY 24 · 3

3. **Availability.** Requires that computer system assets are available to authorized parties.
4. **Authentication.** Requires that a computer system be able to verify the identifier of a user, a device or a process.

24.2.2 Computer System Assets. The assets of a computer system can be categorized as hardware, software, and data.

24.2.2.1 Hardware. The main threat to computer system hardware is in the area of availability. Hardware is the most vulnerable to attack and the least amenable to automated controls. Threats include accidental and deliberate damage to equipment as well as theft. The proliferation of personal computers and workstations and the increasing use of local area networks (LANs) increase the potential for losses in this area. Physical and administrative security measures are needed to deal with these threats. Chapters 22 and 23 in this *Handbook* discuss physical security.

24.2.2.2 Software. The OS, utilities, and application programs are what make computer system hardware useful to businesses and individuals. Several distinct threats need to be considered.

A key threat to software is an attack on availability. Software, especially application software, is surprisingly easy to delete. Software also can be altered or damaged to render it useless or dangerous. Careful software configuration management, which includes making backups of the most recent version of software, can maintain high availability. A more difficult problem to deal with is software modification that results in a program that still functions but that behaves differently from before. A final problem is control or possession of software. Although certain countermeasures are available, by and large the problem of unauthorized copying of software has not been solved.

Chapters 38, 39, and 40 of this *Handbook* discuss software security in some detail.

24.2.2.3 Data. Hardware and software security typically are concerns of computing center professionals or individual concerns of personal computer users. A much more widespread problem is data security, which involves files and other forms of data controlled by individuals, groups, and business organizations.

Security concerns with respect to data are broad, encompassing confidentiality, control or possession, integrity, authenticity, availability and utility. For a sound theoretical treatment of the attributes of information that must be protected through security measures, see Chapter 3 in this *Handbook*.

In the case of availability, the concern is with the *destruction* of data files, which can occur either accidentally or maliciously, and with delays in timely access to data.

The obvious concern with confidentiality is the unauthorized reading of data files or databases, and this area has been the subject of perhaps more research and effort than any other area of computer security. A less obvious secrecy threat involves the analysis of data and manifests itself in the use of so-called statistical databases or *data mining*, which provide summary or aggregate information and potentially lead to discovery of unpublicized tendencies, relations, or trends. Aggregate information does not necessarily threaten the privacy of the individuals involved. However, as the use of statistical databases grows, there is an increasing potential for disclosure of personal information through induction or deduction. In essence, characteristics of constituent

24 · 4 OPERATING SYSTEM SECURITY

individuals may be identified through careful analysis. To take a simple example, if one table records the aggregate of the incomes of respondents A, B, C, and D and another records the aggregate of the incomes of A, B, C, D, and E, the difference between the two aggregates would be the income of E. This problem is exacerbated by the increasing desire to combine data sets. In many cases, matching several sets of data for consistency at levels of aggregation appropriate to the problem requires a retreat to elemental units in the process of constructing the necessary aggregates. Thus, the elemental units, which are the subject of privacy concerns, are available at various stages in the processing of data sets.

Finally, data integrity is a major concern in most installations. Modifications to data files can have consequences ranging from minor to disastrous.

24.2.3 Design Principles. Saltzer and Schroeder identify a number of principles for the design of security measures for the various threats to computer systems. These include:

- **Least privilege.** Every program and every user of the system should operate using the least set of privileges necessary to complete the job. Access rights should be acquired by explicit permission only; the default should be “no access.”
- **Economy of mechanisms.** Security mechanisms should be as small and simple as possible, aiding in their verification. This usually means that they must be an integral part of the design rather than add-on mechanisms to existing designs.
- **Acceptability.** Security mechanisms should not interfere unduly with the work of users. At the same time, the mechanisms should meet the needs of those who authorize access. If the mechanisms are not easy to use, they are likely to be unused or incorrectly used.
- **Complete mediation.** Every access must be checked against the access-control information, including those accesses occurring outside normal operation, as in recovery or maintenance.
- **Open design.** The security of the system should not depend on keeping the design of its mechanisms secret. Thus, the mechanisms can be reviewed by many experts, and users can have high confidence in them.²

24.3 PROTECTION MECHANISMS. The introduction of multiprogramming brought about the ability to share resources among users. This sharing involves not just the processor but also:

- Memory
- I/O devices, such as disks and printers
- Programs
- Data

The ability to share these resources introduced the need for protection. Pfleeger and Pfleeger point out that an OS may offer protection along this spectrum:

- **No protection.** This is appropriate when sensitive procedures are being run at separate times.

PROTECTION MECHANISMS 24 · 5

- **Isolation.** This approach implies that each process operates separately from other processes, with no sharing or communication. Each process has its own address space, files, and other objects.
- **Share all or share nothing.** The owner of an object (e.g., a file or memory segment) declares it to be public or private. In the former case, any process may access the object; in the latter, only the owner's processes may access the object.
- **Share via access limitation.** The OS checks the permissibility of each access by a specific user to a specific object. The OS therefore acts as a guard, or gatekeeper, between users and objects, ensuring that only authorized accesses occur.
- **Share via dynamic capabilities.** This extends the concept of access control to allow dynamic creation of sharing rights for objects.
- **Limit use of an object.** This form of protection limits not just access to an object but the use to which that object may be put. For example, a user may be allowed to view a sensitive document but not print it. Another example is that a user may be allowed access to a database to derive statistical summaries but not to determine specific data values.³

The preceding items are listed roughly in increasing order of difficulty to implement but also in increasing order of fineness of protection that they provide. A given OS may provide different degrees of protection for different objects, users, or applications.

The OS needs to balance the need to allow sharing, which enhances the utility of the computer system, with the need to protect the resources of individual users. This section considers some of the mechanisms by which OSs have enforced protection for these objects.

24.3.1 Protection of Memory. In a multiprogramming environment, protection of main memory (random-access memory, or *RAM*) is essential. The concern here is not just security but the correct functioning of the various processes that are active. If one process can inadvertently write into the memory space of another process, then the latter process may not execute properly.

The separation of the memory space of various processes is accomplished easily with a virtual memory scheme. Either segmentation or paging, or the two in combination, provides an effective means of managing main memory. If complete isolation is sought, then the OS simply must ensure that each segment or page is accessible only by the process to which it is assigned. This is accomplished easily by requiring that there be no duplicate entries in page and/or segment tables.

If sharing is to be allowed, then the same segment or page may appear in more than one table. This type of sharing is accomplished most easily in a system that supports segmentation or a combination of segmentation and paging. In this case, the segment structure is visible to the application, and the application can declare individual segments to be sharable or nonsharable. In a pure paging environment, it becomes more difficult to discriminate between the two types of memory, because the memory structure is transparent to the application.

Segmentation, especially, lends itself to the implementation of protection and sharing policies. Because each segment table entry includes a length as well as a base address, a program cannot inadvertently access a main memory location beyond the limits of a segment. To achieve sharing, it is possible for a segment to be referenced in the segment tables of more than one process. The same mechanisms are available in a paging system. However, in this case the page structure of programs and data is not visible

24 · 6 OPERATING SYSTEM SECURITY

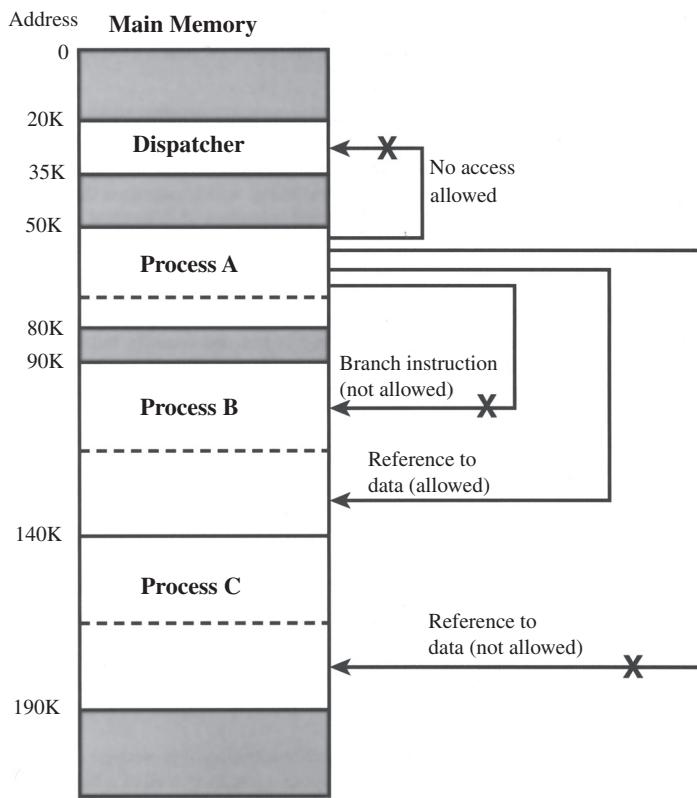


EXHIBIT 24.1 Protection Relationships between Segments

to the programmer, making the specification of protection and sharing requirements more awkward. Exhibit 24.1 illustrates the types of protection relationships that can be enforced in such a system.

An example of the hardware support that can be provided for memory protection is that of the IBM System/370 family of machines, on which OS/390 runs. Associated with each page frame in main memory is a 7-bit storage control key, which may be set by the OS. Two of the bits indicate whether the page occupying this frame has been referenced and changed; these bits are used by the page replacement algorithm. The remaining bits are used by the protection mechanism: a 4-bit access-control key and a fetch-protection bit. Processor references to memory and direct memory access (DMA). DMA I/O memory references must use a matching key to gain permission to access that page. The fetch-protection bit indicates whether the access-control key applies to writes or to both reads and writes. In the processor, there is a program status word (PSW), which contains control information relating to the process that is currently executing. Included in this word is a 4-bit PSW key. When a process attempts to access a page or to initiate a DMA operation on a page, the current PSW key is compared to the access code. A write operation is permitted only if the codes match. If the fetch bit is set, then the PSW key must match the access code for read operations.

24.3.2 User-Oriented Access Control. The measures taken to control access in a data processing system fall into two categories: those associated with the user and those associated with the data.

PROTECTION MECHANISMS 24 · 7

The most common technique for user access control on a shared system or server is the user logon, which requires both a user identifier (ID) and some form of *authentication*, such as providing a password, a token, or biometric attributes. Authentication refers to the binding of a real-world identity (for example, a named employee or a named role in an organization) and the ID being presented. The system will allow a user to log on only if that user's ID is known to the system and if the user knows the password associated by the system with that ID.

Once a user has established a *session*, the operating system can then *authorize* different forms of access (e.g., read, write, append, lock, or execute) to different types of data (e.g., specific files, databases, devices, or communications).

User access control in a distributed environment can be either centralized or decentralized. In a centralized approach, the network provides a logon service, determining who is allowed to use the network and to whom the user is allowed to connect.

Decentralized user access control treats the network as a transparent communication link, and the destination host carries out the usual logon procedure. The security concerns for transmitting passwords over the network must still be addressed.

In many networks, two levels of access control may be used. Individual hosts may be provided with a logon facility to protect host-specific resources and application. In addition, the network as a whole may provide protection to restrict network access to authorized users. This two-level facility is desirable for the common case, currently, in which the network connects disparate hosts and simply provides a convenient means of terminal-host access. In a more uniform network of hosts, some centralized access policy could be enforced in a network control center.

Chapters 28 and 29 of this *Handbook* present more information about identification and authentication.

24.3.3 Data-Oriented Access Control. Following successful logon, the user is granted access to one or a set of hosts and applications. This is generally not sufficient for a system that includes sensitive data in its database. Through the user access-control procedure, a user can be identified to the system. Associated with each user, there can be a profile that specifies permissible operations and file accesses. The OS can then enforce rules based on the user profile. The database-management system, however, must control access to specific records or even portions of records. For example, it may be permissible for anyone in administration to obtain a list of company personnel, but only selected individuals may have access to salary information. The issue is more than just one of level of detail. Whereas the OS may grant a user permission to access a file or use an application, following which there are no further security checks, the database management system must make a decision on each individual access attempt. That decision will depend not only on the user's identity but also on the specific parts of the data being accessed and even on the information already divulged to the user.

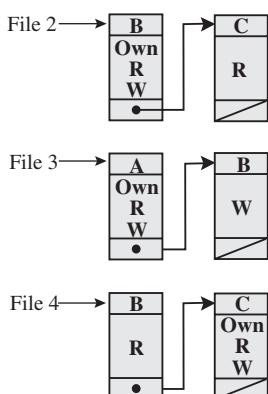
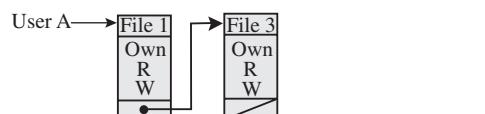
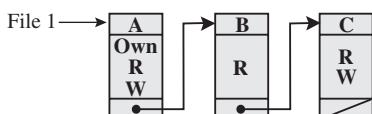
A general model of access control as exercised by a file or database management system is that of an *access matrix* (see Exhibit 24.2a). The basic elements of the model are:

- **Subject.** An entity capable of accessing objects. Generally, the concept of subject equates with that of process. Any user or application actually gains access to an object by means of a process that represents that user or application.
- **Object.** Anything to which access is controlled. Examples include files, portions of files, programs, and segments of memory.
- **Access right.** The way in which an object is accessed by a subject. Examples are read, write, and execute.

24 · 8 OPERATING SYSTEM SECURITY

	File 1	File 2	File 3	File 4	Account 1	Account 2
User A	Own R W		Own R W		Inquiry Credit	
User B	R	Own R W	W	R	Inquiry Debit	Inquiry Credit
User C	R W	R		Own R W		Inquiry Debit

(a) Access matrix



(b) Access control lists for files of part (a)

(c) Capability lists for files of part (a)

EXHIBIT 24.2 Example of Access-Control Structures
Source: Based on a figure in Sandhu (1996).

One dimension of the matrix consists of identified subjects that may attempt data access. Typically, this list will consist of individual users or user groups, although access could be controlled for terminals, hosts, or applications, instead of, or in addition to, users. The other dimension lists the objects that may be accessed. At the greatest level of detail, objects may be individual data fields. More aggregate groupings, such as records, files, or even the entire database, also may be objects in the matrix. Each entry in the matrix indicates the access rights of that subject for that object.

In practice, an access matrix usually is sparse and is implemented by decomposition in one of two ways. The matrix may be decomposed by columns, yielding *access-control lists* (see Exhibit 24.2b). Thus, for each object, an access-control list lists users and their permitted access rights. The access-control list may contain a default, or public, entry. This allows users who are not explicitly listed as having special rights to have a default set of rights. Elements of the list may include individual users as well as groups of users.

Decomposition by rows yields *capability tickets* (see Exhibit 24.2c). A capability ticket specifies authorized objects and operations for a user. Each user has a number

PROTECTION MECHANISMS 24 · 9

of tickets and may be authorized to lend or give them to others. Because tickets may be dispersed around the system, they present a greater security problem than access-control lists. In particular, the ticket must be unforgeable. One way to accomplish this is to have the OS hold all tickets on behalf of users. These tickets would have to be held in a region of memory inaccessible to users.

Network considerations for data-oriented access control parallel those for user-oriented access control. If only certain users are permitted to access certain items of data, then encryption may be needed to protect those items during transmission to authorized users. Typically, data access control is decentralized, that is, controlled by host-based database management systems. If a network database server exists on a network, then data access control becomes a network function.

24.3.4 Protection Based on an Operating System Mode. One technique used in all OSs to provide protection is based on the mode of processor execution. Most processors support at least two modes of execution: the mode normally associated with the OS and that normally associated with user programs. Certain instructions can be executed only in the more privileged mode. These would include reading or altering a control register, such as the program status word; primitive I/O instructions; and instructions that relate to memory management. In addition, certain regions of memory can be accessed only in the more privileged mode.

The less privileged mode often is referred to as the *user mode*, because user programs typically would execute in this mode. The more privileged mode is referred to as the *system mode*, *control mode*, or *kernel mode*. This last term refers to the kernel of the OS, which is that portion of the OS that encompasses the important system functions. Exhibit 24.3 lists the functions typically found in the kernel of an OS.

EXHIBIT 24.3 Typical Kernel Mode Operating System Functions

Process Management

- Process creation and termination
- Process scheduling and dispatching
- Process switching
- Process synchronization and support for interprocess communication
- Management of process control blocks

Memory Management

- Allocation of address space to processes
- Swapping
- Page and segment management

I/O Management

- Buffer management
- Allocation of I/O channels and devices to processes

Support functions

- Interrupt handling
 - Accounting
 - Monitoring
-

24 · 10 OPERATING SYSTEM SECURITY

The reason for using two modes should be clear. It is necessary to protect the OS and key OS tables, such as process control blocks, from interference by user programs. In the kernel mode, the software has complete control of the processor and all its instructions, registers, and memory. This level of control is not necessary, and for safety is not desirable, for user programs.

Two questions arise: How does the processor know in which mode it is to be executing, and how is the mode changed? Regarding the first question, typically there is a bit in the program status word that indicates the mode of execution. This bit is changed in response to certain events. For example, when a user makes a call to an OS service, the mode is set to the kernel mode. Typically, this is done by executing an instruction that changes the mode. When the user makes a system service call, or when an interrupt transfers control to a system routine, the routine executes the change-mode instruction to enter a more privileged mode and executes it again to enter a less privileged mode before returning control to the user process. If a user program attempts to execute a change-mode instruction, it will simply result in a call to the OS, which will return an error unless the mode change is to be allowed.

More sophisticated mechanisms also can be provided. A common scheme is to use a ring-protection structure. In this scheme, lower-numbered, or inner, rings enjoy greater privilege than higher-numbered, or outer, rings. Typically, ring 0 is reserved for kernel functions of the OS, with applications at a higher level. Some utilities or OS services may occupy an intermediate ring. Basic principles of the ring system are:

- A program may access only those data that reside on the same ring or a less privileged ring.
- A program may call services residing on the same or a more privileged ring.

An example of the ring protection approach is found on the VAX VMS OS, which uses four modes:

1. **Kernel.** Executes the kernel of the VMS OS, which includes memory management, interrupt handling, and I/O operations.
2. **Executive.** Executes many of the OS service calls, including file and record (disk and tape) management routines.
3. **Supervisor.** Executes other OS services, such as responses to user commands.
4. **User.** Executes user programs, plus utilities such as compilers, editors, linkers, and debuggers.

A process executing in a less privileged mode often needs to call a procedure that executes in a more privileged mode; for example, a user program requires an OS service. This call is achieved by using a change-mode (CHM) instruction, which causes an interrupt that transfers control to a routine at the new access mode. A return is made by executing the REI (return from exception or interrupt) instruction.

24.3.5 Protection Based on Virtualization. With the growing availability of memory (e.g., tens to hundreds of gigabytes of RAM), disk space (terabytes to petabytes of storage), faster processors (tens of gigahertz), and multicore systems (potentially thousands of processors working in parallel), *virtualization* of computers has progressed to practical and widespread usability. Instantiations of an operating environment can coexist using shared resources without permitting any direct communication among them. Each instantiation is encapsulated and completely protected against

FILE SHARING 24 · 11

intrusion or interference from processes running on other virtual machines sharing the same physical resources.

24.4 FILE SHARING. Multiuser systems almost always require that files can be shared among a number of users. Two issues arise: access rights and the management of simultaneous access.

24.4.1 Access Rights. The file system should provide a flexible tool for allowing extensive file sharing among users. The file system should provide a number of options so that the way in which a particular file is accessed can be controlled. Typically, users or groups of users are granted certain access rights to a file. A wide range of access rights has been used. The next list indicates access rights that can be assigned to a particular user for a particular file.

- **None.** The user may not even learn of the existence of the file, much less access it. To enforce this restriction, the user would not be allowed to read the user directory that includes this file.
- **Knowledge.** The user can determine that the file exists and who its owner is. The user is then able to petition the owner for additional access rights.
- **Execution.** The user can load and execute a program but cannot copy it. Proprietary programs often are made accessible with this restriction.
- **Locking.** The user can change the status of a logical *flag* that indicates temporary restrictions on access to data. Database management systems provide for locking to control concurrent access to records so that different processes can avoid overwriting each other's modifications.
- **Reading.** The user can read the file for any purpose, including copying and execution. Some systems are able to enforce a distinction between viewing and copying. In the former case, the contents of the file can be displayed to the user, but the user has no means for making a copy.
- **Appending.** The user can add data to the file, often only at the end, but cannot modify or delete any of the file's contents. This right is useful in collecting data from a number of sources.
- **Updating.** The user can modify, delete, and add to the file's data. This normally includes writing the file initially, rewriting it completely or in part, and removing all or a portion of the data. Some systems distinguish among different degrees of updating.
- **Changing protection.** The user can change the access rights granted to other users. Typically, only the owner of the file holds this right. In some systems, the owner can extend this right to others. To prevent abuse of this mechanism, the file owner typically is able to specify which rights can be changed by the holder of this extended right.
- **Deletion.** The user can delete the file from the file system.

These rights can be considered to constitute a hierarchy, with each right implying those that precede it. Thus, if a particular user is granted the updating right for a particular file, then that user also is granted these rights: knowledge, execution, reading, and appending.

24 · 12 OPERATING SYSTEM SECURITY

One user is designated as owner of a given file, usually this is the person who initially created the file. The owner has all of the access rights listed previously and may grant rights to others. Access can be provided to different classes of users:

- **Specific user.** Individual users who are designated by user ID.
- **User groups.** A set of users who are not individually defined. The system must have some way of keeping track of the membership of user groups.
- **All.** All users who have access to this system. These are public files.

24.4.2 Simultaneous Access. When access is granted to append or update a file to more than one user, the OS or file management system must enforce discipline. A brute-force approach is to allow a user to lock the entire file when it is to be updated. A finer grain of control is to lock individual records during update. Issues of mutual exclusion and deadlock must be addressed in designing the shared access capability.

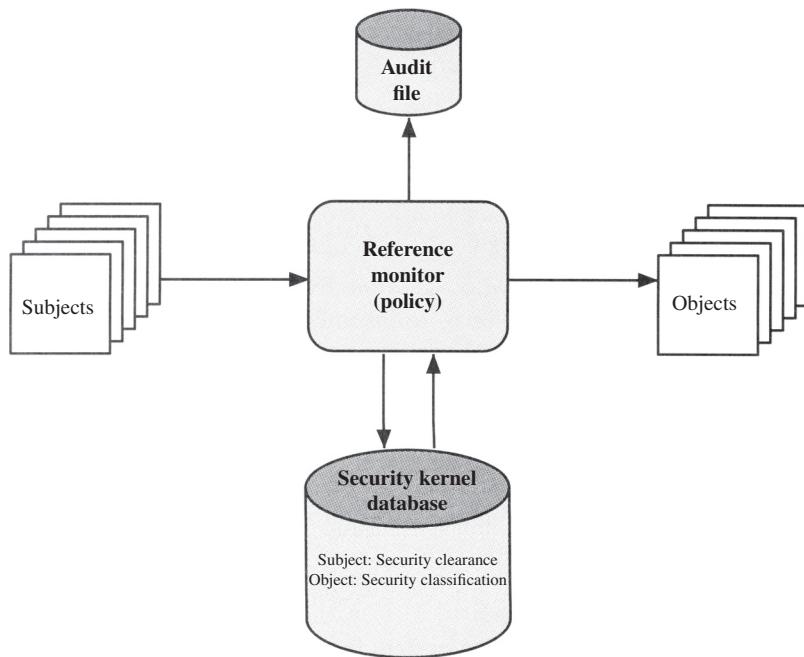
Chapter 52 in this *Handbook* discusses application controls in more detail.

24.5 TRUSTED SYSTEMS. Much of what has been discussed so far has concerned protecting a given message or item from passive or active attack by a given user. A somewhat different but widely applicable requirement is to protect data or resources on the basis of levels of security. This is commonly found in the military, where information is categorized as unclassified (U), confidential (C), secret (S), top secret (TS), or beyond. This concept is equally applicable in other areas, where information can be organized into gross categories and users can be granted clearances to access certain categories of data. For example, the highest level of security might be for strategic corporate planning documents and data, accessible only by corporate officers and their staff; next might come sensitive financial and personnel data, accessible only by administration personnel, corporate officers, and so on.

When multiple categories or levels of data are defined, the requirement is referred to as *multilevel security*. The general statement of the requirement for multilevel security is that a subject at a high level may not convey information to a subject at a lower or incomparable level unless that flow accurately reflects the will of an authorized user. For implementation purposes, this requirement is in two parts and is simply stated. A multilevel secure system must enforce:

1. **No read up.** A subject can only read an object of less or equal security level. This is referred to in the literature as the *simple security property*.
2. **No write down.** A subject can only write into an object of greater or equal security level. This is referred to in the literature as the **-property* (pronounced *star property*).

These two rules, if properly enforced, provide multilevel security. For a data processing system, the approach that has been taken, and has been the object of much research and development, is based on the *reference monitor* concept. This approach is depicted in Exhibit 24.4. The reference monitor is a controlling element in the hardware and OS of a computer that regulates the access of subjects to objects on the basis of security parameters of the subject and object. The reference monitor has access to a file, known as the *security kernel database*, that lists the access privileges (security clearance) of each subject and the protection attributes (classification level) of each

TRUSTED SYSTEMS 24 · 13**EXHIBIT 24.4** Reference Monitor Concept

object. The reference monitor enforces the security rules (no read up, no write down) and has these properties:

- **Complete mediation.** The security rules are enforced on every access, not just, for example, when a file is opened.
- **Isolation.** The reference monitor and database are protected from unauthorized modification.
- **Verifiability.** The reference monitor's correctness must be provable. That is, it must be possible to demonstrate mathematically that the reference monitor enforces the security rules and provides complete mediation and isolation.

These are stiff requirements. The requirement for complete mediation means that every access to data within main memory and on disk and tape must be mediated. Pure software implementations impose too high a performance penalty to be practical; the solution must be at least partly in hardware. The requirement for isolation means that it must not be possible for an attacker, no matter how clever, to change the logic of the reference monitor or the contents of the security kernel database. Finally, the requirement for mathematical proof is formidable for something as complex as a general-purpose computer. A system that can provide such verification is referred to as a *trusted system*.

Chapter 9 in this *Handbook* discusses mathematical models of computer security in more detail.

A final element illustrated in Exhibit 24.4 is an audit file. Important security events, such as detected security violations and authorized changes to the security kernel database, are stored in the audit file.

24 · 14 OPERATING SYSTEM SECURITY

In an effort to meet its own needs and as a service to the public, the U.S. Department of Defense in 1981 established the Computer Security Center within the National Security Agency (NSA), with the goal of encouraging the widespread availability of trusted computer systems. This goal is realized through the center's Commercial Product Evaluation Program. In essence, the center attempts to evaluate commercially available products as meeting the security requirements just outlined. The center classifies evaluated products according to the range of security features that they provide. These evaluations are needed for Department of Defense procurements but are published and freely available. Hence, they can serve as guidance to commercial customers for the purchase of commercially available, off-the-shelf equipment.

24.5.1 Trojan Horse Defense. A *Trojan horse* attack involves software that appears to have acceptable functions but which conceals additional, unauthorized functionality. Chapters 2 and 15 in this *Handbook* provide details of many Trojan horse attacks.

One way to secure against Trojan horse attacks is by the use of a secure, trusted OS. Exhibit 24.5 illustrates an example. In this case, a Trojan horse is used to get around the standard security mechanism used by most file management and OSs: the access-control list. In this example, a user named Bob interacts through a program with a data file containing the critically sensitive character string "CPE170KS". User Bob has created the file with read/write permission provided only to programs executing on his own behalf: that is, only processes that are owned by Bob may access the file.

The Trojan horse attack begins when a hostile user, named Alice, gains legitimate access to the system and installs both a Trojan horse program and a private file to be used in the attack as a "back pocket." Alice gives read/write permission to herself for this file and gives Bob write-only permission (see Exhibit 24.5a). Alice now induces Bob to invoke the Trojan horse program, perhaps by advertising it as a useful utility. When the program detects that it is being executed by Bob, it reads the sensitive character string from Bob's file and copies it into Alice's back-pocket file (see Exhibit 24.5b). Both the read and write operations satisfy the constraints imposed by access-control lists. Alice then has only to access Bob's file at a later time to learn the value of the string.

Now consider the use of a secure OS in this scenario (see Exhibit 24.5c). Security levels are assigned to subjects at logon on the basis of criteria such as the terminal from which the computer is being accessed and the user involved, as identified by password/ID. In this example, there are two security levels, sensitive (gray) and public (white), ordered so that sensitive is higher than public. Processes owned by Bob and Bob's data file are assigned the security level sensitive. Alice's file and processes are restricted to public. If Bob invokes the Trojan horse program (see Exhibit 24.5d), that program acquires Bob's security level. It is therefore able, under the simple security property, to observe the sensitive character string. When the program attempts to store the string in a public file (the back-pocket file), however, the *-property is violated and the attempt is disallowed by the reference monitor. Thus, the attempt to write into the back-pocket file is denied even though the access-control list permits it: The security policy takes precedence over the access-control list mechanism.

24.6 WINDOWS 2000 SECURITY. A good example of the access-control concepts discussed in this chapter is the Windows 2000 (W2K) access-control facility, which exploits object-oriented concepts to provide a powerful and flexible access-control capability.

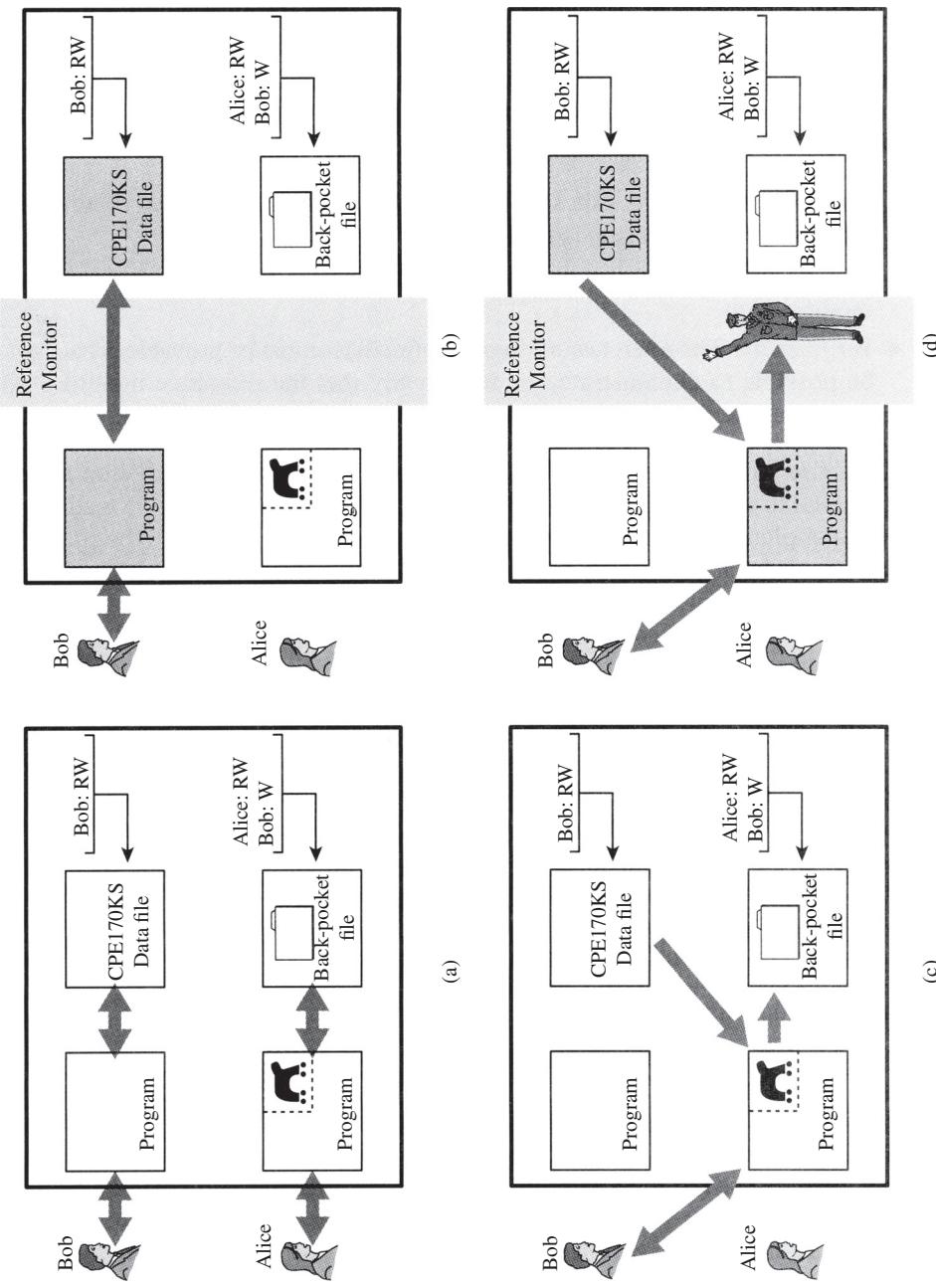


EXHIBIT 24.5 Trojan Horse and Secure Operating Systems

24 · 16 OPERATING SYSTEM SECURITY

W2K provides a uniform access-control facility that applies to processes, threads, files, semaphores, windows, and other objects. Access control is governed by two entities: an access token associated with each process and a security descriptor associated with each object for which interprocess access is possible.

24.6.1 Access-Control Scheme. When a user logs on to a W2K system, W2K uses a name/password scheme to authenticate the user. If the logon is accepted, a process is created for the user and an access token is associated with that process object. The access token, whose details are described later, include a security ID (SID), which is the identifier by which this user is known to the system for purposes of security. When the initial user process spawns any additional processes, the new process object inherits the same access token.

The access token serves two purposes:

1. It keeps all necessary security information together to speed access validation. When any process associated with a user attempts access, the security subsystem can make use of the token associated with that process to determine the user's access privileges.
2. It allows each process to modify its security characteristics in limited ways without affecting other processes running on behalf of the user.

The chief significance of the second point has to do with privileges that may be associated with a user. The access token indicates which privileges a user may have. Generally, the token is initialized with each of these privileges in a disabled state. Subsequently, if one of the user's processes needs to perform a privileged operation, the process may enable the appropriate privilege and attempt access. It would be undesirable to keep all of the security information for a user in one systemwide place, because in that case enabling a privilege for one process enables it for all of them.

A security descriptor is associated with each object for which interprocess access is possible. The chief component of the security descriptor is an access-control list that specifies access rights for various users and user groups for this object. When a process attempts to access this object, the SID of the process is matched against the access-control list of the object to determine if access will be allowed.

When an application opens a reference to a securable object, W2K verifies that the object's security descriptor grants the application's user access. If the check succeeds, W2K caches the resulting granted access rights.

An important aspect of W2K security is the concept of impersonation, which simplifies the use of security in a client/server environment. If client and server talk through a remote procedure call (RPC) connection, the server can temporarily assume the identity of the client so that it can evaluate a request for access relative to that client's rights. After the access, the server reverts to its own identity.

24.6.2 Access Token. Exhibit 24.6a shows the general structure of an access token, which includes these parameters:

- **Security ID (SID).** Identifies a user uniquely across all of the machines on the network. This generally corresponds to a user's logon name.
- **Group SIDs.** A list of the groups to which this user belongs. A group is simply a set of user IDs that are identified as a group for purposes of access control. Each

WINDOWS 2000 SECURITY 24 · 17

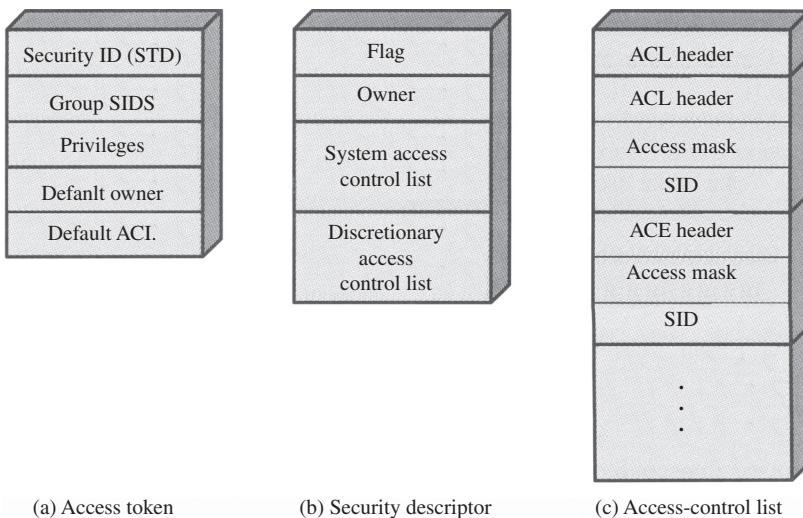


EXHIBIT 24.6 Windows 2000 Security Structures

group has a unique group SID. Access to an object can be defined on the basis of group SIDs, individual SIDs, or a combination.

- **Privileges.** A list of security-sensitive system services that this user may call. An example is create token. Another example is the set backup privilege; users with this privilege are allowed to use a backup tool to back up files that they normally would not be able to read. Most users will have no privileges.
 - **Default owner.** If this process creates another object, this field specifies who is the owner of the new object. Generally, the owner of the new process is the same as the owner of the spawning process. However, a user may specify that the default owner of any processes spawned by this process is a group SID to which this user belongs.
 - **Default access control list (ACL).** This is an initial list of protections applied to the objects that the user creates. The user may subsequently alter the ACL for any object that it owns or that one of its groups owns.

24.6.3 Security Descriptors. Exhibit 24.6b shows the general structure of a security descriptor, which includes these parameters:

- **Flags.** Defines the type and contents of a security descriptor. The flags indicate whether the system access-control list (SACL) and discretionary access control list (DACL) are present, whether they were placed on the object by a defaulting mechanism, and whether the pointers in the descriptor use absolute or relative addressing. Relative descriptors are required for objects that are transmitted over a network, such as information transmitted in an RPC.
 - **Owner.** The owner of the object generally can perform any action on the security descriptor. The owner can be an individual or a group SID. The owner has the authority to change the contents of the DACL.
 - **System access control list (SACL).** Specifies what kinds of operations on the object should generate audit messages. An application must have the corresponding

24 · 18 OPERATING SYSTEM SECURITY

privilege in its access token to read or write the SACL of any object. This is to prevent unauthorized applications from reading SACLs (thereby learning what not to do to avoid generating audits) or writing them (to generate many audits to cause an illicit operation to go unnoticed).

- **Discretionary access-control list (DACL).** Determines which users and groups can access this object for which operations. It consists of a list of access-control entries (ACEs).

When an object is created, the creating process can assign as owner its own SID or any group SID in its access token. The creating process cannot assign an owner that is not in the current access token. Subsequently, any process that has been granted the right to change the owner of an object may do so, but again with the same restriction. The reason for the restriction is to prevent a user from covering his or her tracks after attempting some unauthorized action.

Let us look in more detail at the structure of access-control lists, because these are at the heart of the W2K access-control facility (see Exhibit 24.7). Each list consists of an overall header and a variable number of access-control entries. Each entry specifies an individual or group SID and an access mask that defines the rights to be granted to this SID. When a process attempts to access an object, the object manager in the W2K executive reads the SID and group SIDs from the access token and then scans down the object's DACL. If a match is found—that is, if an ACE is found with a SID that matches one of the SIDs from the access token—then the process has the access rights specified by the access mask in that ACE.

Exhibit 24.7 shows the contents of the access mask. The least significant 16 bits specify access rights that apply to a particular type of object. For example, bit 0 for a file object is File_Read_Data access, and bit 0 for an event object is Event_Query_Status access.

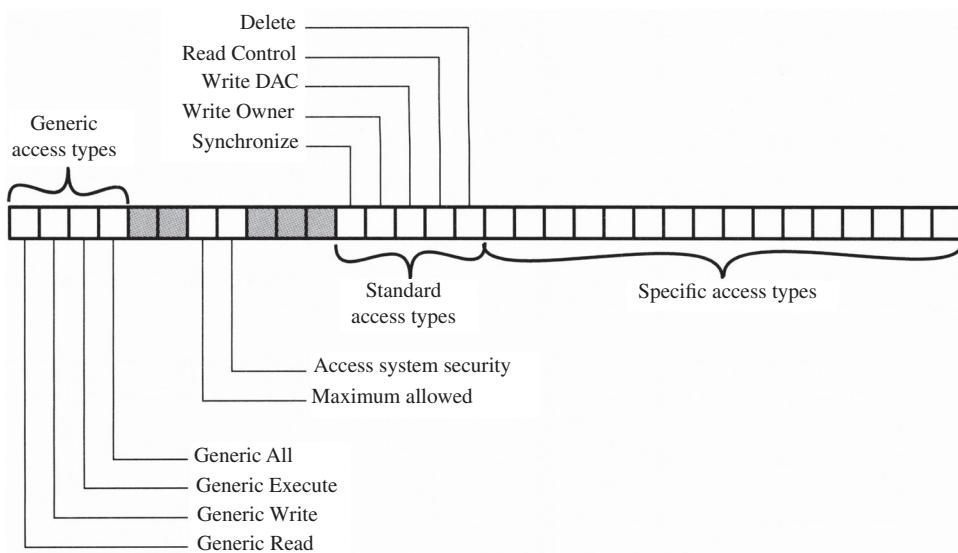


EXHIBIT 24.7 Windows 2000 Access Mask

WINDOWS 2000 SECURITY 24 · 19

The most significant 16 bits of the mask contains bits that apply to all types of objects. Five of these are referred to as standard access types:

1. **Synchronize.** Gives permission to synchronize execution with some event associated with this object. In particular, this object can be used in a wait function.
2. **Write_owner.** Allows a program to modify the owner of the object. This is useful because the owner of an object always can change the protection on the object. (The owner may not be denied Write DAC access.)
3. **Write_DAC.** Allows the application to modify the DACL and hence the protection on this object.
4. **Read_control.** Allows the application to query the owner and DACL fields of the security descriptor of this object.
5. **Delete.** Allows the application to delete this object.

The high-order half of the access mask also contains the four generic access types. These bits provide a convenient way to set specific access types in a number of different object types. For example, suppose an application wishes to create several types of objects and ensure that users have read access to the objects, even though read has a somewhat different meaning for each object type. To protect each object of each type without the generic access bits, the application would have to construct a different ACE for each type of object and be careful to pass the correct ACE when creating each object. It is more convenient to create a single ACE that expresses the generic concept *allow read*; simply apply this ACE to each object that is created and have the right thing happen. That is the purpose of the generic access bits, which are:

- **Generic_all.** Allow all access.
- **Generic_execute.** Allow execution if executable.
- **Generic_write.** Allow write access.
- **Generic_read.** Allow read-only access.

The generic bits also affect the standard access types. For example, for a file object, the Generic_Read bit maps to the standard bits Read_Control and Synchronize and to the object-specific bits File_Read_Data, File_Read_Attributes, and File_Read_EA. Placing an ACE on a file object that grants some SID Generic_Read grants those five access rights as if they had been specified individually in the access mask.

The remaining two bits in the access mask have special meanings. The Access_System_Security bit allows modifying audit and alarm control for this object. However, not only must this bit be set in the ACE for a SID, but the access token for the process with that SID must have the corresponding privilege enabled.

Finally, the Maximum_Allowed bit is not really an access bit but a bit that modifies W2K's algorithm for scanning the DACL for this SID. Normally, W2K will scan through the DACL until it reaches an ACE that specifically grants (bit set) or denies (bit not set) the access requested by the requesting process or until it reaches the end of the DACL, in which latter case access is denied. The Maximum_Allowed bit allows the object's owner to define a set of access rights that is the maximum that will be allowed to a given user. With this in mind, suppose that an application does not know all of the

24 · 20 OPERATING SYSTEM SECURITY

operations that it is going to be asked to perform on an object during a session. There are three options for requesting access:

1. Attempt to open the object for all possible accesses. The disadvantage of this approach is that the access may be denied even though the application may have all of the access rights actually required for this session.
2. Only open the object when a specific access is requested, and open a new handle to the object for each different type of request. This is generally the preferred method because it will not unnecessarily deny access, nor will it allow more access than necessary. However, it imposes additional overhead.
3. Attempt to open the object for as much access as the object will allow this SID. The advantage is that the user will not be artificially denied access, but the application may have more access than it needs. This latter situation may mask bugs in the application.

An important feature of W2K security is that applications can make use of the W2K security framework for user-defined objects. For example, a database server might create its own security descriptors and attach them to portions of a database. In addition to normal read/write access constraints, the server could secure database-specific operations, such as scrolling within a result set or performing a join. It would be the server's responsibility to define the meaning of special rights and perform access checks. But the checks would occur in a standard context, using systemwide user/group accounts and audit logs. The extensible security model should prove useful to implementers of foreign file systems.

24.7 FURTHER READING

- Boebert, W., R. Kain, and W. Young. "Secure Computing: The Secure Ada Target Approach." *Scientific Honeyweller* (July 1985). Reprinted in M. Abrams and H. Podell, *Computer and Network Security*. IEEE Computer Society Press, 1987.
- Bransted, D. (ed.). *Computer Security and the Data Encryption Standard*. National Bureau of Standards, Special Publication No. 500-27, February 1978.
- Denning, P., and R. Brown. "Operating Systems." *Scientific American* 251 (September 1984): 94–106.
- Gasser, M. *Building a Secure Computer System*. New York: Van Nostrand Reinhold, 1988.
- Gollmann, D. *Computer Security*, 3rd edition. Hoboken, NJ: Wiley & Sons, 2011.
- Patterson, D. A., and J. L. Hennessy. *Computer Organization and Design: The Hardware/Software Interface*, 4th ed. Morgan Kaufmann, 2011.
- Pfleeger, C. P., and S. L. Pfleeger. *Security in Computing*, 4th ed. Upper Saddle River, NJ: Prentice-Hall, 2006.
- Saltzer, J., and M. Schroeder. "The Protection of Information in Computer Systems." *Proceedings of the IEEE* (September 1975).
- Sandhu, R., and P. Samarati, "Access Control: Principles and Practice." *IEEE Communications* (September 1994).
- Singhal, M., and N. Shivaratri. *Advanced Concepts in Operating Systems*. New York: McGraw-Hill, 1994.
- Sinha, P. K. *Distributed Operating Systems: Concepts and Design*. Hoboken, NJ: Wiley-IEEE Press, 1996.

NOTES 24 · 21

- Stallings, W. *Cryptography and Network Security: Principles and Practice*, 5th ed. Upper Saddle River, NJ: Prentice-Hall, 2010.
- Stallings, W. *Operating Systems: Internals and Design Principles*, 7th ed. Upper Saddle River, NJ: Prentice-Hall, 2011.
- Viega, J., and J. Voas. “The Pros and Cons of Unix and Windows Security Policies.” *IT Professional* 2, no. 5 (September/October 2000): 40–47.

24.8 NOTES

1. P. Denning and R. Brown, “Operating Systems,” *Scientific American* (September 1984).
2. J. Saltzer and M. Schroeder, “The Protection of Information in Computer Systems,” *Proceedings of the IEEE* (September 1975).
3. C. P. Pfleeger and S. L. Pfleeger, *Security in Computing*, 4th ed. (Prentice-Hall PTR, 2006).

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 25

LOCAL AREA NETWORKS

**N. Todd Pritsky, Joseph R. Bumblis, and
Gary C. Kessler**

25.1 INTRODUCTION	25·1	25.5 NETWORK OPERATING SYSTEM ISSUES	25·9
25.2 POLICY AND PROCEDURE ISSUES	25·1	25.5.1 Windows 9x	25·10
		25.5.2 Windows NT/2000, XP, Vista, Win7, and Win8	25·11
25.3 PHYSICAL SITE SECURITY	25·3	25.5.3 UNIX	25·13
25.4 PHYSICAL LAYER ISSUES	25·4	25.5.4 MacOS	25·15
25.4.1 Sniffers and Broadcast LANs	25·4	25.6 CONCLUSION	25·16
25.4.2 Attacks on the Physical Plant	25·5	25.7 FURTHER READING	25·16
25.4.3 Modems, Dial-Up Servers, and Telephone Lines	25·5		
25.4.4 Wireless LAN Issues	25·7	25.8 NOTES	25·17

25.1 INTRODUCTION. This chapter discusses generic issues surrounding local area network (LAN) security. Securing the LAN is essential to securing the Internet because LANs are where most of the attackers, victims, clients, servers, firewalls, routers, and other devices reside. Compromised LAN systems on the Internet open other nodes on that local network to attack and put other systems at risk on the Internet as a whole. Many of the general issues mentioned herein are described in more specific terms in other chapters of this *Handbook*, such as Chapters 15, 22, 23, and 47 in particular.

25.2 POLICY AND PROCEDURE ISSUES. Thirty years ago, nearly all computer users had accounts on a shared mainframe or minicomputer. A single system manager was responsible for security, backup, disaster recovery, account management, policies, and all other related issues. Today all users are system managers, and, in many cases, individuals have responsibility for several systems. Since the vulnerability of a single computer can compromise the entire LAN, it is imperative that there be rules in place so that everyone can work together for mutual efficiency and defense. But where policies and procedures can be centralized, they should be, because most users do not take the security procedures seriously enough.

25 · 2 LOCAL AREA NETWORKS

The next list, modified from the Internet Engineering Task Force (IETF) Request for Comment (RFC) 2196, is a rough outline of LAN-related security policies and procedures that should at least be considered.¹

1. Administrative Policies Framework
 - a. Information security issues
 - i. Password management procedures
 1. Are passwords assigned or chosen by user?
 2. Are password auditing procedures in place?
 3. Are password policies enforced (e.g., minimum length, allowed and required characters, expiration, blacklisting)?
 4. How many passwords are required to access all systems?
 - ii. Virus protection
 1. Are servers protected?
 2. Is there an e-mail “viruswall”?
 3. Is virus protection a centrally managed activity, or up to each user?
 4. How do users maintain the current virus signature database?
 - iii. Encryption and certificates
 - iv. Security event handling
 - b. Network connectivity issues
 - i. Dial-up access
 - ii. Hubs versus switches
 - iii. Identification and authentication
 1. Passwords and management
 2. Authentication systems
 3. Two-factor authentication?
 4. Single sign-on?
 5. Biometrics
 - c. Physical site security and disaster recovery
 - i. Physical security of systems
 - ii. Dial-up access lines and modems
 - iii. Scheduling backups
 - iv. Access to servers
 - v. Storing and limiting access to backup media
 - vi. Disaster recovery and contingency plans
 - d. Operating system and LAN security
 - i. Operating system-specific issues
 1. Monitoring operating system security vulnerabilities
 2. Applying security patches
 3. Securing the operating system (OS) during installation
 4. Auditing the systems

PHYSICAL SITE SECURITY 25 · 3

- ii. Is Dynamic Host Configuration Protocol (DHCP) employed?
 - iii. Log analysis software and procedures
 - iv. Vulnerability testing
 - v. Intrusion detection tools
2. User Policies Framework
- a. Written, published network security policy
 - b. Network/Internet appropriate use policy (AUP)
 - c. User training and education
 - i. Security (general issues)
 - ii. Importance of protecting customer data, client information, and other private information such as medical claim records and patient information
 - iii. Policies and AUPs
 - iv. Suggestions for how to be safe
- d. Virus protection
 - i. Using antivirus software
 - ii. Maintaining a current virus signature database
 - e. Choosing good passwords
 - f. Best practices
 - i. Email (netiquette and handling attachments)
 - ii. Browser (ensure that Java, JavaScript, ActiveX, and other auto-execution code is current as specified by the manufacturer. See “Notes”² for helpful links.)
 - iii. Microsoft Office Suite
 - 1. Use of macros
 - 2. Implications for document management
 - iv. Protecting files on the server and your own system
 - 1. Use of Windows “shares” and UNIX Network File System (NFS)
 - 2. Use of NetWare, New Technology File System (NTFS), and UNIX access controls
 - i. Spotting a possible compromise (methods for identifying intrusions or other unauthorized activity)
 - ii. What to do and whom to contact if a compromise is suspected

Not every issue will apply to all networks, but each LAN is an ever-evolving entity, and the policies guiding its operation also must evolve.

25.3 PHYSICAL SITE SECURITY. The physical protection of the site is very important but too often overlooked. One reason that site security is often lacking is that some of the policies and procedures can be perceived as messages to employees that they are not trusted. Nevertheless, physical site security includes many aspects of protecting network servers, communications facilities, individual user’s systems, and information.

For more information on facilities security, see Chapters 22 and 23 of this *Handbook*.

25 · 4 LOCAL AREA NETWORKS

25.4 PHYSICAL LAYER ISSUES. The LAN itself has a number of additional vulnerabilities because the systems and media are so widely dispersed. This section discusses securing the LAN infrastructure.

25.4.1 Sniffers and Broadcast LANs. Traditional LAN media access control (MAC) schemes operate assuming a logical, if not physical, broadcast topology. In a broadcast network, every station hears every transmission. When operating in *promiscuous mode*, a LAN station will read every frame that goes by regardless of whether the frame is addressed to the station or not.

When protocol analysis, or *sniffer*, software became available in the 1980s, every station on the LAN became a potential network analysis and management tool—or a surreptitious eavesdropper. Since so many of the network's applications—particularly Transmission Control Protocol/Internet Protocol (TCP/IP)-based applications—transmit passwords and files in plaintext, this type of software is potentially very dangerous.

A number of sources offer very powerful and very flexible commercial packet-sniffing software, such as Network Associates' SnifferPro and Novell's LANalyzer. These packages usually have additional capabilities, such as network monitoring, performance monitoring, and traffic analysis. Prior to 1990, network protocol analysis required a special piece of hardware with a connection to the network. Today, a large number of software packages that do TCP/IP packet sniffing can help an intruder because they can be installed directly on an individual's laptop or desktop computer. Some of these packages include:

- WireShark (Windows, MacOS, and UNIX)
- BUTTsniffer (Windows NT)
- Network Monitor (free with, and for, Windows NT)
- Sniffit (Linux, SunOS, Solaris, FreeBSD, Irix)
- snort (UNIX and Windows)
- Solsniff (Solaris)
- tcpdump (MacOS and UNIX)
- WinDump (DOS)

Relatively few countermeasures can be taken against these kinds of tools. Fortunately, they are effective only on a broadcast network, such as a hubbed LAN. If an Ethernet hub, for example, is replaced with a switch, the only traffic broadcast to all hosts on the network are those frames actually addressed to the LAN's broadcast address. In this case, a station can sniff only that traffic that is going to or coming from the station with the sniffer software. Replacing all hubs with switches may be unreasonable in many environments, but placing all servers on a switch instead of a hub will improve both performance and security.

Other network-based tools that can detect a host with a network interface card (NIC) in promiscuous mode, such as AntiSniff (Windows) and sentinel (UNIX). These tools work by performing a number of tests to detect the promiscuous host; they then check the network host's operating systems, domain name system (DNS) activity, and network and machine latency.

Sniffers can be defeated using cryptography. Use of secure IP (IPsec) and Secure Shell (SSH) for TCP/IP applications can provide privacy and integrity to all

PHYSICAL LAYER ISSUES 25 · 5

communication and applications. Information on IPsec is available on www.ipsec-howto.org/x202.html, and SSH is available from SSH Communications Security (www.ssh.com).

25.4.2 Attacks on the Physical Plant. The most common medium employed today on LANs is copper-based, unshielded twisted pair (UTP) cable. All copper media, including UTP and coaxial cable, emanate a magnetic field because of the changing current on the wire. Van Eck monitoring devices can pick up these emanations remotely and reproduce the frames on the wire or keystrokes at a machine. As far-fetched as this might sound, the vulnerability is very real. The U.S. Government and military have a set of standards to reduce and limit electromagnetic radiation (EMR) called TEMPEST, and the National Security Agency (NSA) has a list of products that are TEMPEST-compliant. Indeed, this vulnerability may not be a major problem in most organizational networks, but there is some set of networks where this is a concern.

An alternative to encasing workstations in Faraday cages (i.e., copper mesh layers surrounding all components) is to generate random electronic noise that masks meaningful radiated data transmissions.³

One way to reduce or eliminate EMR is to reduce or eliminate the amount of copper media on the network. Optical fiber, for example, has no EMR because there is no electricity on the line. It does not eliminate the EMR at the desktop itself, but it does prevent its entry into all interconnecting fiber cables.

Users also can be the source of some types of denial-of-service (DoS) attacks, either purposely or accidentally. Consider coaxial cable Ethernet networks (10BASE-5 or 10BASE-2) where nodes are attached to a common LAN medium. In both cases, the coaxial cable has a terminating resistor at the end of the wire to eliminate reflection of the signal. If the resistor is removed, it allows extra noise on the wire that can block all traffic on the network. End-of-cable resistors should be beyond the reach of users if at all possible.

Similar DoS attacks may occur whenever users modify the wiring scheme of the network. Removal of terminating resistors is only one action that can cause problems. Hubs in common areas might be unplugged or have network connectors removed. A token ring hub-to-hub connection can be detached, breaking the integrity of the ring and thus preventing LAN hosts from communicating with other hosts and denying service to users.

It is important that the physical network should be secured to the greatest extent possible. LAN managers should educate users to avoid the accidental problems that might occur and even how to recognize nefarious attacks.

25.4.3 Modems, Dial-Up Servers, and Telephone Lines. Modems anywhere on the LAN are a potential danger, particularly those that are connected directly to a user's system, with or without official sanction. Modems can provide a back door into the LAN, possibly bypassing the firewall, strong authentication, proxy server, and any other network security. Although less common in today's LAN architectures, dial-up modems are still used to dial into routers/bridges during network troubleshooting exercises. Care, as discussed in this section, should still be taken if dial-up modems are in use to access LAN devices.

In general, all modems should be concentrated at the network's dial-up server. Individual user systems should, in most cases, be banned from having modems on the network. This is a difficult rule to enforce, however, because laptops and an increasing number of desktop systems include preinstalled modems, and a user can always connect

25 · 6 LOCAL AREA NETWORKS

an external modem to the serial port of any system. This is an example of why security managers have to integrate policies into the culture of an organization. Otherwise, users will find a way around what they perceive to be prohibitive and onerous policies, and modems are one way to circumvent the corporate firewall.

Modems in auto-answer mode are particularly dangerous. Although most companies do not advertise their dial-in telephone numbers, these numbers do not remain secret for very long from someone who wants to find them. Anyone with a telephone directory can easily start mapping an organization's block of corporate telephone numbers. For example, if the main number is 802-555-3700, attackers have a place to start. When attackers call the main number and ask the receptionist for the organization's fax number, they obtain even more information. Using war-dialer software, attackers can scan an entire block of telephone numbers (e.g., 555-3700 through 555-3799) and obtain a list of which numbers are active, which respond with a tone, and what the tone represents (e.g., fax, modem, etc.). If a user has an auto-answer modem on a computer, attackers may gain access to the user's system without so much as a password. Security managers should work with their local telephone companies to obtain telephone numbers for modem lines that are not in the organization's telephone block.

The dial-up server, then, is the place to concentrate the modems and the authentication of dial-up users. There are several strategies for authentication at the dial-up server, the strongest being some form of two-factor authentication, such as a combination of passwords and a token. Another strong protection mechanism is to implement a dial-back mechanism, so that once a bona fide user logs in, the system hangs up and calls back to a preconfigured telephone number. This is a very effective scheme that works well with fixed-location telecommuters but not with roaming employees. In addition, attackers have been known to tamper with the central switch of a telephone company to call-forward from an assigned number to the attacker's own modem.

When a user requires two separate logons, one to the dial-up server and then one to a domain server, security restrictions should control what the caller is allowed to do after passing the first test. One of the authors of this chapter worked with a company that had a shared secret (an oxymoron) telephone number for its modem bank and then a single shared username and password for all the users to authenticate to the dial-up server. To access files and shared network resources, the user then had to authenticate to the domain controller. But after passing the identification and authentication for the first server, an attacker was on the organization's LAN and had complete, unfettered access to the Internet and an identity that implicated the company in any possible malfeasance.

Some guidelines for securing dial-up servers include:

- Maintain and monitor telephone logs.
- Configure all software and modems so that a user logout forces the modem to disconnect, and a modem disconnect forces a user logout.
- Configure the modems so that they return to their default configuration after every connection.
- Implement a dial-back mechanism, where possible.
- Use two-factor authentication for roaming users.
- Periodically scan internal telephone numbers for unauthorized modems.
- Prevent the display of any banner message when the user connects to the modem, and certainly do not display any sort of *welcome* message.

PHYSICAL LAYER ISSUES 25 · 7

- Train the organization help desk in social engineering techniques, and prohibit them from giving out modem telephone numbers, user names, or other sensitive information that could help attackers.

25.4.4 Wireless LAN Issues. Wireless LANs (WLANs) have vulnerabilities their wired counterparts do not. The most obvious difference between wired and wireless networks is the medium itself. Although copper-based LANs emit a small amount of radiation that can be intercepted, the entire basis of wireless LANs is transmitting data using relatively strong radiation in some form.

There are WLANs based on infrared signals that cannot penetrate building walls and so achieve some degree of security due to the limited propagation of those signals. Such LANs typically are found in networks requiring high levels of security. However, most WLANs today use radio transmission techniques. In these networks, anyone on a nearby street can use a listening device to intercept data and even capture the network identifiers required to connect to the LAN. The practical range of interception is governed by the inverse-square law for the signal strength (it declines as the square of the distance from the source) and by the sensitivity and signal-to-noise characteristics of the receivers.

Fortunately, there is a certain measure of security within the physical layer itself. The Institute of Electrical and Electronics Engineers (IEEE) 802.11-based LANs employ either direct-sequence spread spectrum (DSSS) or frequency-hopping spread spectrum (FHSS) techniques. As always, depending on range and noise levels, it is possible to eavesdrop. However, interpreting the signals is made more difficult by how DSSS and FHSS work.

To make sense of the transmissions, the receiver must know either the *chipping code* used in a DSSS network or the *frequency-hopping pattern* in an FHSS implementation. Without such information, the signal will appear to be nothing more than background noise in the industrial, scientific, and medical (ISM) radio bands to the illicit receiver. It is not an insurmountable problem for the would-be eavesdropper, but the work factor is greatly increased when compared to narrowband radio techniques. The spread spectrum approach also offers more reliability in the face of interference from denial of service (i.e., intentional jamming), as the signal is spread over a broad range of frequencies. Some vendor equipment also comes with software components that allow for tuning around interference.

For example, the IEEE 802.11g amendment uses the newer orthogonal frequency division multiplexing (OFDM) from 802.11a for higher data speeds, yet is backward compatible with 802.11b using DSSS, which was already using the same ISM frequency band. DSSS data rates of 1, 2, 5.5, and 11 Mbps are supported, as are OFDM data rates of 6, 9, 12, 18, 24, 48, and 54 Mbps. IEEE requires only mandatory data rates of OFDM using 6, 12, and 24 Mbps, regardless whether it is 802.11a or 802.11g OFDM.⁴

For greater privacy than that provided by the physical layer alone, the 802.11 standard includes an optional encryption method called Wired Equivalent Privacy (WEP). This technique is truly optional so not all vendors support the standard. WEP uses a 40-bit form of the RC4 algorithm by default, although some products support stronger, 128-bit versions. It is a good idea to choose a product that offers more than just the 40-bit version, as a 40-bit keyspace does not provide great security given today's computing power.

Because WEP does not offer strong encryption and does not describe a standard key-exchange mechanism, many vendors have implemented layer three tunneling methods, such as those found in virtual private networks (VPNs), to provide greater privacy. These

25 · 8 LOCAL AREA NETWORKS

VPN-based approaches generally employ other encryption processes (e.g., Microsoft Point-to-Point Encryption) as used in the Point-to-Point Tunneling Protocol (PPTP) that use longer keys than WEP and often support Public Key Infrastructure (PKI) or other key-exchange mechanisms. Some implementations also provide authentication through standards such as Remote Access Dial-In User Service (RADIUS) for more flexible client management. The major problem is that these approaches are not all interoperable and are not necessarily multiprotocol capable.

WEP also can be used for authentication to prevent unauthorized access to the WLAN itself. Such authentication adds another layer of protection to the username and password combination employed by typical server software. Before gaining access to information resources on the server, a client first must gain access to the physical medium. Using the shared key scheme, a wireless device must possess the same encryption key as the LAN's access point, the device enabling wireless connectivity to the wired portion of the LAN. Any data transmitted must be encrypted with the key, or the frame will be ignored. Many wireless access products also have the capability to create access control lists based on MAC addresses to filter LAN connections.

IEEE 802.11i-2004 or 802.11i security standard, implemented as Wi-Fi Protected Access II (WPA2), is an amendment to the original IEEE 802.11. WPA2 was originally developed by the Wi-Fi Alliance (www.wi-fi.org). This standard specifies security mechanisms for wireless networks. It replaced the short Authentication and Privacy clause of the original IEEE 802.11 standard with a detailed Security clause. In the process it deprecated the original WEP encryption method. The amendment was later incorporated into the published IEEE 802.11-2007 standard.

In summary, network managers should consider these security items when evaluating wireless network components:

Physical Layer Schemes

- **Infrared.** Cannot penetrate walls and is good for high-security applications.
- **FHSS.** Signal hopping provides good level of security but complexity of technique limits bandwidth.
- **DSSS.** Low “spreading ratio” can increase available bandwidth but also the possibility of interception and jamming.
- **OFDM.** Used for higher data rates in newer WLANs but may be more sensitive to noise and jamming.

Encryption Options

- **WEP.** Deprecated in 2004, commonly employed 40-, 64-, or 128-bit keys with the Rivest Cipher 4 (RC4) algorithm.
- **Wi-Fi Protected Access II.** A more secure alternative to WEP, employing the Advanced Encryption Standard (AES) Counter-Mode/CBC-MAC Protocol (CCMP).⁵
- **Alternative encryption methods.** Often more secure than WEP and WPA2 but not always interoperable.

Authentication Methods

- **Wired Equivalent Privacy.** Uses pre-shared key (PSK) mode, which requires clients to have the same encryption key as the LAN access point, which introduces key management issues.
- **Wi-Fi Protected Access II.** WPA2 supports IEEE 802.1X/EAP authentication or PSK technology.

NETWORK OPERATING SYSTEM ISSUES 25 · 9

- **Access control list.** Allows only certain clients to gain physical access to the LAN, based on the MAC address; this adds complexity to client management.
- **Server-based authentication.** Flexible user authentication with RADIUS for an additional layer of protection from illicit connections.

For more information about security for wireless networks, see Chapter 33 in this *Handbook*.

25.5 NETWORK OPERATING SYSTEM ISSUES. In the early 1990s, it was common to find desktop systems running the Windows operating system *and* the Novell NetWare network operating system (NOS). Desktop applications ran over Windows, and NetWare was used only to move files to and from the shared file space, or to print documents.

Today, the distinction among the desktop operating system, server operating system, and NOS has disappeared. Operating systems such as Linux, MacOS, UNIX, and Windows all provide desktop application suites with networking capabilities, including communications protocols such as TCP/IP. There are some general security considerations for all LANs regardless of the specific operating system:

- Use whatever capabilities are provided by the operating system to employ strong passwords.
- Create password policies that force strong passwords. Change passwords periodically and do not allow reuse. Periodically audit passwords using password-cracking tools such as L0phtCrack (Windows) or crack (UNIX). Ensure that the administrator and root accounts are given passwords that are not widely distributed or guessed.
- Disable (or uninstall) any services that are not being used.
- Keep the operating system and application software up to date, and with the latest security patches installed.
- Carefully manage access control lists for files and other system/network resources.
- Strictly define users, groups, and network/domain trusts.
- Tightly secure any running applications.
- Log on as administrator or root only when necessary; otherwise log on as a regular user.
- Allow operators and administrators to log on only locally at server systems.
- Limit use of guest, demo, or anonymous accounts.
- Where feasible, put boot and system files as well as application files and data on different partitions, hard drives, or input/output (I/O) controllers.
- Regularly audit server systems.
- Monitor log files.
- Remove floppy, CD, DVD, and thumb drives from servers after the system has a stable configuration.
- Implement best industry practices when securing the operating system.
- Use vulnerability assessment tools on a regular basis to scan servers.
- Use intrusion detection tools to monitor potential attacks on the LAN that are launched from the internal network.

25 · 10 LOCAL AREA NETWORKS

- If using the Simple Network Management Protocol (SNMP) for network administration, carefully choose community names and block external access to the SNMP service. Make management information bases (MIBs) read-only, where possible.
- Avoid use of the DNS HINFO (host info) resource record to identify the central processing unit type and the installed operating system.

Specific operating system vulnerabilities are beyond the scope of this chapter; entire books and Websites are devoted to securing some of these individual operating systems. At a minimum, network managers must monitor their operating system vendor's Website and all other sites that cover the NOS's security. The sections that follow provide some general observations and comments about the various network operating systems.

25.5.1 Windows 9x. All of the Windows operating systems (including NT and 2000) support peer-to-peer resource sharing and are vulnerable to exploitation of NetBIOS file and print sharing. In particular, when file and print sharing is enabled, the service is, by default, bound to TCP/IP. Although this does not cause an additional exposure to systems on the local network (since shares can be seen by other nodes on the LAN anyway), it does provide a potential vulnerability for hosts connected to the Internet. File and print sharing can be unbound from TCP/IP using Start, Control Panel, and Network.

Windows 9x (including Windows 95, Windows 98, and Windows ME) systems also have a vulnerability in the way authentication is performed when a user wishes to access a remote share. Windows uses the challenge-handshake authentication protocol (CHAP) for all passwords that need to be sent over the network, so passwords never appear in plaintext on the LAN. However, because Windows 9x uses the *same* challenge during a given 15-minute period, an intruder with a physical access to the LAN and to a sniffer could effect a replay attack by resending a duplicate authentication request and remapping a share on the Windows 9x system. This example illustrates the critical role of physical security in preventing compromise of LANs.

One of the classic Trojan horse programs for Windows 9x is Back Orifice (BO). Advertised as a remote Win9x administrator tool, it can be used by a nefarious user to take total control of someone else's system, including the capability to modify the registry, reboot the system, transfer files, view cached passwords, spawn processes, and create shares. SubSeven and NetBus are other tools that can be used to take control of a remote Windows system. Some commercial virus scanners can detect these programs on individual systems, and several workstation firewalls exclude communications characteristic of these Trojans.

Windows 9x has no particular logon security mechanism. Although every user might be forced to log on to a system, any user can log in with any username and any password to get at least basic access. Several password-cracking programs are available through the Internet to break Windows' .PWL password files, which are accessible once a user has access to the network. If a password-protected screen saver is preventing an attacker from logging in, there is an easy way around this as well: Simply reboot the computer. However, third-party security programs include nonbreakable secure logins and secure screen savers; many include bootlock functions that prevent any access whatever unless a valid password is entered. Current examples of such software can be located easily using buyers' guides such as the annual listing from the Computer Security Institute.⁶

NETWORK OPERATING SYSTEM ISSUES 25 · 11

25.5.2 Windows NT/2000, XP, Vista, Win7, and Win8. From a security perspective, Windows 2000 Millennium Edition (ME) is not significantly stronger than Windows 9x. Network administrators should periodically scan the network's public shares to ensure that they are appropriate. Windows NT Server, NT Workstation, and 2000 Server editions are built for security, and network services and have the software architecture to support these services. However, a security vulnerability announcement related to these operating systems seems to come out almost weekly. Many of the hacking tools available for Win9x also are available for Windows NT and 2000; Back Orifice 2000 (BO2K), for example, is an NT/2000 version of BO and NetBus also can take control of an NT/2000 host.

Scripting holes in Internet Explorer (IE) and Office 2000 make all Windows systems susceptible to a variety of new virus and worm attacks. Although the early viruses, such as Melissa and I LOVE YOU, required users to open e-mail attachments, that is no longer so. Microsoft Outlook and Outlook Express will execute Hypertext Markup Language (HTML) and script code in the body of an e-mail by default. Several ActiveX components also will execute from an e-mail containing HTML and script code; examples of such controls include Scriptlet.typlib (ships with IE 4.x and 5.x) and the UA control (Office 2000). The best protection against these types of vulnerabilities is to define Outlook and Outlook Express to read e-mail in the "Restricted Sites Zone" and disable all Active Scripting and ActiveX related settings in that zone. This vulnerability affects all Windows systems with Office 2000 or Internet Explorer 4.x/5.x installed, even if IE is not used.

Securing Windows NT/2000 systems is well beyond the scope of this chapter, but some of the precautions, in addition to those listed already, follow.

- Format the drive using NTFS rather than FAT.
- Use long file names and disable the DOS 8.3 naming format.
- Disable the Everyone group.
- Rename the administrator account.
- Turn auditing on. (It is off by default.)

All NT-based systems have been given the C2 security rating by the National Computer Security Center, which includes Windows NT 3.5 and 4.0 and Windows 2000 SQL Server version 8.0. This means that when these versions of Windows are installed correctly, they meet evaluation criteria set forth in the National Computer Security Center (NCSC) *Orange Book* of security specifications. C2 certification is not applied to the operating system itself; rather, it is applied to a particular installation. Microsoft provides tools to audit a site so administrator(s) can deploy the correct hardware and software configurations to achieve this level of security in a network.

Windows 2000 introduced a new feature that administrators might want to employ. Windows NT introduced the ability to compress and decompress files on the fly. Windows 2000 introduces the capability to encrypt and decrypt files on the fly. The Encrypting File System (EFS) uses public key cryptography and the Data Encryption Standard (DES) to encrypt files on the hard drive. EFS includes an optional capability for a recovery mechanism in case of key loss. Organizations using EFS on critical systems should consider employing this mechanism to protect against loss or destruction of the private key.

Windows XP, which has sold well over 400 million copies and still enjoys the second largest installed base of any Microsoft OS next to Win7, introduced a number of

25 · 12 LOCAL AREA NETWORKS

security-oriented features. The major improvement was the inclusion of a firewall capability in both the Home Edition and Professional versions. Professional goes a couple of steps further by offering EFS file encryption, Kerberos,⁷ smart card support, and a new software restriction feature allowing administrators to mitigate the impact of viruses and Trojan horses. XP also supports raw sockets, which in itself is not unusual—UNIX and Linux do as well. This feature is intended to increase the functionality of Internet services, but in Microsoft’s implementation, it is available to any user, no matter what privilege level. Thus hackers can possibly gain control of a computer running XP and use it to initiate DoS attacks by commanding the OS to generate a flood of traffic.

Although Microsoft’s Windows Vista was released with much fanfare—and after a number of delays—it has not been as widely implemented as might have been due to concerns about security and stability of the platform, not to mention a perception that many of the security features create a great deal of inconvenience in the user experience (e.g., User Account Control, to be described, and the Digital Rights Management implementation). Regardless, a prime mover of Vista development was Microsoft’s Trustworthy Computing initiative, and as such, a number of new security capabilities were added.

Features include User Account Control (UAC), Bitlocker Drive Encryption, Windows Defender, Data Execution Prevention, Application isolation, Windows Service Hardening, Network Access Protection, and a variety of others. UAC is the most visible from the user perspective, requiring user intervention before allowing any action that requires administrative-level privileges. These include⁸:

- Right-clicking an application’s icon and clicking “Run as administrator”
- Changes to files or folders in %SystemRoot% or %ProgramFiles%
- Installing and uninstalling applications
- Installing device drivers
- Installing ActiveX controls
- Changing settings for Windows Firewall
- Changing UAC settings
- Configuring Windows Update
- Adding or removing user accounts
- Changing a user’s account type
- Configuring Parental Controls
- Running Task Scheduler
- Restoring backed-up system files
- Viewing or changing another user’s folders and files

As an interesting historical aside, the National Security Agency assisted Microsoft in the development of Vista. As reported in the *Washington Post* in January 2007:

For the first time, the giant software maker is acknowledging the help of the secretive agency, better known for eavesdropping on foreign officials and, more recently, U.S. citizens as part of the Bush administration’s effort to combat terrorism. The agency said it has helped in the development of the security of Microsoft’s new operating system—the brains of a computer—to protect it from worms, Trojan horses and other insidious computer attackers. . . .

The NSA declined to comment on its security work with other software firms, but Sager said Microsoft is the only one “with this kind of relationship at this point where there’s an acknowledgment publicly.”

NETWORK OPERATING SYSTEM ISSUES 25 · 13

The NSA, which provided its service free, said it was Microsoft's idea to acknowledge the spy agency's role.

"I kind of call it a Good Housekeeping seal" of approval, said Michael Cherry, a former Windows program manager who now analyzes the product for Directions on Microsoft, a firm that tracks the software maker.

Cherry says the NSA's involvement can help counter the perception that Windows is not entirely secure and help create a perception that Microsoft has solved the security problems that have plagued it in the past. "Microsoft also wants to make the case that [the new Windows] is more secure than its earlier versions," he said.⁹

It is left as an exercise for the reader to decide whether having a spy agency working on a premier OS is a good thing or not.

One of the most important capabilities for the network security manager is to audit the Windows server systems to protect their integrity. These tools are part of the base operating system or the Windows NT Resource Kit:

- *netstat* examines open ports.
- *Event Viewer* examines application, security, and system logs.
- *net start*, *net user*, *net group*, *net local group* display running services, users, groups, and local groups.
- *dumpel* converts Event Viewer logs to simple text files.
- *NetMon* displays network traffic.
- *netsvc* displays local and remote running services and drivers.
- *addusers* displays users and groups.
- *findgrp* displays local and domain groups for a user.
- *local* and *global* show all members of specific local or global groups.
- *dommon* displays trusted domains.
- *xcacs* examines the file Access Control Lists (ACL).
- *perms* examines the ACLs associated with a user.
- *sysdiff* displays changes in the Registry and file system.
- *regdmp* creates an ASCII version of the Registry.
- *ralist* lists a domain's Remote Access Servers (RAS).
- *rasusers* lists users authorized for dial-in access.

According to Tech Review Source ([www.techreviewsource.com /Windows-7/windows-7-security-features#.UL-S3WdN_f0](http://www.techreviewsource.com/Windows-7/windows-7-security-features#.UL-S3WdN_f0)), Windows 7 was the most secure OS Microsoft had released to date. Security features include the new Action Center, improved UAC, Parental Controls, Internet Explorer 8 Security, Windows Defender, and Microsoft Security Essentials. Windows 8 added many changes to the operating system but they were primarily in the user interface and user experience; Win8 capitalized on Win7's newer security capabilities.

25.5.3 UNIX. UNIX is the oldest operating system still in widespread (and growing) use today. Readers needing guidance for Novell Netware should refer to Chapter 18 in the fourth edition of this *Handbook*. Originally developed in 1969 at AT&T Bell Laboratories, UNIX became the first operating system to integrate network communications when TCP/IP was bundled into Berkeley Software Development

25 · 14 LOCAL AREA NETWORKS

(BSD) 4.2 UNIX in 1984. UNIX had traditionally been reserved for server systems and hardcore computer users. With the development of the X-Windows interface for UNIX and the wide deployment of Linux since the mid-1990s, UNIX and its variants represent the only significant competition to Windows in the desktop and server environment.

Like TCP/IP and the Internet itself, UNIX was developed for functionality and use within a trusted user community. As such, while UNIX has many powerful tools, it does not have a cohesive security architecture, nor is it an inherently secure operating system.

UNIX has most of the basic operating system protections: passwords, access control lists, groups, user privilege levels, and so on. But UNIX also comes with a large variety of services (daemons) enabled by default, including *File Transfer Protocol (FTP)*, *Telnet*, *finger*, *echo*, *chargen*, *daytime*, *Remote Procedure Call (RPC)*, *BIND*, and more. In addition, nearly every UNIX daemon has had some sort of security vulnerability reported at one time or another, with buffer overflows being quite prevalent.

There are many things that an administrator should consider when securing a UNIX/Linux system. In addition to the general steps just listed, the security manager might also:

- Disable (or remove) any unused services, particularly *finger*, the BIND name daemon (*named*), *RPC*, *sendmail*, Trivial FTP (tftp), Post Office Protocol (POP), Internet Message Access Protocol (IMAP), *sadmind*, *mountd*, and Network File System (NFS).
- Install the latest version and security patch of all installed software.
- Take great care when configuring access control lists and other sharing.
- Prevent running *Sendmail* in daemon mode (turn off the -bd switch) on machines that are neither mail servers nor mail relays.
- Limit use of the “r” remote access protocols.
- Use shadow password files.
- Implement TCP Wrappers to control access to services.
- Consider using encrypted communication protocols, such as Secure Shell (SSH) or Secure Sockets Layer (SSL), for remote access. Prevent transmission of cleartext passwords over the Internet.

One of the most important capabilities for the network/security manager is to audit the Windows server systems to protect their integrity. These tools are part of the base operating system or the Windows NT Resource Kit:

- *netstat* examines open ports.
- *lsof* displays hidden file space and network connections.
- *tcpdump* displays network traffic.
- *who* displays users that are logged on and the *utmp* log file.
- *last* displays login history, and the *wtmp* log file.
- *lastb* display a history of bad logins (and the *btmp* log file).
- *syslogd* is a central server facility for managing and logging system messages.
- *TCPWrapper* monitors and manages incoming service requests.

NETWORK OPERATING SYSTEM ISSUES 25 · 15

25.5.4 MacOS. The Macintosh operating systems have mostly given way to other platforms. That said, they are still worth mentioning, as Apple's market share across the board has been buoyed somewhat in recent years by its strength in mobile and tablet devices. The Macintosh was the first desktop operating system that included networking and resource sharing via AppleTalk as integral elements. But like UNIX and TCP/IP before it, the MacOS is designed for convenience and usability but not for security.

Because of its peer-to-peer nature, MacOS has traditionally had a number of potential exposures. Although Windows and UNIX also could operate in a peer-to-peer mode, a novice user generally would not know how to share resources and thereby might not inadvertently open holes.

For example, a nefarious Mac network user could quickly see what servers and shares were available on the network by using the original *Chooser* accessory. In desktop use, Macs had relatively little security. Password protection was provided by default only with some laptop systems, and not even a password-protected screen saver came standard with the system. In short, there was very little standing between a determined attacker and a Mac computer. Third-party software was required to provide password protection against access to the system and files, or for data protection with disk encryption.

Mac-based viruses and worms continue to be much less prevalent than their Windows counterparts, but Macs are not totally immune. First, those viruses that depend on Microsoft Office Suite software will work because the Mac versions of Word and Excel employ macros. Second, Internet-based attacks aimed at TCP/IP—such as the Ping-of-Death, Teardrop, and SMURF—can still affect a Mac server. There are more choices for add-on security than in the past, with packages like Norton™ Antivirus for Mac® and Norton™ Internet Security for Mac® being similar to their counterparts for other operating systems.

The Mac OS has undergone a fundamental transformation since the last “classic” release of OS 9 back in 1999. With the advent of OS X, Apple introduced its first UNIX-based platform, which has been preloaded on all Macintosh computers since 2002.

Version 10.5 (Leopard) was released in 2007 and added a number of security enhancements to address traditional platform weaknesses, including:

- Secure guest account
- Application Layer firewall
- Application signing
- Sandboxes
- Full disk encryption
- Library randomization

Apple's intention was to provide for greater pre-emption of attacks as well as resiliency to them, and the company has continued its development of better security throughout OS X's lifecycle. Currently on version 10.8 (Mountain Lion), which was released in 2012, the platform includes many of the same features as previous iterations, with improved administrative controls and privacy options.

One of the biggest security challenges is conceptually not unique to Macintosh, but with the increased integration of Apple's iCloud services in OS X (not to mention iOS in the mobile space), there is a greater risk of files being moved from local to network

25 · 16 LOCAL AREA NETWORKS

storage. Users expect access to data from any device, anywhere at any time, and that convenience makes all lost or stolen laptops, tablets, and phones significant risks.

While the MacOS's fortunes have waxed and waned and waxed again, there still are fewer security incidents reported because the greater popularity of other platforms make it easier to learn about them, because there are more potential targets, and because one single attack can affect more systems. Put another way, there are fewer attacks on Macs because the hacker community is not so familiar with them, and there are fewer attractive targets worth the opportunity cost. Still, it is just as important to keep the version of MacOS up to date as with other operating systems.

25.6 CONCLUSION. A good administrator can secure almost any NOS, although no NOS is secure initially. The network administrator needs continuous vigilance and monitoring, while recognizing that the operating system is only a part of the overall security plan for the LAN and network services. Most network administrators, due to the nature of their job and training, focus exclusively on the computers attached to the LAN and to the LAN's operating system and software. Unfortunately, this approach is too narrow in its scope. Personal firewall software also might be employed to protect individual systems against attack, but almost all of these products are oriented toward IP-based attacks and miss attacks that employ the NOS's native operating system.

Routers, network firewalls, and proxy servers are essential for protecting LAN systems from attack by an external source. The network administrator also must provide tools to protect servers and workstations from other users on the LAN.

25.7 FURTHER READING

This section lists some books, articles, and Websites that cover the issues addressed in this chapter. Administrators should monitor vendor, operating system, and security Websites that will have up-to-date, timely additional information.

- Anonymous. *Maximum Security*, 4th ed. Indianapolis, IN: SAMS: 2003.
- Bathurst, R., R. Rogers, and A. Ghassemlouei, *The Hacker's Guide to OS X*. Amsterdam: Elsevier, 2013.
- Brenton, C., and C. Hunt. *Mastering Network Security*, 2nd ed. San Francisco: SYBEX, 2002.
- Edwards, M. J. *Internet Security with Windows NT*. Loveland, CO: Duke Press, 1998.
- Fraser, B., ed. *Site Security Handbook*. RFC 2196/FYI 8. September 1997.
www.ietf.org/rfc/rfc2196.txt
- Garfinkel, S., G. Spafford, and A. Schwartz. *Practical Unix and Internet Security*, 3rd ed. Sebastopol, CA: O'Reilly & Associates, 2003.
- IEEE Working Group for WLAN Standards. <http://grouper.ieee.org/groups/802/11/>
- The Information Warfare Site (IWS). Rainbow Series Library page. www.iwar.org.uk/comsec/resources/standards/rainbow/rainbow.html
- Landau, T. *Sad Macs, Bombs, and Other Disasters*, 4th ed. Berkeley, CA: Peachpit Press, 2000.
- L0pht Heavy Industries. Overview of AntiSniff. 1999. <http://blockyourid.com/~gbpprorg/l0pht/antisniff/overview-2.html>
- Mann, S., and E. L. Mitchell. *Linux System Security: The Administrator's Guide to Open Source Security Tools*, 2nd ed. Englewood Cliffs, NJ: Prentice-Hall, 2002.
- McClure, S., J. Scambray, and G. Kurtz. *Hacking Exposed: Network Security Secret & Solutions*, 6th ed. Berkeley, CA: Osborne/McGraw-Hill, 2009.

NOTES 25 · 17

- McNamara, J. The Complete, Unofficial TEMPEST Information Page. October 2, 2000. www.jammed.com/~jwa/tempest.html
- Payne, W. H., and T. Sheldon. *The Complete Reference to Netware 5*. New York: McGraw-Hill, 1999.
- Thurrott, P., and R. Rivera. *Windows 8 Secrets*. Indianapolis, IN: John Wiley & Sons, 2012.
- University of Connecticut. (2004, June 17). IT Policies, Standards, and Guidelines: LAN Security Guidelines. University Information Technology Services Website. June 17, 2004. http://uits.uconn.edu/?page_id=499
- Wi-Fi Alliance. WiFi Alliance home page. 2013 www.wi-fi.org

25.8 NOTES

1. See www.ietf.org/rfc/rfc2196.txt
2. See www.asystematics.com/asusys/c3-java-saf.htm, www.oracle.com/technetwork/java/index-jsp-141438.html, www.ehow.com/how_7160897_run-activex-controls-plugins.html, and <http://windows.microsoft.com/en-US/windows-vista/Should-I-install-ActiveX-controls>
3. For details of software-based TEMPEST, see Markus G. Kuhn and Ross J. Anderson, "Soft Tempest: Hidden Data Transmission Using Electromagnetic Emanations," 1998; www.cl.cam.ac.uk/~mgk25/ih98-tempest.pdf
4. See <http://kl2217.wordpress.com/2009/07/31/ieee-802-11-standards-comparison/> and <http://janmagnet.files.wordpress.com/2008/07/comparison-ieee-802-standards.pdf> for WLAN physical layer details.
5. See www.wi-fi.org/files/wp_9_WPA-WPA2%20Implementation_2-27-05.pdf for WPA and WPA2 security details.
6. See www.securityinfowatch.com/directory for SecurityInfoWatch.com's Buyers Guide of Computer Security Products.
7. Microsoft came under fire for its original implementation of Kerberos in Windows 2000. Microsoft's version used a piece of proprietary data in the authentication process, making it impossible for third-party Kerberos servers to work with Windows. The company has since released interoperability information.
8. Ed Bott, "What Triggers User Account Control Prompts?" February 2, 2007; www.edbott.com/weblog/2007/02/what-triggers-user-account-control-prompts
9. Alec Klein and Ellen Nakashima, "For Windows Vista Security, Microsoft Called in Pros," *Washington Post*, January 9, 2007; www.washingtonpost.com/wp-dyn/content/article/2007/01/08/AR2007010801352.html

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER **26**

GATEWAY SECURITY DEVICES

Justin Opatrny

26.1	INTRODUCTION	26·2	26.3	NETWORK-SECURITY MECHANISMS	26·9
26.1.1	Business Requirements		26.3.1	Allowed Paths	26·9
	Outpacing Security	26·2	26.3.2	Tunneling	26·10
26.1.2	Demand-Driven Processing	26·3	26.3.3	Anti-Spoofing	26·10
26.1.3	Ubiquitous Mobility	26·3	26.3.4	Network Address Translation	26·11
26.1.4	Regulatory and Industry Compliance	26·3	26.3.5	Intrusion Detection	26·11
			26.3.6	Intrusion Prevention/ Response	26·11
			26.3.7	Encryption	26·13
			26.3.8	Identity-Based Enforcement	26·14
26.2	BASIC CONCEPTS AND TERMINOLOGY	26·4	26.3.9	Complex Protocol Engines	26·14
26.2.1	General Capabilities	26·4	26.4	Content Control and Data Leakage Prevention	26·15
26.2.2	Unified Threat Management	26·4	26.3.10	IPv6	26·16
26.2.3	Next-Generation Firewall	26·4	26.3.11	Additional Considerations	26·17
26.2.4	Web Application Firewalls	26·4	26.4	DEPLOYMENT	26·18
26.2.5	Firewall Architectures		26.4.1	Zoned Architecture	26·18
	Changing	26·4	26.4.2	GSD Positioning	26·19
26.2.6	Packet Filtering	26·5	26.5	MANAGEMENT AND MONITORING STRATEGIES	26·20
26.2.7	Stateful Inspection	26·5	26.5.1	Monitoring	26·20
26.2.8	Application Layer Gateway	26·6	26.5.2	Policy	26·21
26.2.9	Current Gateway Security Devices	26·6	26.6	MAINTENANCE	26·22
26.2.10	Host Environment Context	26·6	26.6.1	Patching	26·22
26.2.11	Firewall Platforms	26·7	26.6.2	Pattern Updates	26·22
26.2.12	Routing Devices	26·7	26.6.3	Logging and Alerting	26·23
26.2.13	Appliances	26·7			
26.2.14	Virtualization	26·8			
26.2.15	Host-Based Agents	26·9			

26 · 2 GATEWAY SECURITY DEVICES

26.6.4	Secure Baseline Configuration	26·23	26.8.4	Throughput Implications of Added Features	26·26
26.6.5	Default Configurations	26·23	26.8.5	Management	26·27
26.6.6	Harden the Device	26·24	26.8.6	Logging and Alerting	26·28
26.7	DISASTER RECOVERY	26·24	26.8.8	Auditability	26·29
26.7.1	Fail-Over/High Availability	26·24	26.8.9	Disaster Recovery	26·29
26.7.2	Load-Balancing	26·24	26.8.10	Usability	26·30
26.7.3	Backup/Restore	26·24	26.8.11	Learning Curve	26·30
26.8	NETWORK-SECURITY DEVICE EVALUATION	26·25	26.8.12	Features	26·30
26.8.1	Current Infrastructure Limitations	26·25	26.8.13	Price	26·31
26.8.2	New Infrastructure Requirements	26·25	26.8.14	Vendor Considerations	26·32
26.8.3	Performance	26·26	26.8.15	Managed Security Service Providers	26·33
			26.9	CONCLUDING REMARKS	26·34
			26.10	FURTHER READING	26·35

26.1 INTRODUCTION. Once considered sufficient to protect an entire organization from external threats, the firewall is still perhaps the most recognized and deployed network-security devices for Internet-connected operations. However, earlier firewall generations made security decisions with little contextual support other than the origin and destination of the packets traversing a particular allowed path.

As communications capabilities and functionality demands increased, so too did the firewall's need to inspect and enforce allowed paths using more complex protocols and require ever-increasing throughput. This evolution transformed the firewall into a true gateway security device (GSD)—able to provide allowed path enforcement using a combination of techniques which once required additional security devices to accomplish.

Properly selected and deployed GSDs are one security layer designed to handle these increasingly complex scenarios. The GSD is effective only with a full understanding of the capabilities and limitations—both operational and failure conditions—of consolidating multiple security functions into a single device. By providing allowed path enforcement more intelligently and accurately, combined with the added rigor of genuinely understanding expected network flows, GSDs provide sufficient additional defense-in-depth layers throughout the organization.

Although this chapter focuses on the GSD as a combined security device, the concepts covered are useful for understanding and evaluating the functionality of individual network-security devices. Every organization must make risk and performance decisions by weighing this approach against maintaining independent devices that focus on a particular security function.

26.1.1 Business Requirements Outpacing Security. Technological advancement continues to transform an enterprise's ability to manage the data life cycle. The pervasiveness of mobile devices and move toward cloud-computing resources that are no longer solely controlled or consumed by the organization continues to redefine the perimeter. Users are increasing demands for unfettered access to corporate data from anywhere on any device (including those not controlled by the organization).

INTRODUCTION 26 · 3

This dynamic environment increases the need for layered security architectures with deeper awareness of content and context.

26.1.2 Demand-Driven Processing. Enterprises not possessing the requisite internal human or technological expertise to achieve the organization's goals for information technology (IT) have long looked to outsourced solutions to meet their needs. Software as a Service (SaaS) provides offerings such as productivity applications, collaboration, and email (e.g., Microsoft Office365) and customer relationship management (CRM) (e.g., salesforce.com). Infrastructure- and Platform-as-a-Service (IAAS and PAAS, respectively) provide on-demand storage and computing (e.g., Amazon Web Services—S3/EC2 or Rackspace Open Cloud). Outsourced offerings continue to mature and redefine how enterprises develop, manage, and present their information.

Virtualization technology continues to provide opportunities to use internal processing capabilities more efficiently. Although providing better performance, it can also diminish security in that existing offerings are not necessarily as mature as dedicated security infrastructure. When internal and Internet facing systems must operate on the same virtual architecture, this creates additional risk as systems now interact at the hypervisor and virtual switch (vSwitch) level where traditional network-security devices are unable to inspect and enforce traffic at this level.

See Chapter 68 in this *Handbook* for more details about security and outsourcing.

26.1.3 Ubiquitous Mobility. Today's business climate demands the ability for employees to work from anywhere, and this need for mobility and flexibility continues a significant shift in how organizations define and protect their perimeters. Employees may use a variety of systems whether at work, home, or on the road. The level of access and functionality required extends well beyond email into enterprise applications and data. As functional mobility stretches from company-owned to personal devices, organizations must have a method to ensure a compromised mobile device does not weaken the existing internal or outsourced security controls.

26.1.4 Regulatory and Industry Compliance. Federal regulatory requirements continue to have a significant impact on how organizations manage risk through data protection, retention, and privacy activities—many having significant auditing and reporting requirements.

- The Sarbanes-Oxley Act (SOX) focuses on controls and procedures designed to preserve the integrity of publically traded organizations' financial reporting.
- The Gramm-Leach-Bliley Act (GLBA)—specific to financial institutions—concentrates on the protection of customer data and privacy.
- The Health Information Portability and Accountability Act (HIPAA) requires protection of individually identifiable health information such as personally identifiable information (PII) and protected health information (PHI).

Industry-specific requirements continue to appear and evolve in an effort to address the minimum security requirements to operate with or within a specific industry. The Payment Card Industry Data Security Standard (PCI DSS) establishes baseline requirements for the protection of cardholder data during processing, transmission, and storage. PCI DSS requires organizations to determine their level of involvement

26 · 4 GATEWAY SECURITY DEVICES

with the cardholder data. Once established, this determines which requirements are necessary to demonstrate compliance with the standard.

Each critical infrastructure sector—as defined in the United States by the Department of Homeland Security—has a Sector-Specific Plan (SSP), which provides representative organizations and agencies with the risk-management tools and strategies specific to the protection of each industry. One of the more mature programs is North American Electric Reliability Corporation (NERC) Critical Infrastructure Protection (CIP). NERC CIP addresses multiple physical and digital security elements of the North American power system, including gateway protection through the establishment of an electronic security perimeter.

See Chapter 64 of this *Handbook* for more details on GLBA, SOX, and PCI DSS.

26.2 BASIC CONCEPTS AND TERMINOLOGY. Greater demands for mobility, new models for business interaction, and leveraging Internet-based processing capabilities continue to force the evolution of traditional perimeter protections.

26.2.1 General Capabilities. With the substantial processing requirements of specialized network security systems (e.g., IPS or anti-spam), organizations traditionally architected security infrastructures that leveraged dedicated solutions/devices for each function. However, current generation GSD processing capabilities provide the opportunity to combine many of these once dedicated systems into a single device.

26.2.2 Unified Threat Management. Unified threat management (UTM) combines elements such as anti-malware, anti-spam, IDS/IPS, VPN, application proxy, and content filtering—transforming the firewall into the original iteration of the GSD. These added capabilities allow the UTM to provide greater control and inspection at the application layer. However, typically there was only marginal management and performance integration between feature sets.

26.2.3 Next-Generation Firewall. The next-generation firewall (NGFW) is the latest evolution in stated capabilities that complements and surpasses those of the UTM. Instead of just bolting multiple security technologies on top of one another, the NGFW provides tighter integration of each level of security. These new capabilities include greater protocol awareness and more granular allowed path enforcement. The NGFW is able to profile protocols regardless of port chosen. This increases the ability to detect deceptive behavior such as encrypted payloads over protocols that would not normally use encryption. This generation of GSD also has the ability to adjust policy dynamically—extending protection to other parts of the security infrastructure.

26.2.4 Web Application Firewalls. GSDs provide additional capabilities to inspect and enforce allowed paths for Web-based communications. However, the complexity of the current and next-generation Web protocols may outpace the security provided by this device. The Web application firewall (WAF) provides more robust HTTP protocol inspection capabilities. The WAF provides customizable rules to protect against common payload-based attacks such as SQL injection, XSS (cross-site scripting), and command injection. This platform can also serve as a virtual patch for Web applications—legacy (no patches possible) or ones awaiting unreleased patches.

26.2.5 Firewall Architectures Changing. As security vendors work to keep pace with these changes with more functionality and higher performance, it is the

BASIC CONCEPTS AND TERMINOLOGY 26 · 5

customer's responsibility to understand the advantages and disadvantages each proposed protection solution. Such analysis provides the necessary insight for the organization to deploy the most appropriate architecture to meet the necessary security and performance requirements.

26.2.6 Packet Filtering. Routing devices gave rise to first-generation firewalls capabilities. Packet filtering is a set of explicit rules describing the allowed paths network traffic may travel. The rules, in the form of an access control list (ACL), independently evaluate one or more portions of each packet's header to make the allowed path decision.

The packet filter acts on each packet as an individual entity without respect of the packet(s) that come before or after. Although this method provides the security with the least overhead of other firewall architectures, it is vulnerable to several network and transport layer attacks.

Internet Protocol (IP) *spoofing* is specially crafting a packet in an effort to deceive the router into accepting traffic that appears legitimate. The attacker will configure a packet to look like it originated from the internal network even though it is coming into the external interface. A *Land* attack sends a spoofed TCP SYN packet to a host with the source and destination IP and port being equal and can cause a denial of service on a vulnerable network stack. The *Teardrop* attack intentionally fragments a packet and manipulates the fragment offset where the preceding fragment's offset overlaps with the offset of the next fragment. When the receiving system reassembles the overlapping fragments, the network stack will crash, causing a denial of service.

Although both the Land and Teardrop attacks require a vulnerable end-device, initial packet filters did not have contextual understanding to stop such attacks. Mitigation for the Teardrop attack includes implementing packet reassembly—recombining the fragments into the original packet to ensure the recombined packet does not violate basic IP packet specifications—before forwarding to its next hop.

For details of these and other denial-of-service attacks, see Chapter 18 in this *Handbook*.

Current generations GSDs use stateful inspection (covered next) in an effort to overcome many of the limitations of packet filtering. However, many nonsecurity-related devices use packet filters to provide basic protection of their administrative interfaces and monitoring functions.

26.2.7 Stateful Inspection. Stateful inspection identifies and tracks additional parameters within each packet, adding context by representing flows along allowed paths as network connections instead of individual packets. Even though IP and User Datagram Protocol (UDP) are connectionless protocols, the firewall creates a virtual connection to emulate its connection-oriented counterparts. Although the specific parameters tracked in the state table vary across GSD vendors, the table resides in memory for more efficient processing.

Since a packet filter inspects packets individually, it does not affect firewalls in a load-balanced or fail-over architecture. However, to ensure sustained connections in high availability architectures, stateful inspection must establish and maintain a synchronized copy of the state table in each of its high availability members. Although some vendors still maintain that a direct serial connection is the most reliable connection method, highly available solutions are unlikely to be in a physical proximity necessary for serial connection and typically take advantage of existing network connections to maintain the state synchronization.

26 · 6 GATEWAY SECURITY DEVICES

26.2.8 Application Layer Gateway. As threats became more sophisticated across existing allowed paths, firewalls added application-layer gateway (ALG) protection. Each application-specific proxy (e.g., HTTP, RPC, FTP, etc.) acts as the allowed path broker for the connection, creating a separate, backend connection to the other host. Although the full connection is actually two independent sessions, this duality is transparent to the requester and the server. Being a brokered connection adds an additional layer of defense by re-writing potentially malformed requests, validating/enforcing protocol compliance, and not allowing direct communication between client and server.

The detailed payload analysis increases the amount of time the firewall spends evaluating each packet potentially reducing throughput for other network traffic and increasing the need for sufficient processing resource. The ALG may also create issues such as poor performance for protocols requiring minimal overhead in addition to unexpected compatibility issues due to vendor-specific protocol implementations. Active review and management will help to minimize the ALG's impact on operations.

26.2.9 Current Gateway Security Devices. The firewall evolved yet again as attacks continued to increase in complexity and scale. Stateful filtering and application layer gateway functionality being effective at protecting against specific known attack vectors. However, the sophistication of the newer generation attacks quickly revealed those two features were unable to provide the necessary level of protection. The additional capabilities integrated into the firewall gave rise to the GSD.

The initial iteration of GSD is known as unified threat management (UTM). UTM is the consolidation of the firewall with multiple additional security platforms (e.g., network intrusion-prevention systems, content filter, anti-spam, etc.) into a single device. UTM vendors that only had a strong presence in one or two of the protection measures had to either integrate a third-party security product or build their own offering to fill out the full protection suite. This was also an opportunity to consolidate management and monitoring capabilities to create a unified platform. However, the UTM would prove somewhat inefficient, as each packet travels serially through the different security engines, severely impacting performance as the number of protections enabled increases.

The most current generation of GSDs—known as next generation firewalls (NGFWs)—provide greater application detection, awareness, and enforcement in addition to tighter integration of all security layers. Increasing the processing and inspection capabilities alone were not enough to meet the demands of the current generation application protocols. These protocols are more complex in form and function, even though they may follow the traditional allowed path such as HTTP. The NGFW must also be able to decode and understand the application's inner workings regardless of port and enforce on all or subsets of the protocol. NGFWs process packets through selected security engines in parallel to increase performance and focus the protection needs to the specific type of traffic detected.

26.2.10 Host Environment Context. Although this chapter focuses on gateway security devices from a network perspective, comparable host-based protections warrant coverage, as they are an important defense-in-depth component and have certain capabilities that the network-based GSDs cannot provide. With mobility and ubiquitous information access being central themes surrounding today's personal and business environments, the host-level protection is an ever-present fixture. The host's

BASIC CONCEPTS AND TERMINOLOGY 26 · 7

movement from relying on the network security measures to local context and environmentally aware security measures will provide more robust and flexible security wherever the host goes.

Host-based security—aside from basic firewalling—has additional complexity and considerations due to the protection requirements of the additional elements (software, services, etc.) running. The host's ability to understand its local environment allows for more granular protection and visibility. Instead of focusing purely on port-based network access, the increase environmental context defines what applications and services can: access or receive network traffic, access or change other services, or execute code in a virtual machine (VM) to verify expect behavior(s).

26.2.11 Firewall Platforms. GSDs continue to advance in form, function, and contextual understanding to provide allowed path enforcement at all layers of the network stack. Malicious actors continue to develop increasingly targeted and effective code that travels on an otherwise allowed path. These more detailed network inspection requirements dictates that every organization select the most appropriate platform to meet its security and performance needs.

26.2.12 Routing Devices. As the need to transmit data within a network and beyond increased, routing device were adapted to provide native and then modular security services.

26.2.12.1 Access Control Lists. Routing device ACLs typically provide rudimentary allowed path enforcement with minimal performance impact. ACLs permit or deny individual network packets based on a combination of parameters from the header such as source and destination IP addresses and/or TCP/UDP port.

Routing device-based security services matured to provide enhanced capabilities based on their contextual understanding of the network. Routing devices can support several stateful inspection techniques, the most rudimentary being the ability to track connections based solely on the acknowledgement (ACK) flags of an established connection. However, this method is more susceptible to spoofed packets directed at the intended target than traditional stateful inspection based on a table of established connections. Routing platforms can also provide anti-spoofing capabilities based on network topology. When enabled, the routing device will make forwarding decisions based on its understanding of connected and learned routes and only allow source IPs to originate from the appropriate interface.

26.2.12.2 Hardware Modules. With bandwidth increasing into the tens to hundreds of Gigabits per second (and beyond), routing device manufacturer needed to offload the security functions to a dedicated vendor-specific module (or line card). The module provides logical protection for one or more allowed paths while using the backplane to achieve higher data rates and provide protection based on one or more of the firewall architectures described above. Some routing vendors also support modules—similar to a blade server—allowing other security vendors to integrate their technology deeper into the routing infrastructure.

26.2.13 Appliances. The GSD appliance emerged as a way to shed complex security decisions from routing devices onto a dedicated security platform. Appliance form factors, capabilities, and limitations vary among GSD vendors with the GSD application integrating into the chosen hardware platform.

26 · 8 GATEWAY SECURITY DEVICES

26.2.13.1 Proprietary. A proprietary hardware-based GSD appliance manufacturer is able to establish and retain tight controls over all aspects of the platform. The hardware is purpose-built to interoperate with the GSD software and provides different levels of performance and functionality. High-performance components such as application-specific instruction sets are typically proprietary. However, the vendor will balance the cost effectiveness of the platform by choosing commercial-off-the-shelf (COTS) hardware for the utility components (e.g., network interfaces and RAM). The accompanying proprietary, hardware-specific operating system (OS) allows the vendor to maintain tight controls on interoperability, configuration, and software revisions.

26.2.13.2 Hardware Independent. Hardware independent GSD appliances transfer responsibility to the customer to procure hardware that meets a specific set of requirements while allowing the GSD vendor to retain responsibility for developing the GSD application and/or OS. The GSD uses either a customized version of an OS (typically Linux or Berkeley Software Distribution [BSD]) or installs directly onto a customer-provided base OS.

A GSD application installed on a proprietary (or premodified) OS reduces native security exposures through vendor custom hardening of OS mechanisms to remove nonessential features and functionality while maintaining the tight controls described in the section above.

The GSD application may also be able to install on specific versions of the customer's native OS choice, providing the customer the most flexibility in choosing a GSD appliance. The GSD application will make some hardening changes but leaves much of the operating system management to the customer. This increases the risk of insecure and/or underperforming implementations due to the GSD vendor not fully controlling OS configurations, patches, and interoperability.

26.2.13.3 Soft Appliance. The soft appliance augments the hardware independent/proprietary OS GSD into a bootable platform that can run on consumer-grade equipment. This platform is more suited for point-in-time protection needs, such as main GSD device failure, fast-turn deployments, and testing. Before deploying a soft appliance, the customer must temper the potential "flexibility" of this solution with other considerations such as manageability (configuration, logging, etc.), scalability, and required functionality.

26.2.13.4 Embedded Appliance. In addition to high-end, scalable devices, GSD vendors also recognized the opportunity to service other markets such as providing specialized capabilities on platforms scaled to meet the needs of consumers and small-to-medium businesses (SMBs).

SMBs are less likely to have sufficient funding or dedicated security staff to purchase, manage, and monitor multiple security devices. The SMB-grade GSDs provide a more economical all-in-one device, due to its smaller form factor, simplified user interface, and lower performance requirements. In addition to providing the standard wired and/or wireless connectivity options, the device typically provides a full complement of UTM security functionality.

26.2.14 Virtualization. GSD vendors now provide solutions ranging from leveraging the virtual environment to provide independent network security capabilities to embedding granular protections at the hypervisor level. Instead of requiring

NETWORK-SECURITY MECHANISMS 26 · 9

a full installation, some GSD vendors will provide a preinstalled virtual appliance that minimizes total deployment time. It is important to ensure that the virtualized GSD maintains its priority in obtaining and sustaining resources from the VM instance.

A VM-capable GSD acts as an independent network security entity. Instead of having dedicated hardware, the GSD leverages the hypervisor's hardware abstraction layer to gain access to the necessary resources. This implementation provides hardware appliance-like capabilities to control the allowed paths based on what networks route through the GSD VM instance.

A VM-embedded GSD virtual appliance integrates with the hypervisor to gain low-level access to the entire VM host/cluster. The GSD now has visibility into previously inaccessible communication paths and is able to conduct allowed path enforcement for both inter- and intra-VM communications.

26.2.15 Host-Based Agents. As mobile devices demand greater internal network access—including data to be stored local to the device—so did the requirement for those host-based protections to be able to minimize the additional risk of allowing such access. Antivirus capabilities, although arguably only marginally effective against new and emerging threats, still maintain a necessary presence to filter known attacks vectors. Firewalling capabilities at the network and application level reduce the host's network-facing threat profile. The Host Intrusion Prevention System (HIPS) provides protection against known and potentially unknown threats because of its ability to inspect up to the application layer.

When network-security devices fail to detect a threat, host-based protection has additional measures to detect local anomalous behavior due to its increased contextual awareness of how the host—the applications and the operating system—should behave, thus allowing the host an advantage in detecting emerging threats. However, organizations should not underestimate the administrative overhead necessary to manage and monitor this protection.

See Chapter 27 of this *Handbook* for more details on host-based protections.

26.3 NETWORK-SECURITY MECHANISMS. Network security, once synonymous with the firewall, continues its growth in both scope and depth. Government and organizations are requiring security to be baked into the systems and solutions they produce or purchase.

It is entirely possible for an attack to originate from a system not coinciding with a direct perimeter attack. These attack vectors include likely exposures, such as an employee mobile device infected with malicious code while traveling; the new USB flash drive an employee receives at a conference; or a vendor providing legitimate support that plugs his/her laptop into the internal network.

Those who are selecting network security systems must be able to decipher the hyperbole from reality to ensure the device(s) deployed realistically address(es) the specific risks pertinent to the organization. Gateway security devices maintain their basic allowed path control heritage but now are far more advanced in their ability to integrate protection within the payload and beyond.

26.3.1 Allowed Paths. Network-security devices provide allowed-path protections to ensure network traffic only flows in expected and intended ways. Unfortunately, this base level of enforcement is insufficient because application-layer attacks (e.g., client-side browser exploits) leverage these traditional allowed paths. Regulatory and compliance security requirements are also driving the deployment of more

26 · 10 GATEWAY SECURITY DEVICES

capable allowed-path enforcement mechanisms deeper into the internal network to protect intellectual property, customer/partner data, and industrial control systems (ICSs).

To manage allowed-path risk effectively, organizations must define and baseline their expected traffic flows both internally and externally. The baseline is vitally important when the only indication of compromise may be an unusual—though apparently benign—interaction between systems that do not typically interact with one another.

26.3.2 Tunneling. Tunneled access into the internal network for the purposes of remote access, information sharing, and commerce is essential for most organizations. These tunnels, typically in the form of a virtual private network (VPN), use encryption to protect the confidentiality of the transferred data. If the GSD does not terminate the tunnel, it typically only provides network and transport layer allowed path enforcement through to the termination endpoint. This scenario fragments the inspection of the encapsulated traffic, requiring the termination endpoint itself or another security device post-decryption to provide additional protection through allowed path enforcement and application layer inspection.

See Chapter 32 for more information on GSD VPN capabilities

26.3.3 Anti-Spoofing. One of the most readily exploited vulnerabilities of the IP protocol is *spoofing*. Due to this inherent weakness in the IPv4 protocol, a host’s network stack is capable of producing packets that can appear to originate from any IP address and/or port chosen. This attack is an effort to deceive the receiving host and/or any network protection in between into believing the originating host is following an allowed path. Due to the simplicity of executing spoofing attacks, GSD vendors developed and implemented several effective methods to thwart these types of attacks. The network topology is a key element in the GSD’s ability to determine what traffic should be originated from each of its physical and/or logical interfaces.

One common anti-spoofing method is explicit definition of the network topology on a per-interface basis. This provides an efficient method for the GSD to determine whether the packet’s source IP address entering an interface matches what the GSD expects to be the source network. If not, the packet fails to meet the criteria and is dropped. Although simple and effective, the disadvantage is the manual management of the per-interface network definitions.

Another method is to learn the network topology by leveraging the routing table to make the same type of per-interface anti-spoofing decisions. The firewall goes through the same process of evaluating if the packet’s source IP matches what should be originating from that interface. However, this does not require ongoing manual intervention when network changes occur, as the routing engine updates this information. Although this method creates less administrative overhead, two disadvantages include problems if the network uses asymmetric routing and/or falls prey to a route poisoning attack.

Both of these anti-spoofing measures are highly effective when evaluating whether network traffic is originating from the appropriate interface. However, the focus remains solely on preventing internal network address compromise. Since not all of the IP networks in existence are actually allocated and/or in use, these “bogon” networks are another method used to launch spoofing attacks. Since “bogon” network traffic is originating on the external interface, the GSD would not block this based on common anti-spoofing rules. Updated lists are readily available allowing the security administrator to enforce “bogon” protections. Some GSDs have a feature that automatically updates the list on a recurring basis, providing automated protection.

NETWORK-SECURITY MECHANISMS 26 · 11

26.3.4 Network Address Translation. Network address translation (NAT) is a mechanism that maps internal (typically private RFC 1918) network addresses to one or more publicly routable IP addresses. This is an essential function in IPv4 due to the increasing shortage of publically routable address. Due to its ability to mask the internal network addressing, NAT quickly became thought of as a security mechanism. Realistically, NAT only ensures network traffic maintains their translations—but it has no inherent security capabilities such as payload inspection.

To ensure sustained connections in high availability architectures, NAT must establish and maintain a synchronized copy of the translation table with each high availability member. This mechanism follows the same serial or network connection used to maintain the state table.

26.3.5 Intrusion Detection. GSDs inherently provide detective controls at multiple layers. These distinct controls (e.g., firewall, intrusion detection/prevention system, etc.) provide greater visibility into the types and scopes of potential incidents when properly configured and monitored. The logging and alerting capabilities associated with each control are important and necessary mechanisms for security personnel to use to understand the current state of the perimeter. With this information, it is possible to develop actionable intelligence to detect and act on anomalous and/or malicious behavior before, during, and after an incident.

Responding to alerts and reviewing all of the information generated by the GSD can be extraordinarily time consuming. However, it is important that the security administrator take the necessary steps to ensure timely and accurate log reviews. Without such reviews, it is easy to overlook issues such as pre- and post-attack-related traffic and system misconfigurations. Logs also provide an excellent opportunity for GSD tuning, including rule-base optimization and reducing false positives (postverification).

Properly planned and configured alerting and logging mechanisms are essential. When done poorly, these mechanisms can overload the support staff with unnecessary or extraneous information, making them less efficient and effective at addressing actual security issues.

26.3.6 Intrusion Prevention/Response. Although logs and alerts are useful for detecting and investigating security incidents, these mechanisms are passive and cannot provide threat protection. Firewalls evolved to include automated active mechanisms to respond to security events and provide active protection.

See Chapter 27 of this *Handbook* for more details on intrusion detection and prevention.

26.3.6.1 Connection Termination. The essence of perimeter defense is to protect the internal network by blocking connection attempts that do not follow the established allowed path policy. One method for TCP-based protocols is to respond to the initiator (and sometimes the intended recipient) with a reset (RST) packet, severing the connection. The most common method—which works for connection-oriented and connectionless protocols—is to drop the packet without a response. As firewalls evolved into richer functioning GSDs, this protection grew to include blocking based on malicious activity occurring at the application layer.

Inversely, attackers can use connection termination as a method to determine where and what type of network-security devices are in use. This context may provide them additional insight on potential methods to bypass and/or adversely impact the security

26 · 12 GATEWAY SECURITY DEVICES

platform. Though this situation is unavoidable no matter what termination mechanism is employed, it is another reason to apply defenses in layers.

26.3.6.2 Adaptive Threat Mitigation. Although the GSD's policy provides the main enforcement mechanism, vendors quickly realized that this was insufficient, as attackers continually develop additional methods to thwart these protections. GSDs include several enhanced capabilities that allow the device to adapt to specific threat conditions and provide targeted prevention.

Since attackers employ methods to change attack sources to avoid detection, and it is not always practical to follow a rigid whitelisting (least privilege) strategy, organizations typically opt to use blacklists as an additional level of protection. Some GSDs use dynamic object groups to update these lists on a predetermined schedule helping minimize administrative overhead.

The use of thresholds provides the GSD a more efficient method to address consistent attacks that would otherwise follow the same policy-based evaluation for every attempt. One common threshold type is event quantity/type over a given set of time (e.g., total number of packets per second or x number of packets per second across multiple destination IP addresses). Once breached, the GSD will typically automatically block all traffic from the origin IP address for a predetermined time period. This type of adaptive technique requires additional processing resources and requires tuning to ensure this does not impact legitimate network flows.

Being primarily perimeter devices, GSDs typically have minimal visibility into the internal, trusted network. However, modern security infrastructures can leverage multiple threat detection methods across the network to provide more comprehensive protection. Once detected, the device can leverage the security management infrastructure to update other security device policies—automatically providing additional protection or containment. This capability is also integrating with host-based protections to drive protection deeper into the network. If a host detects a threat while on an untrusted network, the infrastructure would learn about this and adapt network protections before the device connects to an internal network.

26.3.6.3 System-Level Actions. As an added layer of defense, gateway security devices typically have built-in mechanisms to detect and respond to threats against the platform itself. Following an implicit deny model, the expected behavior is to fail secure (e.g., shutdown the firewall) and not pass traffic. However, as this automatic response affects the functionality of the security device, it is important for the administrator to understand and test this functionality prior to deployment to ensure this response behavior is consistent with the protection and availability needs of the organization.

26.3.6.4 Application Inspection. The ability to protect allowed paths based on IP and/or port combinations provides marginal benefit when matched against the growing protocol complexity and configurability. As GSDs evolved, so did the ability to evaluate and enforce allowed paths based on the application layer. This capability ranges from basic protocol evaluations such as an RFC conformance check to full application layer gateways. Depending on the type of application layer protocol, a dedicated device may be necessary in high capacity, security, and/or reliability environments.

WAFs. Web application firewalls (WAFs) provide targeted enforcement within an HTTP(S) allowed path. WAFs can more readily detect and prevent common Web server attacks such as XSS and SQL injection and typically have added capabilities to protect

NETWORK-SECURITY MECHANISMS 26 · 13

more complex Web 2.0 (e.g., AJAX) and Web services (e.g., SOAP or JSON) protocols. It is also possible to protect otherwise vulnerable Web servers from exploits by writing protection rules that detect and nullify the application layer attack traffic. The WAFs granular rule development is also better suited to protect custom Web applications.

Proxy Servers. The forwarding proxy enforces outgoing Internet connections by intercepting and controlling access based on criteria such as user ID, time of day, and/or whitelist/blacklist. Depending on the type and configuration, the forwarding proxy can either broker the outbound connection and make requests on behalf of the client or just enforce the particular allowed path. Classic forwarding proxies require client-side configuration to enforced Internet-bound communications. A transparent proxy is a forwarding proxy that operates as a “bump-on-the-wire”—sitting in-line or off a span port—and typically does not require client-side configuration.

Conversely, a reverse proxy provides protection to inbound allowed paths. The inbound connection terminates at the proxy, and the proxy itself establishes the second half of the connection directly with the application server. Being a brokered connection, this adds an additional layer of defense for inbound traffic by rewriting potentially malformed requests and/or not allowing direct communication to the application server.

Application Identity. In addition to source and destination IP address, network firewalls typically rely on the destination port to make a final allowed path determination for common protocols (e.g., HTTP—TCP port 80). However, it is trivial to evade upper layer inspection (e.g., proxy or WAF) by changing the destination port to something that bypasses the application layer inspection but still follows an allowed path. Current generation security devices must be able to profile all network traffic based on application type instead of destination port.

Application identification must be able to dissect the traffic and understand the subsystems/protocols embedded within the main communication flow. For example, Skype not only inherently provides video conferencing but also file transfer and chat services. With add-ons, functionality can extend to remote desktop sharing/control. By having this detailed level of application understanding, it is possible for granular enforcement of one or more protocol subsets. This profiling also includes determining common protocols on unexpected ports and encrypted traffic over a common clear-text protocol—whether encrypted or not.

26.3.7 Encryption. Confidentiality protection for sensitive data-in-motion continues to grow due to the need to minimize the risk of compromising personal/personnel data and/or intellectual property. This creates a distinct need to provide allowed path enforcement when using encryption services to or through the GSD. Due to the computationally complexity, encryption protocols can significantly reduce GSD throughput without offloading encryption services to an add-on, on-board device, or dedicated external devices to meet the encryption needs.

See Chapter 7 of this *Handbook* for more details on encryption.

Inspection. The integration of encryption into network and application flows is essential to protect the confidentiality of the data-in-motion. Although not necessarily malicious, these encrypted sessions reduce the effectiveness of any network-security device(s) attempting to inspect the protected payloads. In addition to the typical VPN and HTTPS Web traffic, the GSD must also be able to identify encrypted communications regardless of port or service. The two main methods for inspecting encrypted traffic are direct termination and on-the-fly decryption.

In the first method, the GSD terminates the encrypted connection. This provides an opportunity to inspect the once tunneled communication. Once inspected and if

26 · 14 GATEWAY SECURITY DEVICES

permitted, the GSD may pass the remained of the communication decrypted or re-encrypt the session to maintain the confidentiality protection to the endpoint.

In the second method, the GSD will use escrowed encryption keys to decrypt network flows passively. Besides meeting the additional performance requirements to conduct passive decryption, another possible weakness is the inability to provide fully synchronous responses to detected threat—potentially allowing some malicious traffic across an allowed path before the enforcement occurs.

VPN. A virtual private network (VPN) is an encapsulated network overlaying an existing set of physical and logical networks. A common VPN implementation is for a remote client connection—allowing an authorized user to access the internal network through an encrypted tunnel. Once established, VPN clients can access internal networks and/or systems while protecting the confidentiality of the connection. Site-to-site VPNs allow authorized remote locations to gain direct internal network access for multiple users over a single connection.

Since there is little to no control when mobile systems leave the confines of the internal network, the fidelity of the mobile device is a concern. The risk of compromise at every level—from physical to operating system to applications and beyond—can expose the internal network to increased risk during a VPN session. Without terminating VPN sessions directly on the GSD, it is still possible to increase visibility and additional allowed path enforcement once the unencrypted traffic leaves the VPN termination device—see the Deployment section for additional details and considerations.

See Chapter 32 of this *Handbook* for more details on VPNs.

Acceleration. Using Secure Sockets Layer (SSL) and Transport Layer Security (TLS) to protect Web-based traffic or IPsec to protect VPN traffic creates a potential barrier to comprehensive allowed path enforcement. Since the encrypted traffic potentially fragments allowed path enforcements, GSDs evolved to be capable of terminating these connections to add an additional layer of inspection. However, encryption/decryption is mathematically intensive and requires sufficient processing resources to ensure these processes do not impact GSD throughput. When it is necessary for the GSD to terminate or inspect encrypted traffic, the use of a hardware encryption module is a common method to off-load this processing burden and minimize impact to other traffic flowing through the GSD.

26.3.8 Identity-Based Enforcement. By integrating with enterprise directory services using mechanisms such as Lightweight Directory Access Protocol (LDAP) or Remote Access Dial-In User Service (RADIUS), GSDs are able to make identity-based decisions. Instead of allowing or denying access to an allowed path by source IP, the GSD can authorize access per user or based on group membership and include this authorization data in logs. This additional detail enhances reporting and auditing capabilities, providing more specific information about a particular connection well beyond just the IP address or DNS name of the host requesting allowed path traversal.

26.3.9 Complex Protocol Engines. Complex protocols, including those used to deliver feature rich Web applications, IP telephony, or specialized protocols, such as those used by industrial control systems (ICSs), are an intrinsic necessity for many organizations. Due to the ever-increasing need to share information and access regardless of mobile device form factor or location, these complex protocols are traversing the GSD. Even with successful detection of an application in-use, it is a continual

NETWORK-SECURITY MECHANISMS 26 · 15

development race for the GSD to be able to understand and enforce the intricacies of each protocol.

Today's Web-based applications continually develop new functionality to meet customer demand. Social media Websites offer viable business models and may necessitate some or all users within an organization to access this service to conduct business. Due to ever-present personal elements (e.g., chat/comments, games, etc.), the GSDs must be able to dissect the connection and selectively permit some or all of the functionality offered by the host application.

IP telephony protocols such as session initiation protocol (SIP) and real-time protocol (RTP) are typically intolerant to latency and jitter—which are introducible by GSDs. ICS protocols introduce vendor-specific protocol dependencies/intricacies and are even less tolerant to network variations. The loss or slowness of ICS communications can be as severe as damage to equipment or death when impacting safety systems.

Although an internal or external network perimeter is a logical point of inspection, the GSD is not necessarily the appropriate device to conduct such detailed protocol inspections. Instead, these situations may warrant additional security layers using specialized devices for the application layer or specialty protocol enforcement.

26.3.10 Content Control and Data Leakage Prevention. Organizational policies establish the behavioral expectations of its employees. The Acceptable Use Policy (or similarly named) typically covers matters specific to use of technology resources; however, the policy itself cannot provide active protection. Content filtering (CF) bridges the gap between technical and nontechnical policy enforcement. This enforcement typically focuses on internal users attempting to access information resources outside of the organization. By filtering heavily used protocols such as HTTP and SMTP, the organization can minimize user access to information not fitting the organization's definition of business use.

This technology has several methods for inspecting, classifying, and filtering content. Inspection takes multiple forms, including residing in-line, out-of-band to monitor passing traffic without impacting the physical connection, or in conjunction with proxy redirection. Classification occurs through multiple methods, including vendor-defined categories, blacklists, and/or reputation scoring based on IP address and/or DNS name. This also allows the organization to choose to use a subset of the provide policy and/or selectively develop a whitelist to provide more granular enforcement. The filtering feature logs and/or severs the connections allowing the organization to customize its response to detected policy violations, including redirection to a “not for business use” Web page and/or automatic administrative and/or managerial notification of repeated violations.

See Chapter 31 of this *Handbook* for more details on content filtering.

26.3.10.1 Information Classification and Enforcement. An overlooked but important aspect of information assurance is the proper classification of information/data into categories. This allows the organization to determine the risk reduction efforts necessary to protect each information level/type. Similar to the classification process used by governments, identifying the sensitive information increases the ability to enforce not only who has access to the information but also where and how this information can be transferred.

Digital rights management (DRM) and data loss prevention (DLP) are two current generation technologies for addressing these information-specific risks. DRM

26 · 16 GATEWAY SECURITY DEVICES

focuses on the definition and enforcement of data-at-rest—typically at the information repository level. GSDs may include DLP protections such as enforcing data-in-motion protections, including denying the transfer of files based on content/designation and/or requiring encryption for sensitive data.

26.3.10.2 Anti-malware. A contemporary attack vector is to target a commonly used Website/service and inject malware. Although the content filter may make an initial decision to allow the connection based on the specific site name, it may also have the added capability to intercept and block the malicious content before it ever reaches the host. This functionality makes it possible to evaluate different types of traffic such as SMTP attachments and HTTP downloads. Although this type of application layer protection can dramatically reduce the GSD throughput, it is possible to redirect the scanning to an additional module within the GSD with dedicated processing power or to a dedicated network anti-malware appliance.

26.3.10.3 Active Content. Active code used on HTTP connections—such as Silverlight, Java, AJAX, and Flash—provides a foundation for a rich Web experience but also increases security risk due to their application-focused functionality. Malicious email attachments—such as documents with embedded code—are common components of phishing and spam campaigns.

In conjunction with host-based protections, the GSD can decrease the risk of active code by preventing specific types from ever reaching the endpoint and/or scanning the active code in a sandbox prior to allowing the active code to traverse the allowed path.

See Chapter 17 of this *Handbook* for more details on mobile code.

26.3.10.4 Caching. Although the aggregate cost of bandwidth continues to fall, organizations continue to look for opportunities to manage these frequently congested circuits more efficiently. Proxy servers typically integrate caching mechanisms to store frequently used Internet content locally, increasing local throughput by reducing Internet-bound traffic. The GSD may also be able to provide similar caching services, given sufficient storage capacity and processing resources.

26.3.11 IPv6. The allocation of the last remaining available IPv4 address blocks to Regional Internet Registries (RIR) occurred in 2011—indicating the world is “out of IPv4 addresses.” Although one may argue that statement is overly dramatic, mass migration to IPv6 is still occurring at a slower pace than one may expect, even though IPv6 was developed in the late 1990s. While organizations continue to develop their individual transition plans, it is likely that IPv6 is already in their environment. Many modern operating systems and devices have dual IPv4 and IPv6 network stacks, with some enabled by default. IPv6 and its associated transition technologies have specific security implications that the GSD must understand and enforce.

26.3.11.1 Perimeter Security Concerns.

Addressing. IPv6 is more flexible in its approach to dynamic addressing. Instead of solely relying on DHCP, an IPv6 device can address itself through stateless address autoconfiguration (SLAAC). The host uses a unique identifier (typically its own Message Authentication Code (MAC) address) in addition to the Neighbor Discovery (ND) protocol to complete the automatic addressing. Since there is no authentication

NETWORK-SECURITY MECHANISMS 26 · 17

requirement, the GSD must prevent external devices from attempting to act as an internal router during the addressing process.

The significant increase of available addresses in any particular IPv6 network makes it infeasible to discover devices and network topology using traditional port scanning methodologies. By using the multicast listener discovery (MLD) protocol, an attacker can send a probe to the link-local multicast address (ff02::1) and listen for responses. The GSD must block this capability at the perimeter to prevent external devices from attempting to discover internal hosts and topologies.

Tunneling. Without ubiquitous end-to-end IPv6 connectivity, there are several IPv6 transition technologies (such as 6to4 and Teredo) that allow IPv6 capable systems to tunnel communication over legacy IPv4 networks. As with other tunneled traffic, to be effective, the GSD must not only be able to enforce the appropriate allowed path for the tunneled traffic but also inspect the encapsulated IPv6 packet.

IPv6 also has native support for IPsec (both AH and ESP). Configuration and use of IPsec in an IPv6 environment requires the same discipline in choosing and configuring cryptographic options as IPv4. It is also possible to use IPv4 IPsec to protect IPv6 transition technology tunnels as unencrypted tunnel sessions would otherwise be vulnerable to interception and/or manipulation.

Global Connectivity. Network address translation (NAT) is an essential element for IPv4 Internet connectivity due to the relatively “small” number of publically routable addresses. The size of the IPv6 address space (2^{128} addresses) provides for global IP interconnectivity without the need (or definition) of a NAT replacement. IPv6 devices will now communicate natively across the Internet—changing perimeter dynamics yet again. With IPv6 enabled, it is essential for GSD to enforce strict ingress *and* egress filtering.

Mobility. Mobile IPv6 (MIPv6) allows hosts to maintain access to the home network while physically roaming to other locations and uncontrolled networks. A user would be able to leave the office network and travel from one appointment to another without losing connectivity to the internal network. Since the device can move network to network without dropping the internal connection, MIPv6 also creates potential issues with stateful inspection if the GSD does not understand the protocol. If the organization chooses not to support MIPv6, the GSD should filter Type 2 Routing Header (RH2) packets.

26.3.12 Additional Considerations

26.3.12.1 Host Protection. Although firewalls and GSDs play a crucial role in protecting the overall network infrastructure, it is neither cost feasible nor possible to deploy these devices at all points inside the network. Individual hosts must have a way to protect themselves from threats independent of network-security devices. Unlike the network, a host has a contextual understanding of what the system can, is, and/or should be doing. Beyond the simple allow and deny functionality, contextual security measures can detect unexpected system configuration changes such as service changes or an application behaving in an unexpected manner.

When a host leaves the internal network, it becomes an extension of the network protection profile. By providing adequate local protections at the network and application levels, the mobile endpoint helps to mitigate potential issues upon reconnecting to the internal network. It is also crucial to determine the level of protection necessary and how each of these additional levels of security will affect system performance, management, and end-user impact.

26 · 18 GATEWAY SECURITY DEVICES

26.3.12.2 Network. Hosts need to be able to determine the types and appropriateness of inbound and outbound traffic. These network restrictions may vary from the internal network to uncontrolled networks. In certain circumstances, all network traffic may be suspect and scrutinized further. For example, when on the internal network, the host allows most inbound and outbound traffic. If the host were on an uncontrolled network, it would not allow any nonestablished network flows. Simple host-based network protections are not enough; additional host protection mechanisms such as intrusion prevention can detect and stop network and application-based attacks.

26.3.12.3 Applications Access. The host's contextual awareness also helps dictate the ability for applications to execute as well as send and/or receive data. The goal is to ensure only the appropriate applications and/or services have network access. The host protection policy may allow applications to establish outbound connections, but never listen (nor accept nonestablished inbound packets). For example, an HTTP server uses a daemon to listen for connection attempts. The HTTP client will attempt to make a connection to the HTTP server daemon. By using a host protection mechanism, it would be possible to prevent one or both of these actions.

26.3.12.4 Hybrid Protections. The host intrusion prevention system (HIPS) functions similarly to a network intrusion prevention system (NIPS) by detecting known attack patterns and/or anomalous behaviors. Hybrid host protections build on the host's contextual awareness and provide the ability to monitor other unusual application level activity such as changes to binaries, service manipulation, and spawned listeners.

26.4 DEPLOYMENT

26.4.1 Zoned Architecture. In a contemporary interpretation of the screened subnet architecture, zones define different types and/or sensitivities of networks, applications, and/or services. This provides the ability to manage allowed paths at the macro level (per zone) in addition to the traditional allowed paths (specific hosts or services). As requirements drive security deeper into the network, the zoning concept is equally effective when used to manage and protect internal and external networks.

26.4.1.1 Perimeter Zones. The border router maintains the architecture's first line of defense against external attacks. The ACL(s) on this router should mirror the basic allowed-path configuration of the external (untrusted) firewall interface and provides several important benefits.

The GSD is able to operate at optimal efficiency, since traffic rejected based on border router's packet-filtering rules normally would never reach the firewall. This permits the firewall to focus, in terms of load, on protocol inspection. If, for example, the firewall receives a packet that should never have made it past the border router's ACL, the firewall can assume that the router is not behaving normally. The firewall is then free to respond appropriately, with such actions as terminating all connections from a specific host.

26.4.1.2 External Service Zones. The necessity for mobility and accessibility places significant demands on Internet-facing systems in addition to increasing the administrative overhead of managing external access. The sections below provide only a sampling of possible external service zone architectures.

DEPLOYMENT 26 · 19

Utility. Instead of lumping systems such as Web, DNS, and email onto a single network, there may be an advantage by implementing zones. Utility servers such as DNS and email could logically be on the same network. Web servers typically demand greater bandwidth, and by using this concept, can protect the entire zone with a simple inbound access rule.

Extranet. Extranet systems create additional complexity since they provide the user interface, while internal systems may provide the relevant content. Zoning provides more flexibility by allowing external connections to reach the Extranet servers while providing those same servers access to internal resources.

VPN. VPN networks are also an opportunity to use zoning. Since the VPN connection device must be Internet facing, this requires two different networks connected to the firewall. The external, Internet-facing (VPN untrusted) interface would only allow a few protocols for inbound and outbound encrypted traffic. The second network (VPN trusted) is for unencrypted network traffic moving to and from the internal network. This architecture also provides extra internal network protection in the event of a compromise of the VPN device as well as creating a traffic inspection point that is unimpeded by encryption.

26.4.1.3 Internal Service Zones. With the increasing prevalence of non-company-owned assets and more complex data flows, organizations continue to look for better methods to protect internal networks. The sections below provide only a sampling of possible internal service zone architectures.

Administrative and Monitoring Systems. Organizations typically limit direct access to network and security devices. By requiring all administrative and monitoring traffic to originate from a specific host and/or network, a zone can effectively reduce each device's threat profile. This zone would not only allow minimal noninitiated inbound traffic but would also limit outbound connections to the managed and monitored systems.

High-Value Systems. Organizations relying heavily on intellectual property and/or other protect data types (e.g., personal and financial information) have an intrinsic need to provide higher level protections to protect their investments. Zoning provides an additional layer of protection by minimizing unnecessary information flows to/from these high-value systems.

Industrial Control Systems. Incidents such as Stuxnet and Night Dragon are stark reminders that even though certain system types are considered extremely complex and less accessible, it is only a matter of time before successful compromise is possible. Manufacturing organizations with significant investments in industrial control and supervisory control and data acquisition (SCADA) systems can (or have requirements to) isolate these environments as an additional layer of protection and/or compliance.

26.4.2 GSD Positioning. The increased use of encrypted protocols such as SSL/TLS and IPsec can blind network protections. Certain GSDs have the ability to terminate encrypted sessions, though the increased processing and bandwidth requirements may exceed the limits of the device. If concerned, the security architecture should deploy appropriate countermeasures at strategic locations that avoid encrypted traffic. This way, the GSD can focus on its primary role of detecting and preventing malicious activity.

26.4.2.1 Inline. Placing the GSD inline creates a choke point for active enforcement on all network traffic that flows through it. When a malicious packet enters

26 · 20 GATEWAY SECURITY DEVICES

the GSD, protocol analysis will detect the anomaly and will not allow it to flow out the other interface. Although bandwidth limitations are a typical concern, improperly configured inline devices may also present a denial-of-service condition. With proper infrastructure planning and deployment, it is possible to minimize these risks.

26.4.2.2 Controlling Encrypted Traffic. Since mobile devices frequently venture outside of the controlled network, one logical place to evaluate traffic is on the unencrypted side of the connection. This may be on the backside of a SSL terminator (in some cases on the server itself) or on the unencrypted side of a VPN connection. The second option is to use the GSD to inspect (e.g., termination and/or passive decryption) the traffic before forwarding the traffic onto its final destination. This level of functionality requires substantial processing resources but may be a necessity if the security layers downstream are insufficient.

26.5 MANAGEMENT AND MONITORING STRATEGIES. Regardless of vendor claims, network-security devices are never a plug-and-play endeavor. It is essential to take additional steps to define the security requirements for managing and monitoring GSD components. This approach helps ensure a well-rounded security posture.

26.5.1 Monitoring. Firewalls and GSDs provide complex functionality; monitoring such systems must go beyond just verifying system availability and cover device health, availability, and integrity.

26.5.1.1 Health. Metrics such as processor utilization, available RAM, and number of connections all have an impact on overall functionality. A centralized management console may provide the ability to monitor and alert on these metrics. If this functionality is unavailable, it may be necessary to use monitoring protocols such as Simple Network Management Protocol (SNMP) and/or Remote Monitoring (RMON) to gather these statistics. The GSD must tightly restrict the systems able to poll using these methods because of the inherent insecurities of the aforementioned monitoring protocols. By trending these metrics, it may be possible to determine when it is time to increase bandwidth or purchase systems that are capable of meeting the new throughput or processing needs.

26.5.1.2 Availability. When GSDs are unavailable, the network functionality can dramatically diminish. A simple test of system availability is using ICMP to “ping” one or more interfaces to ensure the device itself is responding. However, this approach can be deceptive. Just because the device itself responds, does not mean it is properly forwarding traffic. It is also advisable to send probes (e.g., ICMP, traceroute, or other queries) to something on the other side of each interface to ensure the other device is actually receiving the packets to ensure valid results. This approach provides a better overall picture of the GSD availability.

26.5.1.3 Integrity. The ability to trust network security systems components is vital. The possibility of a root kit compromising a firewall or GSD is now a reality. These systems must have the ability to protect against modification of system components such as ceasing operation and/or alerting the change. If this embedded functionality were unavailable, it is possible to write a script to generate cryptographic hashes of critical system components and verify against a known trusted version.

MANAGEMENT AND MONITORING STRATEGIES 26 · 21

26.5.2 Policy. The GSD policy is the core definition for providing and protecting allowed paths. Most security systems process packets starting at the beginning of the policy and continue until there is either a match or reaching the end of the rule base (which should be an explicit “deny any”). As discussed later in this chapter, there are situations where rules may process before or after the main rule base.

Centralized management consoles provide intuitive GUIs to configure and easily manage one or more firewall and GSD policies. Certain platforms also provide the ability to manage policies directly from the device.

26.5.2.1 Defining Allowed Paths. Allowed paths identify specific protocols used to implement communication. In a typical Internet environment, business services require allowed paths such as HTTP(S), SMTP, and DNS. These requirements will vary, but for an environment, each allowed path should directly relate to the required service.

Starting from an implicit or explicit (depending on the platform) “deny any” rule, allowed paths will be added as “allow” rules, such as PERMIT HTTP, with specifics determined by the following sections.

Although network addressing does not provide effective authentication of systems or users, restrictive endpoints can make it much more difficult for an attacker to exploit an otherwise straightforward vulnerability. It is also important to identify the endpoints carefully, particularly in cases where these endpoints might reside on internal rather than extranet or utility zones.

The direction of traffic, indicated by the source of the connection initiation, is useful for the rule definitions for several reasons. First, rules can be written so that only responses to internally originated allowed paths are allowed in from the untrusted zone, rather than explicitly permitting the protocol bidirectionally. In addition, the firewall may process rules at different times based on design and/or configuration.

26.5.2.2 Complexity of GSD Policies. Standard firewall rules operate on simple Boolean principles. For example, allow or deny network traffic that is going from host or network X to Y on port Z. The complexities required of GSDs evaluating network traffic are dramatically higher. For example, this evaluation could be a combination of a Boolean test to verify an inbound email address is from a trusted source, verify message contents are acceptable, and scan an attachment for viruses. Administrators must understand the higher-level protocols to ensure that the GSD policy matches the types of protections expected and required. As the number and complexity of the rules increases, so do the processing requirements for the GSD.

Beyond the basic firewall capabilities, GSD policies typically include per-rule enforcement of the additional security measures. For example, the GSD only uses network and/or transport layer filtering to restrict access to the organization’s VPN terminators. The policy may include WAF protections, but only for the utility zones. Inbound and outbound Web and file transfer traffic may have the additional requirement for NIPS inspection. Although this level of customization adds some administrative overhead, selectively (rather than broadly) applying additional protections in this manner can reduce the performance impact of the added protections.

26.5.2.3 Change Management. Whether managing one or one hundred policies, an essential element is to have a process to track policy changes. Change management can be cumbersome but has several advantages. First, this provides back-out information if a change were to cause issues. Second, it provides an audit trail of

26 · 22 GATEWAY SECURITY DEVICES

who requested the change (and why) in addition to who changed the policy. Lastly, it provides a method for streamlining change requests by providing a single method for accepting, reviewing, implementing, and validating policy changes.

26.5.2.4 Secondary Evaluation. Making policy changes tends to be a simple process, such as adding HTTPS access for a new extranet server. However, the more complicated the change, the more likely for mistakes to occur. Having another administrator evaluate the proposed change brings a fresh perspective that may catch a small discrepancy that could have negative consequences. Another step in the change management process is validation across multiple operational groups—providing an opportunity to detect potential conflicts with other parts of the infrastructure due to the change.

26.5.2.5 Auditing/Testing. Is your firewall or GSD working as expected? Is the GSD catching the most recent malware-infected attachment? Would you bet your organization's future on it? These questions are simple yet powerful reminders that you must verify that the network-security devices are working as expected. In addition to on-going and in-depth log reviews, a proper regimen of auditing and testing will help to answer these questions. In addition, management consoles should provide an audit log that tracks who makes changes to any part of the environment.

Vulnerability assessment and penetration testing are effective methods to determine the validity of policy rules. For example, the firewall should not allow an HTTP connection to a system with only SSL/TLS access configured. If the HTTP access succeeds, is the failure on the tested rule, is there another rule higher in the rule base erroneously allowing this access, or is something unexpected occurring? The tools and processes must be in place to answer these questions before someone else does.

26.6 MAINTENANCE. Patching workstations and updating malware definitions/signatures are common practices, but this process is even more crucial for systems protecting the network. It is essential to test each update's validity and stability to ensure faulty or rogue changes do not interrupt operations.

26.6.1 Patching. No system is inherently and infinitely secure, and it may be possible to subvert a security device due to system vulnerability. It is crucial to monitor GSD Websites for the most recent operating system and system component patches. By also monitoring vulnerability and exploitation information sources, it may be possible to ascertain more timely information. This additional research could result in being able to determine and implement a temporary solution until the vendor can provide a permanent fix.

Vendors can take weeks and months to develop and release a patch. In certain cases, third-party patches will become available for un-patched, actively exploited vulnerabilities. Although using these is an option, it is inadvisable due to the inability to verify the patch's integrity and safety.

26.6.2 Pattern Updates. Threats evolve at will. Automatic updates provide the smallest delta of exposure for pattern-based malware signatures. However, blindly trusting these updates in a production environment may have adverse effects. To avoid issues, establish a procedure that enables new signatures in monitor only mode to ensure there are no adverse effects before implementing full protection. If automatic updates

MAINTENANCE 26 · 23

are not a viable option, then testing requires a lab environment or use of noncritical systems to vet the new patterns/signatures.

26.6.3 Logging and Alerting. Whether remaining local or transferring to a centralized management system, logs take up a large amount of disk space, and if left unviewed, are worthless to keep. There must be a log review process. Instead of trying to determine the anomalies from an entire log set, it is helpful remove known traffic to help expose the unknown and potentially malicious. For example, if you only expect outbound SMTP traffic from a single network, you would filter out those known items and more easily see if the other unauthorized hosts or networks are originating SMTP traffic.

Organizations with a large Web-based footprint will generate a large number of HTTP/HTTPS-related events because each new connection creates a log entry. By having a GSD do protocol inspection and logging, there is the potential of gaining a better understanding whether packets are expected or malicious. If there are other systems tracking and/or correlating these connections, it may be useful to turn off logging on selected GSD rules to reduce the overabundance of logs that are unnecessary to review.

A security incident and event management (SIEM) system provides an additional method for collection, aggregation, and consolidation of logs from many types of devices. The SIEM leverages baselining and configurable rules to correlate the logs and provide real-time incident-based alerting.

Alerting is complementary to and usually depends on the logging mechanism. Once the firewall or GSD generates a log entry, the administrator can configure alerting options when certain log conditions or threshold changes exist. For example, the administrator can configure an alert to send a notification if an internal system process changes state. By carefully determining alerting thresholds and notification methods, there is less opportunity to burden the administrator(s) with nuisance alerts.

26.6.4 Secure Baseline Configuration. Another crucial step to protecting network-security devices is to create a secure baseline configuration. In certain cases, the GSD vendor provides a hardened and/or proprietary version of an operating system. This should not be a “green light” that the device is secure and ready for production. Rather, this provides even more reason to take the time to evaluate the security posture of the device thoroughly.

Once a secure configuration is set, this then becomes the baseline configuration for the remainder of the systems. This provides a standard configuration helping to ensure each device functions and secures in the same manner. There should also be a process to ensure the integrity of the secure configuration as well as verifying and testing when valid updates occur.

26.6.5 Default Configurations. Default system and policy configurations vary widely; some implicitly deny, while others may implicitly allow. In no case should any network-security device deploy into production with a default configuration. Ensure that you determine whether the network-security device fails in an insecure or secure posture. Inline devices may fail insecure in an attempt to maintain traffic flow in the case of failure. Depending on the network, this may or may not be the optimal response. Although potentially disruptive, a fail insecure posture reduces the likelihood of anomalous or malicious traffic passing undetected.

26 · 24 GATEWAY SECURITY DEVICES

Implied Rules. Implied rules are separate from but are part of the default configuration. The network-security device will process the implied rules prior to packets getting to the policy rules. The implied rules may allow the GSD to process specific and known administrative traffic without using processing cycles to go through the entire rule base. It is essential to determine whether these rules exist as well as whether they match the desired security posture. If not, disable or modify the implied rules to fit specific protection needs.

26.6.6 Harden the Device. Protecting the network-security device is a critical part of the security of the overall infrastructure. The two most common ways to reduce these exposures are through the administrative console and validation using vulnerability assessment. The management or system console can provide information such as default listening ports, services running, and the like. If a service is not critical to the functioning of the firewall, disable it to remove any threat of it becoming a compromise vector. Even with the service disabled, ensure you continue to apply patches for those services. Only specific administrative hosts should have direct access to the system. As with the GSD policies, it is useful to validate that configuration changes are actually making the device more secure. By conducting a vulnerability assessment, you are able to determine whether unexpected services are still available or additional vulnerabilities exist.

26.7 DISASTER RECOVERY. The impact of a GSD outage can range from an annoyance to critical event depending on implementation, necessity, and disaster recovery planning. Since these systems represent part of the backbone of the security infrastructure, it is a necessity to provide continual, protected network access.

26.7.1 Fail-Over/High Availability. High availability (HA) architectures are essential in situations where reliable access is a necessity. The HA deployment will typically have an Active/Standby pair; where the active member is processing live packets and the standby member is ready to take over in the case of a failure of the active member. If the standby member is continually synchronizing connection information from the active member, there would be little to no loss of connectivity if the active member fails because the standby would then become active and already have the connection information to allow continued processing of the current sessions and service new sessions.

26.7.2 Load-Balancing. Distributing the load between multiple systems is another way to reduce an availability exposure. If the load equally distributes between two or more systems, the failure of one does not eliminate all access. One caveat is that the load balancer must continue to route each connection to the same system on the back-end to maintain connection information. Otherwise, it may be possible for the load balancer to route traffic to another system that does not know about the established connection (e.g., asymmetric routing).

26.7.3 Backup/Restore. Backup and restore functionality is the most rudimentary disaster recovery method. If the system crashes or suffers a compromise, there must be a way to recover the system in an efficient manner. By having a reliable backup process, it becomes easier to restore the device configuration in the event of a failure. Depending on the vendor, backups may consist of a text file or something more complex, like a GZIP file containing the crucial configuration information. The backup

NETWORK-SECURITY DEVICE EVALUATION 26 · 25

process must include a method to move the backup to another system to eliminate a single point of failure. The restore process may involve restoring the operating system and/or the security device configuration. This may be as simple as uploading the backup file or may require an out-of-band connection to do prework before uploading a preserved configuration.

26.8 NETWORK-SECURITY DEVICE EVALUATION. As digital threats continue to evolve, so must network-security devices. Internet connectivity costs continue to decline and so does the ability to detect complex threats. It is a mistake to think the current security systems and infrastructure will be viable for an extended length of time. In essence, the infrastructure is continually under review. One crucial time to revisit these items is when it comes time to replace security devices.

Effective security lies at the intersection of protection/functionality, usability, and cost. Every situation has unique challenges and opportunities—the goal is to develop the best security solution to meet the organization’s ongoing risk requirements. The following sections provide a framework for evaluating network-security devices currently in-service or during a Request for Information (RFI). Please note that this section addresses the process over the lifecycle of the GSD architecture instead of providing detailed guidance on security-inspection capabilities. These guidelines provide a good basis for thoughtful analysis of site-specific requirements.

26.8.1 Current Infrastructure Limitations. It may be technically possible to do a one-for-one replacement of an existing firewall with a GSD, but it is important to review the current infrastructure to ensure maximum effectiveness of the new device(s).

- Are the current network-security devices past useful life?
 - This includes devices that are out of warranty, no longer vendor supported, and unable to meet current requirements.
- Is the Internet-facing architecture limiting GSD deployment options?
 - Having only a few publically routable IP addresses can limit the number of external service zones without extensive use of NAT.
 - Placing the GSD in an allowed path of all encrypted traffic may relegate the GSD to being more packet filter than extensible security solution.

26.8.2 New Infrastructure Requirements. The dynamics of network and security infrastructures shift constantly. Trends of rapidly escalating bandwidth usage or an increased reliance on encrypted communications are common. The prospect of adding the GSD capabilities have a dramatic influence on these decisions.

- Are there new demands for increased bandwidth, encryption, or application layer enforcement?
 - Determining new demands include observed or anticipated/projected growth.
- Is Internet-access redundancy becoming a necessity?
 - It is important to consider a GSD with built-in redundancy capabilities, such as Active/Passive failover or clustering.
- Is your infrastructure able to support diverse carriers and/or routes to the Internet?

26 · 26 GATEWAY SECURITY DEVICES

26.8.3 Performance. Technological obsolescence continues to haunt organizations attempting to grow to meet internal and external demands.

- Does the existing device have excess capacity for future needs?
- Security devices running at high processor and/or memory usage are good indicators that the system is at or reaching its maximum performance threshold.
- Without extra processing capacity, a modest increase in network traffic or enabling a new feature could cripple an already resource starved device.

26.8.4 Throughput. The need for bandwidth is relentless. Bandwidth of 10 Gigabit is no longer sufficient in the core since the price and availability of 40 and 100 Gigabit is already increasingly obtainable. Standard 802.11ac is ready to eclipse the established 802.11n wireless standard and offer claimed speeds of a gigabit per second (or greater), leading organizations to question their use of wired connectivity. Organizations should take a realistic approach before paying for higher bandwidth solutions.

- Does the device meet bandwidth requirements?
 - The GSD must be able to sustain the necessary peak traffic loads without interruption or service degradation.
 - Are aggregated network interfaces supported to support modest increased bandwidth rather than paying for higher speed interfaces (e.g., multiple Gig interfaces versus 1 × 10 Gig)?
- Does the device have excess capacity for future needs?
 - Modular devices are able to increase memory, add network interfaces, or increase bandwidth by adding or replacing a hardware module.
 - Peaks in network traffic are inevitable. Having additional bandwidth capacity for existing GSD network segments will reduce the impact of these situations.

26.8.5 Implications of Added Features. Many vendors tout throughput at (or near) wire speeds. However, these statistics do not necessarily mirror reality as the standard throughput tests always consider the varying packet sizes and protocol mixtures. In actuality, packet types and sizes vary constantly; be sure to review vendor claims through independent testing organization and, if possible, internal testing.

- What performance impact does activating additional security features create?
 - The greater the number of additional functions—typically resulting in more detailed payload inspection—the increased negative impact on the device's processing power and bandwidth.
 - If requiring encryption, hardware acceleration, if available, will reduce much of this specific performance impact.
 - Be sure to evaluate the vendor's and independent documentation showing measured performance against varied traffic flows such as IMIX (Internet Mix) and EMIX (Enterprise Mix).

NETWORK-SECURITY DEVICE EVALUATION 26 · 27

- What level of additional administrative overhead will the features create?
 - Consolidating management of multiple security functions does not necessarily make them easier to manage. This also adds complexity depending how each subsystem integrates.

26.8.6 Management. GSDs typically require substantial planning, configuration, monitoring, and maintenance. Having a robust management platform in distributed GSD environments is essential to the success of the implementation.

- Does the current environment have the management features needed or required?
 - Self-evaluation of the current management environment is fairly easy, as there are already established likes, dislikes, nice-to-haves, must-haves, and must-nots.
- Is there the need/want for distributed or centralized management?
 - The size and type of GSD deployment is important to consider. It probably does not make sense to use an extensive management platform for a few devices. With larger geographically distributed environments, having centralized control of all devices is necessary for efficient and effective command and control.
- Are redundant management systems needed?
 - Redundant management systems ensure administrative control and logging of GSD devices in the case of a failure. Be sure to review management systems redundancy options to ensure the secondary system has functional parity.
- Is the ability to manage the system from both a centralized management console and/or directly from the device a necessity?
 - Certain types of GSDs allow system and policy changes to be made directly from the device. This feature can be useful during testing or troubleshooting but may cause confusion or misconfiguration if not properly understood and managed.
- Is encryption a requirement to protect management, policy, and logging functions?
 - The need to encrypt sensitive information is an important consideration for transmission of GSD management and logging functions. If sent clear-text, it is possible for someone to monitor these communications and gain significant intelligence about the network security architecture.
- Is detailed reporting mechanism a necessity?
 - Most logging systems provide the ability to filter data, but it may be useful to do more in-depth reporting. Some example reports include rule base hit percentages (e.g., which rules receive the highest percentage of hits) or traffic usage statistics (e.g., percent usage by protocol, network segment, host, and/or user).
- What level of granularity is necessary for management permissions?
 - With only a few administrators, this may not be as significant, but larger deployments may have distributed administrators. The GSD management system must provide role-based security with the ability to restrict access to certain tasks, areas, and/or logs. It may be useful to compartmentalize the systems and allow other support teams view-only access for troubleshooting and auditing purposes.

26 · 28 GATEWAY SECURITY DEVICES

- Are both GUI and/or command-line policy or device management options required?
- Certain GSD brands have the option of configuring devices using a more traditional command line interface and can be useful for quick troubleshooting and scripted changes. A drag-and-drop graphical environment uses the central management systems to easily create, change, and view policies and logs. If both options are available, ensure both options have feature parity.

26.8.7 Logging and Alerting. Although the GSD's primary responsibility is denying malicious traffic, logging and alerting play a critical role in the overall security posture. Logging is essential to investigating incidents, viewing trends, and may even be a regulatory requirement.

- How do you want to be able to filter logs?
 - Collecting, searching, and reviewing logs in a text-based format has advantages, including consolidation and scripting for log reviews. Text-based log require scripts to parse through logs but do have the advantage of being able to use complex queries and the ability to analyze logs on virtually any system.
 - Graphical user interfaces (GUIs) provide the ability to filter logs based on one or more criteria but are weak when attempting complex or recursive queries. A GUI increases the aesthetic appeal of the logs and can ease viewing by grouping or filtering different types of events, as well as adding color and/or graphics that help catch the eyes' attention more quickly. Although useful, the GUI may be slower to load and may restrict the ability to create complex queries.
- Do you want compressible logs?
 - Even with well-maintained logging standards and grooming processes, the sheer volume of GSD logs alone can be overwhelming. Text-based logs typically provide high compression ratios and easy archival. Some of the proprietary formats may automatically provide some compression but may not be able to match that of nonproprietary formats.
- How do you want to store logs?
 - GSD logging options include storage on the local device, on an integrated logging and management platform, or on a separate storage device. Local storage creates the risk of data loss if the GSD fails or potential for system impact if the logs fill the local storage. Sending the logs to other systems will increase network traffic, but this provides easier log management over the information's lifecycle.
- What types and granularity of alerts do you need?
 - Alerting helps to bring information to life without the need to watch logs every minute of the day. It may be possible to create email or SMS alerts an event occurs breaches a configured threshold (e.g., the GSD network throughput is more than 85 percent of the maximum possible or a host is attempting to attack the GSD).
 - It may be possible to develop detailed alerting conditions (e.g., criteria dependencies) to minimize erroneous alerts.

NETWORK-SECURITY DEVICE EVALUATION 26 · 29

- Do you have a legal and/or regulator requirement to keep logs for a set period of time?
- Internal and/or regulatory requirements may make it necessary to keep this data ranging from months to many years.

26.8.8 Auditability. Audit records provide an accounting of what changes occurred, when, by whom, and so forth. Having an integrated audit trail of GSD environment changes can be invaluable during an incident.

- Do you want to be able to audit device changes?
 - Changes to the device itself do not occur nearly as often as changes to policies. Having the ability to collect information about these changes can help uncover system instability, accidental misconfiguration, and/or malicious intent.
- Do you want to be able to audit policy changes?
 - GSD policies tend to change rapidly as new access requirements develop, new capabilities or requirements appear, and the ever-changing threat landscape progresses. Policy audit trails help to resolve mistakes and provide archival information of what changes occurred and when.
- What level of audit granularity is necessary?
 - Data such as date, time, type of change, to what system(s), and from what administrator are all valuable pieces of information. Increasing auditing granularity can increase the overall log volume and processing overhead.
- Should the audit logs be independent from the firewall event logs?
 - Having a logging system that automatically separates audit logs from other event logs helps ease viewing and searching. This also creates the ability to create more granular permission levels (e.g., who can see which type of logs).

26.8.9 Disaster Recovery. With GSDs providing consolidated security functions, incorporating disaster recovery (DR) options into the GSD environment helps ensure it does not become a single point of failure for the entire organization.

- What are your specific disaster recover requirements?
 - As mentioned earlier, different types of redundancy help during DR scenarios. High availability setups such as active/passive failover and clustering can help to increase overall availability.
- Does the device have easy-to-use backup and restore mechanisms?
 - The backup mechanism should provide a simple interface to complete backup and restores as well as transfer them off the local device for storage and archival.
- Do you want the ability to do automated restores?
 - Certain GSDs have the ability to do an automated or semi-automated restore using scripts or through pushable configurations from the management system.
- Should the system fail secure or insecure?
 - It is important to determine if the uptime advantage of a fail insecure system outweighs the risk of network traffic passing uninspected or the downtime associated with a fail secure system.

26 · 30 GATEWAY SECURITY DEVICES

- Do you need out-of-band access?
 - Having an out-of-band access solution such as a direct serial connection and/or dial-in solution may provide the ability to monitor and/or change the system during the network outage.

26.8.10 Usability. Having a management environment that has every feature imaginable is worthless if it is too complex and convoluted to operate and maintain.

- Is the management console intuitive?
 - The management console should have a consistent and simple interface to manage the devices and policies. The more intuitive the interface, the less time the administrators will need to spend figuring out how to complete otherwise easy tasks or troubleshooting issues.
- Are the primary functions easy to accomplish?
 - The interface should make it easy to follow and edit policies as well as view and filter logs.

26.8.11 Learning Curve. By choosing a GSD vendor with an intuitive interface and functionality, it should be much easier to learn the environment.

- Is training required to learn the new device(s) and management platform?
 - Even with the best interface and features, it may be necessary for the administrator(s) to get training to use the system effectively. Check with the vendor to see what training (basic through advanced) is available and at what cost. If possible, determine the effectiveness of the instructor(s) through referrals.
- If some of the security features are new to the organization, expect a steeper initial and increased ongoing learning curve due to the breadth and depth for each additional protection measure.
- How does the vendor approach security architecture?
 - Every GSD vendor has a little different view and implementation of security features. Some focus primarily on allowing or denying traffic, although others treat network traffic as flows between zones. Moving from one vendor type to another can be confusing and could slow or confuse the conversion.

26.8.12 Features. Although the GSD platform may be able to provide a single source for network-security protections, it is important to remember that it may not be cost or resource effective to meet the organization's security needs and posture.

- Are you looking for a device to be an all-in-one solution?
 - The ability to provide a single solution security platform has merits, but each additional feature has an impact on performance and availability, and when implemented improperly, increases the potential to cause issues.
- What features do you want the GSD to have?
 - Common UTM features may include firewall, IPS, VPN, anti-malware, anti-spam, DLP, and Quality of Service (QoS).
 - Common NGFW features include port-independent application identification and enforcement for sub-features of a particular protocol.
- Does the platform support full IPv6 enforcement across all security functions?

NETWORK-SECURITY DEVICE EVALUATION 26 · 31

26.8.13 Price. Be sure to understand the cost implications of the entire GSD infrastructure over its lifetime.

- What will be the *total* cost of ownership (TCO) of the upgrade or replacement?
 - Areas for considerations include: hardware purchase price; price for requested security features (additional features, if added later, will increase the overall costs and may require new hardware); management infrastructure (hardware and software); hardware and software maintenance fees; training costs; direct or third-party support contract; potential conversion downtime; and learning curve.

Initial Cost

- How much will the new GSD equipment cost?
 - This includes the price of the GSD, expansion cards, and management devices.
 - An appliance may be a more cost-effective alternative long term if there are multiple deployment options as the vendor is solely responsible for issues.
- How much extra will shipping cost?
 - If this is a distributed deployment, there may be extra shipping costs and tariffs associated with getting the equipment to international locations.
- How much will it cost to purchase the features I need?
 - This is highly dependent on the GSD. Certain GSDs come with everything out of the box, while others provide some basic functionality including the firewall, VPN, and basic protocol inspection but charge extra to use the full functionality of the existing items and/or to add new functionality.
- How much is it necessary to invest in training to learn the new solution?
 - Training often includes attending offsite courses (which may require additional travel expenses). In some cases, it may be more cost effective to bring a trainer in-house to provide a more tailored training program. A train-the-trainer methodology can be a money-saving option, depending on team size and funding.

Ongoing Costs

- How much is the yearly hardware and/or software maintenance?
 - The maintenance contract may be negotiable based on length of term and costs (potentially based on a percentage of the purchase).
- What level of service do you require?
 - GSD vendors provide different service levels to correspond with response times, escalation ability, and access to additional information. In accordance, the higher level of service, the more this will cost.
- Can the provider support you at each of your locations?
 - If necessary, major issues may require a vendor or vendor partner to come onsite to complete a repair. If available, determine the cost (time and travel) before using this type of service.
 - If you have international deployments, be sure to understand the vendor's ability to provide general warranty repair and/or onsite service for those locations.

26 · 32 GATEWAY SECURITY DEVICES

- Is there a requirement for an onsite spare?
 - If the DR plan requires quick recovery, having an onsite spare may help. The disadvantage to having a spare (not a redundant pair) is that it most likely will become an idle asset. Be sure to weigh this versus the cost of a higher-level service contract or local availability of replacement parts.
- Do you need additional consulting/support time?
 - Vendors and vendor partners may provide blocks of time where a dedicated resource can come onsite for additional support, planning, and platform reviews.

26.8.14 Vendor Considerations. Being able to trust that a vendor's product will protect the organization is crucial. Be sure to feel comfortable with every aspect of the product and service before turning over the network security reins. In addition to getting the features required, the vendor providing the solution should have a solid security foundation, financial stability, and adequate support resources.

- Will this vendor meet the organization's current and future needs?
 - This key question focuses on the vendor's ability to provide a quality product, support that product, and grow with your needs. Be sure to evaluate the vendor's product roadmap to see how the product (hardware, functionality, cost, etc.) will evolve and when currently unavailable features will be available in the product.
- Does the vendor already support or have experience with deployments similar to yours?
 - If the vendor already has customers with an infrastructure and needs similar to yours, this may provide additional ideas and/or information to make the deployment more successful.

26.8.14.1 Reputation. Although infighting and opinions about the best security platforms abound, the selected vendor should be an active member of the security community.

- Does this vendor have a good reputation in the security community?
 - There are no shortage of product reviews and comparisons to help determine whether a product or vendor has a solid reputation for quality and service. These are not the ultimate arbiter of a good product but are a foundation for additional research.
- Will they provide references?
 - Check with the vendor to see if another company (or two) will provide a reference as to their experience with the product and service.
- What is the vulnerability history of the current and previous systems?
 - Check vulnerability and exploit monitoring sites to determine vulnerability history.
- What is the mean time to correct vulnerabilities?
 - Although the optimal time to resolve the vulnerability is immediately (or prevent it before it happens), check to see how long it takes vendors to respond to and resolve vulnerabilities. This includes verifying the fix corrected the problem the first time.

NETWORK-SECURITY DEVICE EVALUATION 26 · 33

- What is the mean time to correct application identity issues and/or gaps?
 - The complexity and changing nature of application layer protocols requires active vendor participation. Determine the process, timelines, and potential costs associated with the vendor correcting (or supporting) application identity needs.

26.8.14.2 Support Options. Due to the complexity and breadth of functionality these devices possess, the selected vendors must have a robust support infrastructure to resolve issues.

- How experienced are the front-line and higher-tier technicians?
 - Ask the vendor to provide information regarding training and experience of the different support technician levels.
- Are there different tiers of initial call support?
 - Determine whether all support calls originate through the main support/service desk or if it is possible to have a dedicated (typically shared among other organizations supported) resource.
- Do you have the ability to escalate-on-demand?
 - Determine whether it is possible to escalate an issue immediately (and at what cost) if the local administrators complete standard troubleshooting steps.
- Are you willing to work with systems during beta testing?
 - Depending on risk tolerance, being a beta partner with the vendor may provide early access to fixes and enhancements.
- What are current and previous clients' experiences with support quality and availability?
 - Again, ask for references to get real-world perspective.

26.8.15 Managed Security Service Providers. Organizations have the ability to transfer varying levels of internal network security responsibilities to a Managed Security Service Provider (MSSP). Use of the MSSP may be an opportunity to supplement off-shift log reviews, consolidation, and/or alerting. In addition, there may be an opportunity to transfer maintenance and change control to alleviate internal resource constraints and/or knowledge gaps. Choosing this type of service requires detailed investigation.

- What is the MSSP's reputation/experience/workload?
- The MSSP should have trained and experienced personnel to support your GSD infrastructure.
- Determine the number of clients per support resource and if other MSSP locations can continue support and operations if a failure occurs at one MSSP location.
- Is it possible to do this securely?
 - Ensure the MSSP complies with secure storage of collected information and access to the GSD management infrastructure.

26 · 34 GATEWAY SECURITY DEVICES

- If necessary, how will the MSSP follow change control?
 - Determine the change control process before turning over operational responsibility. This includes determining change reviews, change windows, change approvals, service-level agreements (SLAs) regarding time to complete, etc.
- Can the MSSP meet your SLA requirements?
 - Set clear expectations during contract negotiations when establishing SLAs. The more stringent the SLA, the more likely the cost is to increase. Develop metrics to assess the MSSP's effectiveness in achieving the SLA requirements.
- At what cost?
 - The cost of doing business with the MSSP increases based on the number of services and/or devices. Be sure to investigate the contract to understand all potential fees or additional requirements that may drive up the cost after the signing of the initial contract.

26.9 CONCLUDING REMARKS. The requirements for network security are changing constantly. The organizations are under continued pressure to meet customer and business partner demands for information and access. Technologies such as mobile devices and extranets continue to blur the previous accepted reality of true network borders. Perimeter security is no longer as easy as installing a firewall or requiring proxy services for outbound connections. Reliance on a specific technology or security ideology is insufficient to protect the information and systems integral to the success of an organization. The frequency and veracity of both external and internal attacks requires a more robust and flexible mechanism to combat these increasing threats, and the GSD leapt forward to provide the current generation of perimeter protection technologies.

The GSD retains all of the previous firewall-based functionality such as allowed-path control, VPN services, and network address translation while providing the flexibility to integrate additional services designed to provide visibility into and protection from a multitude of threats. The addition of anti-malware capabilities provides a new layer of protection by relieving hosts from being the only location providing full detection, prevention, and remediation services. The integration of proxy services, application and content control, and intrusion prevention allows GSD deployments to simplify the network-security architecture and opens up new inspection and protection opportunities.

Today's threats provide little insight into how protection measures will need to evolve to meet the next-generation attacks and attackers. To remain a viable security option, the GSD security vendors must remain agile. The integration and implementation of new protection measures will need to be simple and seamless. As processing power continues to increase, so will the capabilities of the GSD to take on greater workloads and complexity. Basic content inspection will give way to a greater understanding of and protection for data context, value, and flows. The support and protection of worldwide networks will require providing the full existing set of GSD functionality for IPv6 traffic.

No matter the threat, no matter the network, no single device is a silver bullet to complete security. Each organization must evaluate all avenues of protection to ensure that the technologies deployed meet the security functionality required.

FURTHER READING 26 · 35

26.10 FURTHER READING

- Forouzan, Behrouz A. *TCP/IP Protocol Suite*, 4th ed. McGraw-Hill, 2009.
- McClure, Stuart, Joel Scambray, and George Kurtz. *Hacking Exposed*, 7th ed. McGraw-Hill Osborne Media, 2012.
- Scarfone, Karen and Paul Hoffman. "Guidelines on Firewalls and Firewall Policy: Recommendations of the National Institute of Standards and Technology." NIST Special Publication 800-41 Revision 1, 2009. <http://csrc.nist.gov/publications/nistpubs/800-41-Rev1/sp800-41-rev1.pdf>
- Thomson, S., T. Narten, and T. Jinmei. "RFC 486: IPv6 Stateless Address Autoconfiguration." Internet Engineering Task Force, 2007. <https://tools.ietf.org/html/rfc4862>
- Narten, T., E. Nordmark, W. Simpson, and H. Soliman. "Neighbor Discovery for IP version 6 (IPv6)." Internet Engineering Task Force, 2007. <http://tools.ietf.org/html/rfc4861>
- Deering, S., W. Fenner, and B. Haberman. "Multicast Listener Discovery (MLD) for IPv6." Internet Engineering Task Force, 1999. <http://tools.ietf.org/html/rfc2710>
- Johnson, D., C. Perkins, and J. A. Ericsson. "Mobility Support in IPv6." Internet Engineering Task Force, 2004. <http://tools.ietf.org/html/rfc3775>

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER **27**

INTRUSION DETECTION AND INTRUSION PREVENTION DEVICES

Rebecca Gurley Bace

27.1 SECURITY BEHIND THE FIREWALL	27·2	27.4.3 Application Monitoring	27·7
27.1.1 What Is Intrusion Detection?	27·2	27.4.4 Other Types of Monitoring	27·7
27.1.2 What Is Intrusion Prevention?	27·3	27.4.5 Issues in Information Sources	27·8
27.1.3 Where Do Intrusion Detection and Intrusion Prevention Fit in Security Management?	27·3		
27.1.4 Brief History of Intrusion Detection	27·4		
		27.5 ANALYSIS SCHEMES	27·8
		27.5.1 Misuse Detection	27·8
		27.5.2 Anomaly Detection	27·9
		27.5.3 Hybrid Approaches	27·10
		27.5.4 Issues in Analysis	27·10
27.2 MAIN CONCEPTS	27·4		
27.2.1 Process Structure	27·5	27.6 RESPONSE	27·10
27.2.2 Monitoring Approach	27·5	27.6.1 Passive Responses	27·10
27.2.3 Intrusion Detection Architecture	27·5	27.6.2 Active Responses: Man-in-the-Loop and Autonomous	27·11
27.2.4 Monitoring Frequency	27·5	27.6.3 Automated Response Goals	27·11
27.2.5 Analysis Strategy	27·6	27.6.4 Investigative Support	27·12
		27.6.5 Issues in Responses	27·12
27.3 INTRUSION PREVENTION	27·6		
27.3.1 Intrusion Prevention System Architecture	27·6	27.7 NEEDS ASSESSMENT AND PRODUCT SELECTION	27·13
27.3.2 Intrusion Prevention Analysis Strategy	27·6	27.7.1 Matching Needs to Features	27·13
		27.7.2 Specific Scenarios	27·14
		27.7.3 Integrating IDS Products with Your Security Infrastructure	27·14
27.4 INFORMATION SOURCES	27·6		
27.4.1 Network Monitoring	27·7		
27.4.2 Operating System Monitoring	27·7		

27 · 2 INTRUSION DETECTION AND INTRUSION PREVENTION DEVICES

27.7.4 Deployment of IDS Products	27.15	27.9 FURTHER READING	27.16
27.8 CONCLUSION	27.16	27.10 NOTES	27.17

27.1 SECURITY BEHIND THE FIREWALL. Even today, when asked how they would go about securing a computer or computer network, most people mention firewalls, the first widely accepted network security devices. As network security has become a nonoptional facet of system management, firewall mechanisms of various sorts have become a standard fixture in many networks.

As in any complex protection function, firewalls are necessary but not sufficient to completely protect enterprises from security breaches. Expecting firewalls to provide complete protection is tantamount to expecting guards at the gates of corporate campuses to prevent destructive behavior on the part of those operating vehicles within the facility. A lot can happen as traffic flows between hosts on internal networks; even items that appear benign to gatekeepers can be subverted to carry destructive payloads. Thus, most modern network security architectures include intrusion detection and intrusion prevention systems.

Intrusion detection systems (IDSs) are software or hardware systems that automate the monitoring of events occurring within a computer system or network. IDSs not only collect and synchronize records of these events; they also analyze them for signs of security violations. In strictest terms, vulnerability assessment systems (VASs) are a special class of IDSs in which the system relies on static inspection and attack reenactment to gauge a target system's exposure to specific security vulnerabilities. Intrusion prevention systems (IPSs) are another special class of IDSs in which the system is designed to react to certain detected attacks in a prespecified way.

In recent years, the realm of intrusion detection has broadened and deepened, driven by a number of factors. Security personnel, architectural and operational alike, have gained more experience with IDS technologies, working with commercial product providers to expand product capabilities and management schemes to fit current needs. As operational personnel have become more comfortable in using these systems, they have moved certain IDS capabilities to core network management venues. Finally, threats have evolved, driving needs that IDSs are uniquely qualified to address.

The evolution of IDS has resulted in some changes in nomenclature and tradecraft associated with it. One of these changes is that vulnerability assessment is considered a stand-alone discipline, driven by market needs that are often different from those influencing IDS. Another is that intrusion prevention has evolved as a stand-alone product category, offering the ability to respond automatically to certain classes of attacks. Thus, vulnerability assessment is treated as a separate topic in this *Handbook* (covered in Chapter 46), while both IDS and IPS are discussed in this chapter.

27.1.1 What Is Intrusion Detection? Intrusion detection is the process of collecting information about events occurring in a computer system or network and analyzing them for signs of *intrusions*. *Intrusions* are defined as violations of security policy, usually characterized as attempts to affect the confidentiality, integrity, or availability of a computer or network. These violations can come from attackers accessing systems from the Internet or from authorized users of the systems who attempt to overstep their legitimate authorization levels or who use their legitimate access to the system to conduct unauthorized activity.

SECURITY BEHIND THE FIREWALL 27 · 3

Intrusion detection systems are software or hardware products that automate this monitoring and analysis process.

27.1.2 What Is Intrusion Prevention? Intrusion prevention is the process of coupling intrusion detection (as defined) with specified responses to certain detected intrusion scenarios. The triggering events can be viewed as a special subset of intrusions, and are often characterized in richer quantitative and qualitative terms than more generic IDS triggers. For example, a specific IPS might focus on monitoring certain types of network traffic. When the rate of a particular traffic type exceeds the anticipated threshold, the IPS would react in a prespecified way (e.g., limiting the rate of subsequent traffic of that type).

27.1.3 Where Do Intrusion Detection and Intrusion Prevention Fit in Security Management? Intrusion detection is a necessary function in most system security strategies. It (and its offshoot, vulnerability assessment) is the primary security technology that supports the goal of *auditability*. “Auditability” is defined as the ability to independently review and examine system records and activities to:

- Determine the adequacy of system controls
- Ensure compliance with established security policy and operational procedures
- Detect breaches in security
- Recommend any indicated changes¹

The presence of a strong audit function, in turn, enables and supports several vital security management functions, such as incident handling and system recovery. Intrusion detection also allows security managers a flexible means to accommodate user needs while retaining the ability to protect systems from certain types of threats.

There is some debate over whether IPS will displace IDS in the security management lineup. This displacement is unlikely for the time being because of the binding between IDS and audit. As security operations become more tightly integrated with traditional system administration and operations, audit functions are necessary to support root cause analysis in diagnosing and addressing system failure. Badly tuned security devices can create problems in operations; it is important to be able to identify and correct such issues quickly and appropriately. Audit functions are also necessary to measure the effectiveness of security measures in mitigating security threats. Although it may appear that transparently detecting and blocking attacks is the optimal security process, such transparency—that is, lack of an audit trail—interferes with demonstrating effectiveness of security measures against real threats. In other words, in order for security management to justify a budget for security measures, it must be able to document and quantify the suitability and effectiveness of such measures. Therefore, regardless of the effectiveness of IPS in blocking attacks, IDS will likely always be necessary to support these audit functions. IDS provides the baseline information on the attack profiles and frequencies that allow managers to demonstrate the effectiveness and return on investment of preventive mechanisms. Without hard evidence of attack frequencies, security managers are left in the position of the man waving a dead chicken around his head while standing on a street corner. Asked why he is doing that, he answers, “To keep the flying elephants away.” “But there are no flying elephants,” protest the observers. “See? It works!” replies the lunatic. For more information on auditing, see Chapter 54 in this *Handbook*.

27 · 4 INTRUSION DETECTION AND INTRUSION PREVENTION DEVICES

Although intrusion detection and intrusion prevention systems are necessary as system security functions, they are not sufficient to protect systems from all security threats. IDS and IPS must be a part of a more comprehensive security strategy that includes vulnerability assessment, security policy and procedural controls, network firewalls, strong identification and authentication mechanisms, access control mechanisms, file and link encryption, file integrity checking, physical security measures, and security training.

27.1.4 Brief History of Intrusion Detection. Intrusion detection is the automation of manual processes that originated in the earliest days of data processing. Joseph Wassermann of the Bell Telephone Company documented the origin of system and security audit as early as the mid-1950s, when the first computerized business system was being designed and implemented.²

Auditability was a key security feature from the earliest days of computer security, as proposed in J. P. Anderson's 1973 research study chartered by the U.S. Air Force.³ Anderson proposed a scheme for automating the review of security audit trails in 1980, in a research report considered by many to be the seminal work in intrusion detection.⁴ Dorothy Denning and Peter Neumann led a study of intrusion detection, conducted from 1984 to 1986, producing another seminal work in intrusion detection in 1986, in which Denning proposed a model for intrusion detection.⁵

An instantiation of Denning's intrusion detection model was prototyped as the Intrusion Detection Expert System (IDES) by a team at SRI International. IDES was a hybrid system that constructed statistical profiles of user behaviors as derived from operating system kernel audit logs and other system data sources. IDES also provided a rules-based expert system that allowed users to specify patterns of events to be flagged as intrusions.⁶ IDES and the Next Generation IDES (NIDES) system that followed it marked an era in which numerous intrusion detection research projects and prototype systems were developed, including Haystack (Haystack Labs and U.S. Air Force), NADIR (Los Alamos National Laboratory), Wisdom and Sense (Los Alamos National Laboratory and Oak Ridge National Laboratory), ISOA (PRC, Inc.), TIM (Digital Equipment Corporation), ComputerWatch (AT&T), and Discovery (TRW, Inc).⁷

In the late 1980s, researchers at the University of California, Davis, designed the first network-based intrusion detection system (initially called the Network Security Monitor, but later renamed NID), which functioned much the same as many current commercial network-based intrusion detection products.⁸ A subsequent U.S. Air Force–funded research product called the Distributed Intrusion Detection System (DIDS), explored coordinating network-based and host-based intrusion detection systems. DIDS was prototyped by teams at the University of California, Davis, Haystack Laboratories, and Lawrence Livermore National Laboratory.⁹

Intrusion prevention was proposed as a logical next step to intrusion detection almost from the start of intrusion detection research. Support for certain models of intrusion prevention grew when concerns regarding attacks on the TCP/IP network infrastructure (e.g., packet flooding and malformed packet attacks) grew in the mid- to late-1990s. Other types of IPS were proposed to deal with kernel-level hacks and information leakage issues.

27.2 MAIN CONCEPTS. Several strategies used in performing intrusion detection serve to describe and distinguish specific intrusion detection systems. These affect the threats addressed by each system and often prescribe the environments in which

MAIN CONCEPTS 27 · 5

specific systems should be used. As noted, intrusion prevention relies first on intrusion detection strategies; thus, the differentiators will be highlighted as appropriate.

27.2.1 Process Structure. Intrusion detection is defined as a monitoring and alarm generation process, and, as such, it can be described using a simple process model. This model is outlined here and will be used to illustrate the fundamental concepts of intrusion detection.

27.2.1.1 Information Sources. The first stage of the intrusion detection process comprises one or more information sources, also known as event generators. Information sources for intrusion detection may be categorized by location: network, host, or application.

27.2.1.2 Analysis Engine. Once event information is collected, it is passed to the next stage of the intrusion detection process, in which it is analyzed for symptoms of attack or other security problems.

27.2.1.3 Response. When the analysis engine diagnoses attacks or security problems, information about these results is revealed via the response stage of the intrusion detection process. Responses span a wide spectrum of possibilities, ranging from simple reports or logs to automated responses that disrupt attacks in progress. The presence of these automated responses defines an intrusion prevention system.

27.2.2 Monitoring Approach. The first major classifier used to distinguish intrusion detection systems is the monitoring approach of the system. Monitoring is the action of collecting event data from an information source and then conveying that data to the analysis engine.

The monitoring approach describes the perspective from which intrusion detection monitoring is performed. The primary monitoring approaches found in intrusion detection systems today are *network based*, *host based*, and *application based*.

27.2.3 Intrusion Detection Architecture. Even in the early days of manual security audit, researchers noted that in order for audit information to be trusted, it should be stored and processed in an environment separate from the one monitored. This requirement has evolved to include most intrusion detection approaches, for three reasons:

1. To keep an intruder from blocking or nullifying the intrusion detection system by deleting information sources
2. To keep an intruder from corrupting the operation of the intrusion detector in order to mask the presence of the intruder
3. To manage the performance and storage load that might result from running intrusion detection tasks on an operational system

In this architecture, the system running the intrusion detection system is called the *host*. The system or network being monitored is called the *target*.

27.2.4 Monitoring Frequency. Another common descriptor for intrusion detection approaches is the timing of the collection and analysis of event data. This is

27 · 6 INTRUSION DETECTION AND INTRUSION PREVENTION DEVICES

usually divided between *batch-mode* (also known as *interval-based*) and *continuous* (also known as *real-time*) approaches.

In batch-mode analysis, the event data from the information source are conveyed to the analysis engine in a file or other block form. As the name suggests, the events corresponding to a particular interval of time are processed (and results provided to the user) after the intrusion has taken place. This model was the most common for early intrusion detection because system resources did not allow real-time monitoring or analysis.

In real-time analysis, event data from the information source are conveyed to the analysis engine as the information is gathered. The information is analyzed immediately, providing the user with the opportunity to respond to detected problems quickly enough to affect the outcome of the intrusion.

27.2.5 Analysis Strategy. In intrusion detection, there are two prevalent analysis strategies, *misuse detection* and *anomaly detection*.

In misuse detection, the analysis engine filters event streams, matching patterns of activity that characterize a known attack or security violation. In anomaly detection, the analysis engine uses statistical or other analytical techniques to spot patterns corresponding to abnormal system use. Anomaly detection is based on the premise that intrusions significantly differ from normal system activity. In general, intrusion detection systems rely more heavily on anomaly detection and quantitative measures to detect and block attacks.

27.3 INTRUSION PREVENTION. As discussed, IPSs are often considered a special case of IDS in which automated responses are specified. However, with the adoption of the first generation of network IPS products, additional specifications for IPSs have evolved in addition to those assigned to IDS.

27.3.1 Intrusion Prevention System Architecture. As in intrusion detection, most IPSs separate the monitoring and analysis platform from the target platform being monitored. Additional differentiations are drawn between those IPSs that separate the monitoring and analysis platform from the response platform (these are labeled “stand-alone” IPSs) and those that integrate all the functions in a single unit, usually a firewall, network switch, or router (these are labeled “integrated” IPSs.).

27.3.2 Intrusion Prevention Analysis Strategy. IPSs generally use the same structural approach to data analysis as IDS, but the nomenclature for the analysis strategies differs. IPS analysis schemes fall into two general categories, rate based and content based.

Rate-based IPS analysis makes the decision to block network traffic based on indicators of network load, as measured by statistics such as connect rates and connection counts. This category of analysis is especially useful for detecting packet flood distributed denial-of-service (DDoS) attacks.

Content-based IPS analysis makes the decision to block network traffic based on indicators of anomalous packets and specific content (often represented as IDS signatures). This approach is useful for detecting malformed packet DDoSs and other types of attacks not readily spotted by quantitative measures.¹⁰

27.4 INFORMATION SOURCES. Information sources represent the first stage of the intrusion detection process. They provide event information from monitored

INFORMATION SOURCES 27 · 7

systems on which the intrusion detection process bases its decisions. Information sources encompass both *raw* event data (e.g., data collected directly from system audit and logging mechanisms) as well as data output by system management utilities (e.g., file integrity checkers, vulnerability assessment tools, network management systems, and even other intrusion detection systems). In this section, information sources for intrusion detection are classified by location: network, host, or application.

27.4.1 Network Monitoring. The most common monitoring approach utilized in intrusion detection systems is *network based*. In this approach, information is gathered in the form of network packets, often using network interface devices set to promiscuous mode. (Such a device operating in promiscuous mode captures all network traffic accessible to it—usually on the same network segment—not just traffic addressed to it.) Other approaches for performing network-based monitoring include the use of *spanning ports* (specialized monitoring ports that allow capture of network traffic from all ports on a switch) on network switches or specialized Ethernet network taps (e.g., sniffers) to capture network traffic.

27.4.2 Operating System Monitoring. Some monitors collect data from sources internal to a computer. These differ from network-based monitoring in the level of abstraction at which the data is collected. *Host-based* monitoring collects information from the operating system (OS) level of a computer. The most common sources of operating system-level data are operating system audit trails, which are usually generated within the OS kernel, and system logs, which are generated by OS utilities.

27.4.3 Application Monitoring. *Application-based* monitoring collects information from running software applications. Information sources utilized in application-based approaches include application event logs and application configuration information.

Application-based information sources are steadily increasing in importance as systems complexity increases. The advent of object-oriented programming techniques introduces data object naming conventions that nullify much of an analyst's ability to make sense of file access logs. In this situation, the application level is the only place in the system in which one can “see” the data accesses at an appropriate level of abstraction likely to reveal security violations.

One special case of application-based monitoring comprises an entire product category in security. This type of system (sometimes called extrusion detection) monitors data transfers, looking for anomalies associated with data movement across policy boundaries. Such data monitoring is very popular as a compliance mechanism for enterprises dealing with consumer, financial, or other regulated data.

27.4.4 Other Types of Monitoring. As noted, intrusion detection information sources are not limited to raw event data. In fact, allowing intrusion detection systems to operate on results from other systems often optimizes the quality of the intrusion detection system's results. When the data are provided by other parts of the system security infrastructure (e.g., network firewalls, file integrity checkers, virus scanners, or other intrusion detection systems), the sensitivity and reliability of the intrusion detection system's results can increase significantly.

27 · 8 INTRUSION DETECTION AND INTRUSION PREVENTION DEVICES

27.4.5 Issues in Information Sources. There are several issues involving information sources for intrusion detection. The major ones that have persisted over the history of intrusion detection include:

- In host-based systems, there must be a balance between collecting enough information to accurately portray security violations and collecting so much information that the collection process cripples the monitored system.
- The fidelity of the intrusion detection process is dependent not only on collecting the appropriate information but on collecting it from appropriate vantage points within the monitored system or network.
- If the IDS is expected to produce event records that will be used to support legal processes, the system must collect and handle event information in a way that complies with legal rules of evidence.
- The information collected by IDSs often includes information of a sensitive nature. This information must be secured and handled in a way that complies with legal and ethical standards.

27.5 ANALYSIS SCHEMES. Once information sources and sensors are defined and placed, the information so gathered must be analyzed for signs of attack. The *analysis engine* serves this purpose in intrusion detection, accepting event data from the information source and examining it for symptoms of security problems. As mentioned earlier, intrusion detection systems typically provide analysis features that fall into two categories, *misuse detection* and *anomaly detection*.

27.5.1 Misuse Detection. *Misuse detection* is the filtering of event streams for patterns of activity that reflect known attacks or other violations of security policy. Misuse detectors use various pattern-matching algorithms, operating on large databases of attack patterns or *signatures*. Most current commercial intrusion detection systems support misuse detection.

Misuse detection presumes that there is a clear understanding of the security policy for the system, which can be expressed in patterns corresponding to desirable activity and undesirable activity. Therefore, signatures can be described in terms of “this should never happen” as well as “only this should ever happen.” Signatures also can range from simplistic *atomic* (one-part) checks to rather complex *composite* (multipart) checks. An example of an atomic check is a buffer overflow signature, in which one looks for a particular command, followed by a string exceeding a particular length. An example of a composite check is a race condition signature, in which a series of carefully timed commands occur. Signatures are gathered and structured in some way to optimize the filtering of event data against them.

The next requirement for misuse detection is that the event data collected from information sources be encoded in a way that allows it to be matched against the signature data. There are various ways of doing this, ranging from regular expression matching (sometimes called “dirty word” matching) to complex coding schemes involving state diagrams and Colored Petri Nets. State diagrams are a graphical scheme for modeling intrusions. They express intrusions in terms of *states*, represented by nodes or circles, and *transitions*, represented by lines or arcs. Colored Petri Nets are an extension of the state diagram technique that add colored *tokens*, which occupy state nodes, and whose color expresses information about the context of the state.

ANALYSIS SCHEMES 27 · 9

In some IDSs and content-based IPSs, significant resources are devoted to identifying malformed network packets, especially those in which the format of the content of the packet does not match the format of the service (e.g., SMTP) or of the associated port number of the packet. This malformed packet scheme represents one of the major categories of DDoS attacks, which seek to deny network access to legitimate users.

27.5.2 Anomaly Detection. *Anomaly detection* is the analysis of system event streams, characterizing them using statistical and other classification techniques in order to find patterns of activity that appear to deviate from normal system operation. This approach is based on the premise that attacks, and other security policy violations, are a subset of abnormal system events.

Several common techniques are used in anomaly detection:

- **Quantitative analysis.** Most modern systems that use anomaly detection provide quantitative analysis, in which rules and attributes are expressed in numeric form. The most common forms of quantitative analysis are triggers and thresholds, in which system attributes are expressed as counts occurring during some time interval, with some level defined as permissible. Triggers and thresholds can be simple, in which the permissible level is constant, or heuristic, in which the permissible level is adapted to observed levels. Network intrusion prevention systems targeting DDoS attacks often use heuristic triggers and thresholds to characterize the normal bandwidth loads, connect rates, and connection counts in network traffic.
- **Statistical analysis.** Most early anomaly detection systems used statistical techniques to identify abnormal data. In statistical analysis, profiles are built for each user and system resource, and statistics are calculated for a variety of user and resource attributes for a particular interval of time (usually a “session,” defined as the time elapsed between login and logout).
- **Learning techniques.** There has been a great deal of research interest in using various learning techniques, such as neural networks and fuzzy logic, in performing anomaly detection. Despite encouraging results, there remain many practical impediments to using these techniques in production environments. The practical impediments arise due to the mismatch between those attributes that are suitable for characterization by neural networks and fuzzy logics and those attributes that are actionable by operational personnel and systems. The value of using neural networks (most of them utilizing fuzzy logic) is that they can characterize and recognize very subtle signs of trouble in systems. This is of value in situations where the problems being detected are not subtle. However, in security breaches, the difference between normal behavior and security attack is often influenced by the system context; in these scenarios, few, if any, neural networks can provide insights regarding how they reached their decisions regarding the suspicious events on which they trigger. A security person in an operational context usually requires this sort of insight in order to devise an appropriate reaction to the detected intrusion.
- **Advanced techniques.** Anomaly detection as applied to intrusion detection remains an active research area. Recent research efforts include the application of such advanced analytic techniques as genetic algorithms, data mining, autonomous agents, and immune system approaches to the problem of recognizing new attacks

27 · 10 INTRUSION DETECTION AND INTRUSION PREVENTION DEVICES

and security violations. Again, these techniques have not yet been widely fielded in commercial IDSs, although they have appeared in special purpose products.

27.5.3 Hybrid Approaches. There are significant issues associated with both misuse detection and anomaly detection approaches to event analysis for intrusion detection; however, combining both approaches provides considerable benefit. The anomaly detection engine can allow the IDS to detect new or unknown attacks or policy violations. This is especially valuable when the target system protected by the IDS is highly visible on the Internet or other high-risk network. In IPS, anomaly detection applied to network traffic attributes allows one of the only means to deal with packet-flood DDoS attacks, a growing concern in today's networks.

The misuse detection engine, in turn, protects the integrity of some anomaly detection engines by assuring that a patient adversary cannot gradually change behavior patterns over time in order to retrain the anomaly detector to accept attack behavior as normal. Thus, the misuse detector mitigates a significant deficiency of anomaly detection for security purposes.

27.5.4 Issues in Analysis. Here are a few of the many issues in intrusion detection analysis:

- Misuse detection systems, although very effective at detecting those scenarios for which detection signatures have been defined, cannot detect new attacks.
- Anomaly detection systems are capable of detecting new attacks but usually have false positive rates so high that users often ignore the alarms they generate.
- Anomaly detection systems that rely on artificial intelligence (AI) techniques often suffer from a lack of adequate training data. (Data are used to define the detector's logic for distinguishing "normal" from "abnormal" events.)
- Malefactors with access privileges to the system, while anomaly-detection systems are being trained, can covertly teach the system to accept specific patterns of unauthorized activities as normal. Later, the anomaly-detection systems will ignore the actual misuse.

27.6 RESPONSE. The final stage of intrusion detection, response, consists of the actions taken in response to the security violations detected by the IDS. Responses are divided into *passive* and *active* options. The difference between passive and active responses is whether the user of the IDS, or the system itself, is responsible for reacting to the detected violations. As discussed, the former option is associated with classic IDS; the latter is associated with IPS.

27.6.1 Passive Responses. When passive responses are selected, the IDS simply provides the results of the detection process to the user, who must then act on these results, independent of the IDS. In this option, the user has total control over the response to the detected problem. In some IDSs, the information provided to the user regarding detection results can be divided into *alarms* and *reports*.

27.6.1.1 Alarms. Alarms are messages that are communicated immediately to users. Commercial IDSs use a variety of channels for conveying these alarms to security personnel. The most common is a message screen or icon written to the IDS control

RESPONSE 27 · 11

console. Other alarm channels include pagers, email, wireless messaging, and network management system traps.

27.6.1.2 Reports. Reports are messages or groups of messages that are generated on a periodic basis. They typically document events that have happened in the past and often include aggregate figures and trends information. Many commercial IDS products support extensive reporting features, allowing a user to set up automatic report generation with several versions, each targeting a different level of management.

27.6.2 Active Responses: Man-in-the-Loop and Autonomous. When an IDS provides active response options, these usually fall into two categories. The first requires the IDS to take action, but with the active involvement of an interactive user. This option is sometimes called a man-in-the-loop mechanism. This option is preferred for critical systems, as it allows an operator to track an attacker or intervene in a sensitive situation in a flexible, exacting way.

The other active response option, which usually defines an IPS, provides for preprogrammed actions taken automatically by the system with no human involvement. The automated response option is required when dealing with certain sorts of automated attack tools (viruses or worms) or DDoS attacks. These attacks proceed at machine speed and therefore are outside the reach of a human-controlled manual intervention. As automated and DDoS attacks are the tool of choice for online extortionists, the number of IPSs fielded in commercial networks has grown rapidly over the past few years.

27.6.3 Automated Response Goals. Automated responses support three categories of response goals:

1. Collecting more information about the intrusion or intruder
2. Amending the environment (e.g., changing a switch or router setting to deny access to an intruder)
3. Taking action against the intruder

Although the last of these groups, sometimes labeled *strike back* or *hack back*, occasionally is discussed in security circles, the other options are far more productive in most situations. At this time, taking action against the intruder is considered inappropriate in almost all situations, and should be undertaken only with the advice and counsel of a legal authority.

Amending the environment and collecting more information can occur in either stand-alone or integrated fashions.

27.6.3.1 Stand-Alone Responses. Some automated responses are designed to use features that fall entirely within the intrusion detection system. For instance, an intrusion detection system may have special detection rules, more sensitive or detailed than those provided in normal modes of operation. In a stand-alone adaptive response, the IDS would use the more sensitive rules when evidence of the preamble of an attack is detected. This allows the IDS to turn sensitivity levels up only when the additional detection capabilities are needed, so that false alarm rates are reduced.

27 · 12 INTRUSION DETECTION AND INTRUSION PREVENTION DEVICES

27.6.3.2 Integrated Responses. The response option often considered the most productive is that of using integrated measures that change the system settings to block the attacker's actions. Such responses can affect the configuration of the target system, the IDS/IPS host, or the network on which both reside. In the first case, the IDS/IPS might change the settings of the logging mechanisms on the target host to increase the amount or type of information collected. The IDS also might change its analysis engine so that more subtle signs of attack are recognized. In another response option reflected in commercial products, the IDS responds to an observed attack signature by querying the target system to determine whether it is vulnerable to that specific attack. Should the vulnerability be present, the IDS directs the target system to correct that vulnerability. In effect, this process provides the target system with an immune function and permits it to "heal" itself, either blocking an attack outright or else interactively repairing any damage done in the course of the attack. Finally, some systems may use special-purpose decoy systems, called *honey pots* or *padded cells*, as diversions for attackers. When these systems are provided, the IDS may be configured to divert attackers into the decoy environments.

In a special case of integrated response seen in many commercial IPS offerings, the IPS is integrated with a switch or router. When an attack is detected, the switch or router is reconfigured on the fly to block the source of the attack. Other IPS offerings that are designed to deal with DDoS attacks use multiple IDS/IPSSs to detect the attacks, then manipulate the switching fabric¹¹ to divert the attacks from the targeted systems. Over time, such IPS features will likely be integrated with network infrastructure devices, as many firewall features already have done.

27.6.4 Investigative Support. Although the primary design objective of intrusion detection systems is detecting attacks and other possibly problematic system events, information collected and archived by IDSs also can support those charged with investigating security incidents. This functional requirement may levy additional technical requirements on IDSs. For instance, if investigators plan to use IDS monitoring features to perform a targeted surveillance of an attack in progress, it is critical that the information sources be "silent" so that adversaries are not aware that they are being monitored. Furthermore, the IDS monitors must be able to convey information to the investigators through a trustworthy, secure channel. Finally, the IDS itself must be under the control of the investigators or other trusted parties; otherwise, the adversaries may mask their activities by selectively spoofing information sources. Perhaps the most important thing for investigators to remember about IDSs is that the information provided should be corroborated by other information sources (e.g., network infrastructure device logs), not necessarily accepted at face value.

27.6.5 Issues in Responses. As in information sources and analysis strategies, certain issues associated with IDS response features have endured over the history of intrusion detection. The principal issues are:

- Users' needs for IDS response capabilities are as varied as the users themselves. In some systems environments, the IDS response messages are monitored around the clock, with real-time action taken by system administrators based on IDS alarms. In other environments, users may use IDS responses, in the form of reports, as a metric to indicate the threat environment in which a particular system resides. It is important to consider the specific needs of the user when selecting an IDS.

NEEDS ASSESSMENT AND PRODUCT SELECTION 27 · 13

- Given false-positive error rates for IDSs, response options must be tunable by users. Otherwise, users will simply tune out the IDS responses. This nullifies the value of the IDS.
- When the IDS provides automated responses to detected problems, there is a risk of the IDS itself launching an effective denial-of-service attack against the system it is protecting. For instance, suppose an IDS is configured with rules that tell it “upon detecting an attack from a given IP address, direct the firewall to block subsequent access from that IP address.” An attacker, knowing this IDS is so configured, can launch an attack with a forged IP source address that appears to come from a major customer or partner of the organization. The IDS will recognize the attack and then block access from that organization for some period of time, effecting a denial of service.

27.7 NEEDS ASSESSMENT AND PRODUCT SELECTION. The value of intrusion detection products within an organization’s security strategy is optimized by a thorough needs assessment. These needs and security goals can be used to guide the selection of products that will enhance the security stance of the organization.

27.7.1 Matching Needs to Features. The needs most often addressed by intrusion detection and intrusion prevention systems include:

- Prevention of problem behaviors by increasing the risk of discovery and punishment for system attackers.
- Detection of security violations not prevented (or even, in some cases, not preventable) by other security measures.
- Documentation of the existing level of threat to an organization’s computer systems and networks.
- Detection and, where possible, mitigation of attack preambles. (These include activities such as network probes, port scans, and other such “doorknob rattling.”)
- Diagnosis of problems in other elements of the security infrastructure (i.e., malfunctions or faulty configurations).
- Granting system security personnel the ability to test the security effects of maintenance and upgrade activities on the organizational networks.
- Providing information about those violations that do take place, enabling investigators to determine and correct the root causes.
- Providing evidence of compliance with a given regulatory requirement for information protection. This represents a significant need for members of various regulated industries, such as banking and health care.

Regardless of which of these specific needs are relevant to the user, it is important to consider the ability of the intrusion detection system to satisfy the needs of the specific environment in which it is installed. A critical part of this determination is considering whether the intrusion detection system has the ability to monitor the specific information sources available in the target environment. What is even more important is whether the organizational security policy translates into a monitoring and detection policy that can be used to configure the IDS (or in the case of an IPS, a monitoring, detection, and response policy.) The structure of the security policy is especially critical to the success of an IPS.

27 · 14 INTRUSION DETECTION AND INTRUSION PREVENTION DEVICES

27.7.2 Specific Scenarios. There is no universally applicable description for computer networks or the IDSs that protect them. There are, however, some common scenarios, given current trends in networking and system usage.

A popular justification for using IDSs early in an organization's security life cycle is to establish the threat level for a given network enclave. Network-based IDSs often are used for this purpose, with monitors placed outside the organizational firewall. Those who are responsible for winning management support for security efforts often find this use of IDSs to be quite helpful.

Many organizations use IDSs to protect Web servers. In this case, the nature of the interactions that the Web server has with users will affect the selection and configuration of the IDS. Most Web servers serve two types of functions: (1) informational (e.g., Web servers that support simple HTTP and FTP queries from users) and (2) transactional (e.g., Web servers that allow user interaction beyond simple HTTP or FTP traffic). Transactional Web servers are usually more difficult to monitor than informational servers, as the range of interactions between users and servers is wider. For critical transactional Web servers, security managers may wish to consider multiple IDSs, monitoring the servers at multiple levels of abstraction (i.e., application, host, and network).

The third scenario involves organizations that wish to use IDSs as additional protection for specific portions of their networked systems. An example of this is the medical organization that wishes to protect the patient record database systems from privacy breaches. In this situation, as in the Web server example just given, it may be advisable to use multiple IDSs, monitoring interactions at multiple levels of abstraction. The output of these multiple systems can be synchronized and inconsistencies noted for a reliable indication of threat levels. Another example that is increasingly common is the organization that is concerned about wireless connectivity. In this case, WiFi monitoring products are commercially available, with information collection and monitoring features that are similar to those of classic IDSs.

In recent years, the expanding use of wireless local area networks (WLANs) in buildings and campuses has stimulated the development of wireless intrusion detection and prevention systems (WIDPSs). The basic principles are the same as for other IDSs and IPSs, with the addition of an interesting wrinkle: The WIDPSs are often used to discover unauthorized access points installed on an organization's WLANs by rogue employees or by intruders. Karen Scarfone and Peter Mell, in their July 2012 edition of NIST SP 800-94, point out that WIDPS sensors should be placed in areas where there should be no wireless network activity. Some security managers walk through and drive around their facilities with WIDPS tools on their laptop computers to identify unauthorized access points. However, completely passive eavesdropping on WLAN traffic cannot be detected.¹²

In extensive networks, should the security architect decide to layer multiple IDSs and IPSs, a security event monitor/security information manager (SIM/SEM) may be required. Such a system would be necessary to consolidate and integrate the results of each IDS/IPS into a coherent set of conclusions.

27.7.3 Integrating IDS Products with Your Security Infrastructure.

As mentioned, an IDS is not a substitute for a firewall, virtual private network, identification and authentication package, or any other security point product. However, an IDS can improve the quality of protection afforded by the other point products by monitoring their operation, noting signs of malfunction or circumvention. Furthermore,

NEEDS ASSESSMENT AND PRODUCT SELECTION 27 · 15

an IPS can interact in concert with the rest of the point products to help block an attack in progress.

27.7.4 Deployment of IDS Products. The first generations of IDS installations have yielded some insights associated with deployment of IDSs. The key points include the location of sensors, scheduling the integration of IDSs, adjusting alarm settings, and outsourcing IDS/IPS services.

27.7.4.1 Location of Sensors. There are four general locations for IDS sensors:

1. Outside the main organizational firewall
2. In the network DMZ (inside the main firewall, but outside the internal firewalls)
3. Behind internal firewalls
4. In critical subnets, where critical systems and data reside

As mentioned, IDS sensors placed outside the main organizational firewall are useful for establishing the level of threat for a given network. Sensors placed within the DMZ¹³ can monitor for penetration attempts targeting Web servers. IDSs monitors for internal attacks are placed on internal network segments, behind internal firewalls. And for critical subnets, IDS sensors usually are placed at the choke points at which the subnets are connected to the rest of the corporate network. In the case of wireless networking, specialized IPS devices serve to discover and report unauthorized access to networks via WLAN access points or open software access points, through misconfigured WLAN interfaces on laptops and other WLAN devices. These IPS devices are usually placed behind firewalls and distributed across the physical space occupied by the organization and its users.

27.7.4.2 IDS Integration Scheduling. Early generations of intrusion detection products proved that integration processes must not be rushed. IDSs still rely on operator interactions to screen out false alarms and to act on legitimate alarms. Hence, it is critical that the processes provide adequate time for operational personnel to learn the behavior of the IDS on target systems, developing a sense of how the IDS interoperates with particular system components in different situations. This wisdom applies even more to IPS installation, where a miscue in specifying a response can have disastrous effects on the function of critical networks.

27.7.4.3 Alarm Settings. IDSs have significant false alarm rates, with false positive rates as high as 80 percent in some situations. Many knowledgeable IDS integrators advise that alarms be suspended for a period of weeks, even months, as operators gain familiarity with the IDS and target systems. It is especially wise to delay activation of automated responses to attacks until operators and system administrators are familiar with the IDS and have tuned it to the target environment.

27.7.4.4 Outsourcing of IDS/IPS. Any discussion of IDS and IPS strategies would be incomplete without mention of outsourcing these security services. There are significant advantages associated with this approach, especially in enterprises too small to afford an extensive security staff. As in other areas of IT outsourcing, one must have an extremely clear idea of the specific security and operational goals desired

27 · 16 INTRUSION DETECTION AND INTRUSION PREVENTION DEVICES

in order for this approach to be effective. A clearly worded security policy that reflects current concerns is essential. The advantages of outsourcing are many: Managed security service providers usually have considerable experience in dealing with IDS and IPS equipment, their staffs are often well trained and experienced in the use of the equipment, monitoring personnel are usually in attendance around the clock, and contract terms often include specific levels of service agreements.

In the words of a wise CISO, “outsourcing isn’t offloading.” That means that outsourcing security functions does not relieve you of the responsibility for system security. It also means that you must exercise due diligence in selecting the service provider, tasking and managing it, and monitoring it to ensure that your policy goals are being well served by the provider.

27.8 CONCLUSION. Intrusion detection and intrusion prevention are valuable additions to system security suites, allowing security managers to spot, and sometimes block, those security violations that inevitably occur despite the placement of preventive security measures. Although current commercial products are imperfect, they serve to recognize many common intrusion types, in many cases quickly enough to allow security personnel and IPSs to block damage to systems and data. Furthermore, as research and development in intrusion detection and prevention continues, the quality and capabilities of available IDSs and IPSs will steadily improve.

27.9 FURTHER READING

- Crosbie, M., and E. H. Spafford. “Defending a Computer System Using Autonomous Agents.” Proceedings of the 18th National Information Systems Security Conference. Baltimore, MD: October 1995.
- Jackson, K. A., D. DuBois, and C. Stallings. “An Expert System Application for Network Intrusion Detection.” Proceedings of the 14th National Computer Security Conference. Washington, DC: October 1991.
- Fadia, A., and M. Zacharia. *Intrusion Alert: An Ethical Hacking Guide to Intrusion Detection*. Vikas Publishing House, 2010.
- Flegel, U. *Privacy-Respecting Intrusion Detection*. New York: Springer, 2007.
- Hämmerli, B., and R. Sommer, eds. *Detection of Intrusions and Malware, and Vulnerability Assessment: 4th International Conference*, DIMVA 2007 Lucerne, Switzerland, July 12–13, 2007 Proceedings. New York: Springer, 2007.
- Kumar, Sandeep, and E. Spafford. “A Pattern Matching Model for Misuse Intrusion Detection.” Proceedings of the 17th National Computer Security Conference. Baltimore, MD: October 1994.
- Lunt, T., et al. “A Real-Time Intrusion Detection Expert System (IDES).” *Computer Science Lab*, Menlo Park, CA: SRI International, May 1990.
- Mukherjee, B., L. T. Heberlein, and K. N. Levitt. “Network Intrusion Detection,” *IEEE Network* 8, No. 3 (May–June 1994).
- Paxson, V. “Bro: A System for Detecting Network Intruders in Real Time.” *7th USENIX Security Symposium*. San Antonio, TX: January 1998.
- Porras, P., and P. Neumann. “EMERALD: Event Monitoring Enabling Responses to Anomalous Live Disturbances.” *Proceedings of 20th National Information System Security Conference*. Baltimore, MD: October 1997.
- Scarfone, K., and P. Mell. *Guide to Intrusion Detection and Prevention Systems (IDPS)* Revision 1 (Draft). NIST Special Publication 800-94.

NOTES 27 · 17

- Gaithersburg, MD: National Institute of Standards and Technology, 2012.
Available: http://csrc.nist.gov/publications/drafts/800-94-rev1/draft_sp800-94-rev1.pdf
- Schaefer, M., et al. "Auditing: A Relevant Contribution to Trusted Database Management Systems." *Proceedings of the 5th Annual Computer Security Applications Conference*. Tucson, AZ: December 1989.
- Shostack, A., and S. Blake. "Towards a Taxonomy of Network Security Assessment Techniques." *Proceedings of 1999 Black Hat Briefings*. Las Vegas, NV: July 1999.
- Trost, R. *Practical Intrusion Analysis: Prevention and Detection for the Twenty-First Century*. Addison-Wesley Professional, 2009

27.10 NOTES

1. See the Telecom Glossary 2000 from the American National Standards Institute, Inc.: www.its.bldrdoc.gov/projects/telecomglossary2000 (URL inactive).
2. J. J. Wassermann, "The Vanishing Trail," *Bell Telephone Magazine* 47, No. 4 (July/August 1968).
3. J. P. Anderson, "Computer Security Technology Planning Study Volume II," ESD-TR-73-51, Electronic Systems Division, Air Force Systems Command, Hanscom Field, Bedford, MA 01730 (October 1972).
4. J. P. Anderson, *Computer Security Threat Monitoring and Surveillance* (Fort Washington, PA: James P. Anderson Co., April 1980).
5. D. Denning, "An Intrusion Detection Model," *Proceedings of the 1986 IEEE Symposium on Security and Privacy* (Washington, DC: IEEE Computer Society Press, 1986).
6. T. Lunt et al., "A Real-Time Intrusion Detection Expert System (IDES)," *Computer Science Lab*, SRI International, May 1990.
7. B. Mukherjee, L. T. Heberlein, and K. N. Levitt, "Network Intrusion Detection," *IEEE Network* 8, No. 3 (May–June 1994).
8. L. T. Heberlein, K. N. Levitt, and B. Mukherjee, "A Network Security Monitor," *Proceedings of the 1990 IEEE Symposium on Research in Security and Privacy*, Oakland, CA: May 1990.
9. S. Snapp et al., "DIDS (Distributed Intrusion Detection System) Motivation, Architecture, and an Early Prototype," *Proceedings of the 14th National Computer Security Conference*, Washington, DC: October 1991.
10. IPS Quadrant, *Information Security Magazine* (July 2004); http://infosecuritymag.techtarget.com/ss/0,295796,sid6_iss426_art880,00.html
11. A. Freedman, *Computer Desktop Encyclopedia*, v21.1, 2008 (www.computerlanguage.com) provides this definition: "Switch fabric—(1) The inter-connect architecture used by a switching device, which redirects the data coming in on one of its ports out to another of its ports. The word 'fabric' comes from the resulting criss-crossed lines when all the inputs on a switch with hundreds of ports are connected to all possible outputs. (2) The combination of interconnected switches used throughout a campus or large geographic area, which collectively provide a routing infrastructure."
12. K. Scarfone and P. Mell, *Guide to Intrusion Detection and Prevention Systems (IDPS)*, Revision 1 (Draft). NIST Special Publication 800-94

27 · 18 INTRUSION DETECTION AND INTRUSION PREVENTION DEVICES

(Gaithersburg, MD: National Institute of Standards and Technology, July 2012);
http://csrc.nist.gov/publications/drafts/800-94-rev1/draft_sp800-94-rev1.pdf, pp. 36–48.

13. The DMZ is a reserved area in some network architectures in which Web servers are often placed, separated from the Internet by one firewall system and separated from the internal corporate network by another firewall.

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER **28**

IDENTIFICATION AND AUTHENTICATION

**Ravi Sandhu, Jennifer Hadley,
Steven Lovaas, and Nicholas Takacs**

28.1	INTRODUCTION	28·2	28.4	TOKEN-BASED AUTHENTICATION	28·13
28.2	FOUR PRINCIPLES OF AUTHENTICATION	28·2		28.4.1 Card Entry Systems	28·13
28.2.1	What Only You Know	28·3		28.4.2 Proximity and Touch Cards	28·14
28.2.2	What Only You Have	28·3		28.4.3 Smart Cards and Dongles	28·15
28.2.3	What Only You Are	28·4		28.4.4 Soft Tokens	28·16
28.2.4	What Only You Do	28·5		28.4.5 One-Time Password Generators	28·17
28.3	PASSWORD-BASED AUTHENTICATION	28·5		28.4.6 Authentication Using Mobile Devices	28·17
28.3.1	Access to User Passwords by System Administrators	28·5	28.5	BIOMETRIC AUTHENTICATION	28·18
28.3.2	Risk of Undetected Theft	28·6			
28.3.3	Risk of Undetected Sharing	28·7	28.6	CROSS-DOMAIN AUTHENTICATION	28·18
28.3.4	Risk of Weakest Link	28·7			
28.3.5	Risk of Online Guessing	28·8	28.7	RELATIVE COSTS OF AUTHENTICATION TECHNOLOGIES	28·18
28.3.6	Risk of Off-Line Dictionary Attacks	28·9			
28.3.7	Risk of Password Replay	28·10	28.8	CONCLUSIONS	28·18
28.3.8	Risk of Server Spoofing	28·11	28.9	SUMMARY	28·19
28.3.9	Risk of Password Reuse	28·12	28.10	FURTHER READING	28·20
28.3.10	Authentication Using Recognition of Symbols	28·13	28.11	NOTES	28·21

28 · 2 IDENTIFICATION AND AUTHENTICATION

28.1 INTRODUCTION. *Authorization* is the allocation of permissions for specific types of access to restricted information. In the real world, authorization is conferred on real human beings; in contrast, information technology normally confers authorization on *user identifiers* (IDs). Computer systems need to link specific IDs to particular authorized users of those IDs. Even inanimate components, such as network interface cards, firewalls, and printers, need IDs. *Identification* is the process of ascribing an ID to a human being or to another computer or network component. *Authentication* is the process of *binding* an ID to a specific entity. For example, authentication of a user's identity generally involves narrowing the range of possible entities claiming to have authorized use of a specific ID down to a single person.

The focus of this chapter is on person-to-computer authentication. In practice, we also need computer-to-person authentication to prevent spoofing of services on a network. This type of authentication is increasingly important, especially on open networks such as the Internet, where users may be misled about the identity of the Websites they visit. For example, some criminals send unsolicited email messages in Hypertext Markup Language (HTML) to victims; the messages include links that are labeled to suggest an inoffensive or well-respected Website, but the underlying HTML actually links to a fraudulent site designed to trick people into revealing personal information, such as credit card numbers or details to support theft of identity. More generally, computer-to-computer mutual authentication, typically in both directions, is essential to safeguard critical transactions such as those of interbank transfers and business-to-business electronic commerce.

In the early decades of computer usage, most computers authenticated users who accessed mainframes from within a single enterprise. User IDs therefore could be assigned in a centralized and controlled manner. Even so, identifiers have never necessarily been unique, for there is no obligatory one-to-one relationship between a user ID and a human being's real-world identity. For example, several people could share an account such as *inventory clerk* without interference from the computer; at most, the operating system might be configured to prevent simultaneous sharing of an ID by limiting the number of sessions initiated with a specific ID to one.

Conversely, a single user often has many user IDs. For example, there may be unique identifiers for each of dozens of Websites for music clubs, book clubs, enterprise email, and so on. Even on the same computer, a given user might have several accounts defined for different purposes; *jane_doe* and *jdoe* might be identifiers for two different application packages on a system. These multiple identifiers cause problems for administrators if they do not know that the same user is associated with the different IDs; they also cause practical problems for users who have to use different authentication methods for a range of IDs. One of the critical goals of today's identification and authentication (I&A) research and development is to develop reliable and economical methods for *single signon*, whereby users would not have to reidentify and reauthenticate themselves when accessing different computer systems linked into an Internet. For details of I&A in facilities security, see Chapter 23 in this *Handbook*.

28.2 FOUR PRINCIPLES OF AUTHENTICATION. Authentication of a claimed identity can be established in four ways:

1. What only you know (passwords and passphrases)
2. What only you have (tokens: physical keys, smart cards)

FOUR PRINCIPLES OF AUTHENTICATION 28 · 3

3. What only you are (static biometrics: fingerprint, face, retina, and iris recognition)
4. What only you do (dynamic biometrics: voice, handwriting, and typing recognition)

In each approach, the assumption is that no one else but the authorized user of an identifier has access to the password or token and that the probability of simulating static or biometric data is acceptably low.

These methods can be combined; for example, passwords often are combined with tokens or biometrics to provide stronger authentication than is possible with either one alone. A familiar example of this *two-factor authentication* occurs with automatic teller machine (ATM) cards. Possession of the card (the token) and knowledge of the personal identification number (the PIN, corresponding to a password) are required to access a user's bank account.

This chapter introduces each of these four authentication methods and provides additional details for each. For discussions of methods of bypassing and subverting identification and authentication techniques, see Chapter 15 in this *Handbook*.

28.2.1 What Only You Know. Password- or passphrase-based authentication is so widely used that any person who has had any contact with computers and networks probably has had several passwords. Although password technology often is poorly administered and insecure (and frustrating) for users and administrators, passwords can be deployed much more securely and conveniently than they usually are. Many security professionals have felt and hoped for years that passwords would eventually be phased out, to be replaced by tokens or biometrics, but the consensus today is that passwords are not likely to disappear soon and that they will continue to be the dominant authentication technique for years to come.

Demonstrating knowledge of a password does not directly authenticate a human being. It simply authenticates knowledge of the password. Unauthorized knowledge of, or guessing at, a password can lead to impersonation of one user by another; this is called *spoofing*. The theft of a password can be difficult to detect since it is not a tangible asset. Passwords are also very easy to share. It is common for senior executives to give their passwords to their assistants to facilitate their work, even though assigning proxy privileges would be as effective and more secure. Students of security—including criminal hackers—quickly learn that poor password management can be widespread unless organizations provide consistent awareness, education, and training; for example, naïve users sometimes place their passwords on sticky notes under their keyboards or even in plain view on the side of their monitors.

28.2.2 What Only You Have. Authentication based on possession of a token is used where higher assurance of identity is desired than is possible by passwords alone. As with passwords, possession of a token does not directly authenticate a human being; rather it authenticates possession of the token and ability to use it. Sometimes a password or PIN is required to use the token, thus establishing two-factor authentication; the theory is that the requirement to have both elements decreases the likelihood of successful spoofing.

Tokens can take on a variety of forms. The oldest token is the physical key for a physical lock, but these are not often used for securing computer systems. *Soft tokens* are carried on transportable media or even accessed over a network from a server. Soft tokens contain only data; they typically require a password to access the contents.

28 · 4 IDENTIFICATION AND AUTHENTICATION

Modern tokens are usually implemented in self-contained hardware with computing capability. Examples include:

- Credit card-size devices with a liquid crystal display (LCD) that display pseudo-random numbers or other codes.
- LCD devices in the shape of a key fob using the same algorithms as the credit card-shape devices.
- Hardware devices called *dongles* that plug into input-output ports on computers. Examples include dongles for serial ports, parallel ports, Universal Serial Bus (USB) ports, and PC-card interfaces.

All tokens used for computer authentication require software to process information residing in or produced by them. The most significant distinction is whether the tokens require electronic contact with the authentication system. *Contactless* tokens are easier to deploy because they do not require specialized readers. For example, the credit card and key fob pseudorandom number generators simply require the user to enter the visible code in response to a prompt from the authentication software. Contactless tokens are more limited in function than *contact tokens*. For instance, a contact token can be used to create digital signatures whereas a contactless token cannot do so practically.

In cyberspace, a token does not authenticate by means of physical characteristics. Rather the token has some secret, either exclusive to itself or possibly shared with a server on the network. Authentication of the token is really authentication of knowledge of the secret stored on the token. As such, authentication based on possession of a token is tantamount to authentication based on what the token knows. However, this secret can be longer and more random than a secret that a user has to retain in human memory, such as a password. Unfortunately, building cost-effective and secure tokens from which the secret cannot be extracted by tampering or by brute-force guesswork has proven much more difficult than initially anticipated. In the early 1990s, many security professionals believed that tokens would replace passwords; in fact, however, although tokens continue to be an attractive authentication technology, they probably will not become pervasive soon because of the consistent (but highly debatable) belief that passwords are less expensive to implement and manage than other methods of authentication (see Section 28.7).

28.2.3 What Only You Are. Biometrics takes authentication directly to the human being. This topic is covered more extensively in Chapter 29 in this *Handbook*, but the basics can be mentioned here.

As humans, we recognize each other by a number of characteristics. Biometric authentication seeks to achieve a similar result in cyberspace. A *static biometric* is a characteristic of a person such as fingerprint, hand geometry, or iris pattern; more dramatically, it could be the DNA of an individual. The likelihood of two individuals having identical fingerprints, iris patterns, or DNA is minuscule (with exceptions for genetically identical siblings). Biometrics requires specialized and expensive readers to capture the biometric data, making widespread deployment difficult.

Biometrics also suffers from the problems of replay and tampering. Thus, the biometric reader must itself be trusted and tamper-proof; this reduces the likelihood of an attacker capturing the data input and replaying it at a later time, or creating false biometric profiles to trick the system into accepting an imposter. Moreover, the biometric

PASSWORD-BASED AUTHENTICATION 28 · 5

data themselves must be captured in proximity to the user to reduce the likelihood of substitution, such as in the case of stolen blood used to fool a DNA-based biometric system. If the data are transmitted to a distant server for authentication, the transmission requires a secure protocol, with extensive provisions for time-stamping and rapid expiration of the data.

28.2.4 What Only You Do. *Dynamic biometrics* captures a dynamic process rather than a static characteristic of a person. A well-known example is that of signature dynamics. Signature dynamics involves recording the speed and acceleration of a person's hand as a signature is written on a special tablet. Rather than merely the shape of the signature, it is the dynamic characteristics of motion while writing the signature that authenticates the person—motions that are extremely hard to simulate. Another possibility is to recognize characteristics of a person's voice as he or she is asked to read aloud some specified text. Keystroke dynamics of a person's typing behavior is another alternative.

As in all other forms of authentication, dynamic biometrics depends on exclusion of capture and playback attacks, in which, for example, a recording of someone's voice might be used to fool a voice-recognition system. Similarly, a signature-dynamics system might be fooled by playback of the data recorded from an authentic signature. Encryption techniques help to make such attacks more difficult.

Security experts agree that biometrics offer a stronger guarantee of authentication, but deployment on a large scale remains to be demonstrated. Whether this technology becomes pervasive ultimately may be determined by its social and political acceptability as much as by improved technology.

28.3 PASSWORD-BASED AUTHENTICATION. Passwords are the pervasive technology for authentication in cyberspace today. At a conservative estimate, there are close to a billion password-based authentications per day. Examples include the vast number of Internet users and the number of passwords each one uses every day. However, the current deployment of password technology needs to be improved in many ways. Today users must remember too many identities and corresponding passwords. Also, the deployed technology is more fragile than it needs to be; for example, many users choose passwords that can be guessed easily. Passwords are never going to be as secure as the strongest biometric systems, so one would not use them as the sole basis for, say, launching nuclear missiles. However, their use can be made strong enough for many less critical transactions.

The next sections review the major risks of password use and their mitigation by technical, social, and procedural means.

28.3.1 Access to User Passwords by System Administrators. One of the most dangerous practices in use today is the storage of unencrypted user passwords accessible to system administrators. In some sites, new users receive passwords that are assigned and written down by system administrators. If these passwords are used only once, for the initial logon, the user can be forced to choose or create a truly secret password that no one else knows. However, in many such sites, administrators keep control of a paper or electronic record, usually for quick access when users forget their own passwords. Such access completely destroys an important element of I&A: *nonrepudiation*. If someone else has access to a password, then authorized users can reasonably *repudiate* transactions by claiming that their identities were spoofed. It is difficult to counter such repudiation, especially in a court of law considering

28 · 6 IDENTIFICATION AND AUTHENTICATION

an accusation of malfeasance by the authorized user of that password. In general, passwords that will be used repeatedly should not be written down, and they should not be accessible to system administrators. Critical passwords can be written down, stored in tamper-proof containers, and locked away where at least two signatures will be required for retrieval in case of emergency.

28.3.2 Risk of Undetected Theft. Perhaps the biggest intrinsic risk with passwords is that they can be stolen without knowledge of the user. Observation of someone typing in a password is sufficient to leak it. This can happen surreptitiously without the victim's explicit knowledge. A related risk is disclosure of a password to an attacker who persuades the legitimate user to reveal it by posing as a systems administrator who needs the password to do something beneficial for the user. Loss of a physical token eventually may be discovered, since it is missing, although the possibility of cloning these devices remains. Loss of a password, however, can be discovered only by detecting its misuse or by finding it in the possession of an unauthorized user (e.g., in a list of passwords cracked by using a dictionary-based *password-cracking* program, as described in Section 28.3.6).

There are several mitigations of this risk. First, user education and awareness are critically important. People need to treat important secrets with the care they deserve. In an unsafe environment, a password should be typed in discreetly. Efforts to be discreet should be positively reinforced while negligence in exposing passwords during entry should be considered akin to bad social behavior.

User education and awareness, although extremely important, can never be the whole solution. People will inevitably slip up and make mistakes. Some of us are more negligent than others. Others will be observed surreptitiously. In some cases passwords will be revealed to computers with Trojan horses (see Chapter 16 in this *Handbook*) that capture them. Technologists must pursue technical and human solutions to mitigate these risks.

Since some losses of control over passwords are inevitable, it logically follows that password-based authentication should be used only in situations where misuse detection is not only feasible but actually convenient to do in real time. To make this possible, the system architecture should centralize the information needed for misuse detection in one place. If the required information is dispersed across many servers, it will be difficult to coordinate the different audit trails. Traditionally, users of password systems have not considered the need for misuse detection. However, modern security is firmly based on a mix of prevention and detection techniques. Security professionals should apply similar thinking to authentication systems. Ease of misuse detection should be an important criterion in the design of any authentication system. For password-based systems, misuse detection capability should be considered an essential requirement. (For information on intrusion-detection systems, see Chapter 27 of this *Handbook*.)

What else can system designers do to mitigate this risk? It should be made easy for users to change their passwords themselves. Having a system administrator change a password that will be used more than once is illogical.

If a user feels that a password may have been compromised, changing it should be a simple matter. In particular, the system should never prevent a user's attempt to change a password. Some deployed systems will deny change of a password if the password was changed recently, say in the past 24 hours. Although there are reasons for this kind of restriction, it may create a bigger risk than the one it purports to prevent.

Users should be encouraged to change their passwords fairly often; a typical allowable lifetime for a password is between 30 and 90 days. Without occasional changes,

PASSWORD-BASED AUTHENTICATION 28 · 7

a compromised password could be held until the malicious attacker finds opportunity to use it. Frequent changes to passwords reduce the window of opportunity for such attackers.

28.3.3 Risk of Undetected Sharing. Another major risk of passwords is the ease with which they can be shared. There are many examples of sharing between executives and their secretaries, between physicians and office staff or nurses, between professors and their secretaries or students, and among coworkers in any activity. User education and strict policies against password and account sharing are obvious first steps to deter this possibility. Strict policies can be effective within an organization, but their deterrent effect may not carry over to large consumer populations. Misuse detection also can be employed to enforce a strict policy.

The root cause of password sharing within an organization is the lack of effective delegation mechanisms whereby selected privileges of one user can be delegated to another. Better authorization mechanisms could eliminate much of the perceived need for password sharing. It should be possible for secretaries to read their bosses' email under their own identities and passwords. In fact, the bosses should be able to segregate the email that the secretaries can read while denying access to more sensitive email. Moreover, reading the boss's email should be possible without allowing the secretary to send email under the boss's identity; a proxy privilege could allow secretaries to answer their boss's email while signing the replies with their own names. In the nonelectronic world, secretaries routinely answer mail for other people without impersonating them, and this should be the practice with computers as well.

Sharing of passwords among consumers is likely to occur when the cost to consumers is minimal. Although consumers are unlikely to share passwords for an online bank or brokerage account with others, they may be willing to share passwords for an online subscription service, possibly with many friends. A dishonest consumer may even make a business of reselling the service. One way to deter such piracy would be to tie exposure of the password to exposure of a sensitive secret of the consumer, such as a credit card number. Few people, criminal or not, would hand over a password that includes their own credit card number.

Another approach that reduces account sharing is one-time passwords that are generated by inexpensive tokens that have recently been distributed to consumers by Web merchants and banks (e.g., Citibank starting in May 2006) as a method for authenticating identity.¹ These tokens generate random passwords that change every minute or so and that can be traced to the specific unit that creates them—and that unit can be tied precisely to the original recipient. Not only does such a system make password sharing virtually impossible, but a shared password can be traced directly to the violator of the terms of use and result in legal action and fines. For more on one-time passwords, see Section 28.4.1.

28.3.4 Risk of Weakest Link. One of the frustrations of passwords is that users have to remember too many. Thus, users tend to repeat selection of the same password at multiple sites. This is a serious risk. Exposure of a user password at a poorly maintained site can lead to penetration of the user's account at numerous other sites. It is not easy to deploy technical measures to protect directly against this risk. A particular site can force a user to pick a complex password or can even choose the password for the user. However, it cannot prevent use of the same password elsewhere. This is one area where user education, awareness, and self-interest are paramount.

28 · 8 IDENTIFICATION AND AUTHENTICATION

Malicious attackers can set up rogue Websites easily, to entice users to register for attractive services, whereupon the user's password for other sites may be revealed.

Password safes are available in which difficult-to-guess passwords can be stored and accessed once a master password is supplied. These tools can be helpful, but they increase the importance of creating a suitable master password and of safeguarding it against discovery.

A technical solution to mitigate this problem is to avoid the requirement that a user has to register at multiple sites with user IDs and passwords. Instead, the user should register at a few trusted sites, but the user ID should be usable at multiple sites by secured sharing of assurances that the user has in fact been identified and authenticated sufficiently for business to continue. This is essentially what public key infrastructure (PKI) seeks to do. Once a user has identified him- or herself to a provider of such electronic credentials (client certificates), the certificate becomes a method for authentication. For example, some banks provide a service that allows a user to create a unique number that can substitute for their credit card number in a particular transaction. The user authenticates to the bank site, obtains a unique certificate that substitutes for the credit card number, and gives it to the merchant. The merchant can verify that it generates a valid transaction authorization, but never knows the original credit card number. With authentication based on client certificates, it is not necessary to expose a user's password to multiple sites. An effective marriage of passwords and PKI would reduce the exposure to the weakest link. For more details of PKI, see Chapter 37 in this *Handbook*.

A similar approach stores sensitive information in one place and then directs businesses to that place for payment information. For example, today a number of systems (e.g., PayPal) allow a user to register credit card information once, with a trusted service, and then pay online retailers (e-tailers) via that service.

28.3.5 Risk of Online Guessing. Authentication systems are susceptible to guessing attacks. In *online guessing*, an attacker tries to authenticate using a valid user ID and a guessed password. If the password has been poorly selected, the attacker may get lucky. The attacker also may be able to exploit personal knowledge of the victim to select likely passwords. This approach exploits the documented tendency of naive users to select passwords from lists of obvious words, family, friends, pets, sports, commercial brands, and other easily obtained information. For example, studies of password files consistently show that the most frequently selected password in the world is "password"; the second most frequent is the user ID itself or the user ID backward. An account with the same password as the user ID is often called a *Joe account*, as in User ID: joe; Password: joe.

Another kind of password vulnerable to guessing is a password assigned by default; for example, many software installations create accounts with the same password on all systems. Documentation usually warns users to change those *canonical passwords*, but many people ignore the warning. Canonical passwords are particularly dangerous when they grant access to powerful accounts, such as root accounts or to support functions.

The first line of defense against online attacks is to enforce password complexity rules, in addition to user education and awareness. Many systems today require a minimum of eight-character passwords with a mix of upper- and lower-case letters, numerals, and possibly special characters. Nonetheless, online guessing attacks are still possible, and system logging (see Chapter 53) or application logging (see Chapter 52) can be helpful in identifying successful impersonation. For example, log files may

PASSWORD-BASED AUTHENTICATION 28 · 9

show that a particular user has never logged on to a particular account outside working hours, yet someone has logged on as that user in the middle of the night.

Some systems react to online attacks by a simple rule that locks the account after a certain number of failed attempts. This rule may have been borrowed from a similar rule with ATM cards. The rule actually makes sense in the context of ATMs, with two-factor authentication based on possession of the card and knowledge of the PIN. However, in a password-only scheme, the “three strikes and out” rule can lead to denial of service to legitimate users. An attacker can easily lock up many accounts by entering three wrong passwords repeatedly, as discussed in Chapter 18 about denial of service in this *Handbook*. A more graceful rule slows down the rate at which password guessing can be attempted, so that a legitimate user may be perceptibly slowed down in authentication but not denied. For example, locking an account for a couple of minutes after three bad passwords suffices to make brute-force guesswork impractical.

In addition, intrusion-detection systems can be configured to alert system administrators immediately upon repeated entry of bad passwords. Human beings then can intervene to determine the cause of the bad passwords—user error or malfeasance.

28.3.6 Risk of Off-Line Dictionary Attacks. The paramount technical attack on password-based authentication systems is the *dictionary attack*. Such attacks start with copying the password file for a target system and placing it on a computer under the attacker’s control. The password file normally uses *one-way encryption* that allows the system to encrypt an entered password and compare it to the encrypted form of the legitimate password. If the two encrypted strings match, the entered password is presumably correct, and so the system authenticates the user ID.

The dictionary attack is described as off-line because the attacker obtains the necessary information to carry out the attack and then performs computations off-line to discover the password from this information. It is a guessing attack because the attacker tries different likely passwords from an extensive list of possible passwords (the *dictionary*). The list of likely passwords is called a dictionary because it includes words from one or more natural languages, such as English and Spanish; specialized versions used with *password-cracking* programs may sort words by frequency of use rather than alphabetically to speed up successful guesses.

The initial response to dictionary attacks was to stop users from selecting passwords that could be cracked via a dictionary attack. In essence, the system would try a dictionary attack; if it succeeded, it would prohibit the user from selecting this password. This is not a productive approach because attackers’ dictionaries are often ahead of the system’s dictionaries. The productive approach is to prevent the attacker from collecting the information necessary to carry out the dictionary attack.

Designers of password-based authentication systems were slow to recognize the risk of dictionary attacks. It has long been understood that passwords should not be stored on a server in cleartext because this becomes a single point of catastrophic failure. Time-sharing systems of the early 1970s stored passwords in a “hashed” form. Knowledge of the hashed form of a password did not reveal the actual password. Authentication of passwords was achieved by computing the hash from the presented password and comparing with the stored hash. The UNIX system actually made the hashed form of user passwords easily readable, since reversing the hash was correctly considered computationally infeasible. However, knowledge of the hashed form of a password is sufficient for dictionary attacks. The attacker guesses a password from a list, or dictionary, of likely passwords, computes its hash, and compares it with the stored hash. If they match, the attacker’s guess is verified; otherwise the attacker tries

28 · 10 IDENTIFICATION AND AUTHENTICATION

another guess. Since the late 1980s, UNIX systems have stopped making the file of hashed passwords easy to read, so this vulnerability has been reduced.

UNIX also introduced the concept of a *salt* to make dictionary attacks more difficult. The user password and a random number called the salt are hashed together and stored on the server. The salt itself is also stored on the server. To authenticate a user, the presented password and the stored salt are hashed and compared with the stored hash value. Use of a salt means that a separate dictionary attack is required for every user, since each password guess must be hashed along with the salt. Otherwise, the same attack could be run simultaneously against multiple presented passwords.

28.3.7 Risk of Password Replay. If a password is transmitted in cleartext from client to server, it is susceptible to being picked up on the network by an intruder. This is called *password sniffing*. Many systems require the password to be sent to the server in cleartext. Others require transmission of a hash of the password (usually without a salt). Transmitting the hash of a password is risky for two reasons:

1. The hash is sufficient for a dictionary attack unless a salt is used and kept secret.
2. The attacker does not even need to recover the password. Instead, the attacker can replay the hash of the password when needed.

Many existing systems are susceptible to sniffing of passwords on the network in cleartext or hashed form. Fortunately, technical solutions to this problem do exist.

One approach to the replay threat is to use the server's public key to encrypt any transmission of password-related information to the server: Thus, only the server can decrypt the information by using its private key. This is essentially what server-side Secure Shell (SSH) and Secure Sockets Layer (SSL) do. The server-side mode of both SSL and SSH require the server to have a public key certificate. The client-side mode of these protocols requires that the client also have a public key certificate. This approach can be effective but has its own risks.

An alternate approach is to avoid transmitting the password but instead to employ a protocol that requires knowledge of the password to run successfully. One of the earliest and best-known systems to take this approach is Kerberos. In this system, a user's password is converted to a secret key on the client machine and also stored on the Kerberos server. When the user requests authentication, the Kerberos server sends the user's machine a secret session key encrypted using the shared secret key derived from the user's password. The ability to decrypt this message correctly demonstrates knowledge of the password without actually transmitting it, in cleartext, hashed, or encrypted form. Unfortunately, the Kerberos protocol is susceptible to dictionary attacks; any client machine can pretend to be any user and can obtain the necessary information required for a dictionary attack.

Kerberos also does not use a salt, so the same dictionary attack can be applied to multiple users at one time. Kerberos Version 5 provides for a preauthentication option, which makes it somewhat harder to gather the information for a dictionary attack. The data are no longer available by simply asking the Kerberos server for them; instead they must be snuffed on a network. Recent experiments have shown that dictionary attacks on Kerberos are very practical, so this is a serious vulnerability of a widely deployed password-based authentication system.

Since the early 1990s, many password-based authentication protocols have been published that do not suffer from the dictionary attacks to which Kerberos is so

PASSWORD-BASED AUTHENTICATION 28 · 11

vulnerable. In particular, *zero-knowledge password proofs* are based on the idea that two parties (computers and people) can demonstrate that they know a secret password without revealing the password. These methods depend on the ability to establish that the two parties both independently selected the same number—but without knowing what the specific number is.² One popular conceptual model of this process is the zero-knowledge password proof, which runs as follows:

1. Two people want to test whether they share a secret number (in this thought experiment, a single digit between 1 and 10). In this example, the shared number is 3.
2. The two people have a deck of 10 blank cards.
3. The first person counts down to the third card and makes a mark on the right edge of that card.
4. The deck of cards is arranged so that the second person can mark the left edge of the cards but cannot see the right edge.
5. The second person also counts down to the third card (in this example) and marks the left edge.
6. The card deck is shuffled so that the sequence order is lost and then displayed to both parties.
7. If a single card has a mark on both the right edge and the left edge, then the two parties share the secret number, but neither had to reveal exactly which number it was.³

It will be interesting to see if this approach to authentication can be implemented on actual computers and if it is commercially used on a significant scale in the coming years.

28.3.8 Risk of Server Spoofing. As mentioned earlier, one widely used approach to preventing password exposure in transit on a network is to send passwords from client to server encrypted using the server’s public key. The server, which has the corresponding private key, can decrypt the password to recover it. Knowledge of the public key is not sufficient for a spoomer to determine the private key. Naïve protocols for protecting the private key can be susceptible to replay attacks. However, there are two well-designed protocols in widespread use today.

Server-side SSL is the protocol that has been used by most Web surfers. In this protocol, the server’s public key is used to secure transmission of the user’s password from client to server. Like all public key-based schemes, the Achilles’ heel of this protocol lies in authentic knowledge of the server’s public key. The technology of public key certificates seeks to provide public keys with good assurance of the identity of the server to which they belong. A full discussion of issues with Public Key Infrastructure technology appears in Chapter 37, but it suffices to observe that there are pitfalls with the use of certificates for authentication of servers. A rogue server can collect a user’s password by pretending to be something other than what it is. Relying on the look and feel of a Web page for server authentication is hardly sufficient, since it is easy to copy an entire Web page for use as a decoy and to establish confidence before capturing confidential information. An improvement on Website authentication is to associate a specific image and identifying strings (e.g., a picture of a hippopotamus labeled “Archie’s Favorite Critter”) with a user ID; the Website authenticates itself to

28 · 12 IDENTIFICATION AND AUTHENTICATION

a specific user (e.g., Archie) by displaying that user's chosen image.⁴ Authenticity of the server's certificate can be spoofed in many ways that are hard for the user to detect, and manipulation of the trusted root certificates that are configured in the user's Web browser is possible. Moreover, while trust ultimately chains up to a root certificate, the owner of a single certificate below a trusted root is capable of considerable mischief. Server-side SSH is a similar protocol, typically used to provide secure remote access to UNIX servers. Server-side SSL and server-side SSH share the same fundamental vulnerabilities: Both can be spoofed by certificate manipulation. The use of server-side SSL to protect transmission of passwords from client to server is prevalent on the Internet today, but it is important for customers of authentication products to understand the risks inherent in this approach.

In the client-side mode of these protocols, there is no need for a password to be transmitted from client to server, since client-to-server authentication is based on the client's use of its own private key to generate a digital signature. Hence, client-side protocols are not vulnerable to password capture by server spoofing.

28.3.9 Risk of Password Reuse. The need to change passwords with some reasonable frequency is well recognized, but what is reasonable frequency? And how draconian should the enforcement be? It seems that security administrators have pushed too far on these questions. Forcing users to change passwords every month, and enforcing such rules ruthlessly, actually could lead to less security rather than more because so many frustrated users write down and store their ever-changing passwords in nonsecure places. There is a real risk here, created by well-meaning security administrators who have made the problem worse than it inherently is.

Systems that choose a password for the user have their own set of problems and are generally too user unfriendly to be viable in the Internet age. This discussion focuses on systems that allow users to select their own passwords.

How does exposure of a password increase with time? Even the strongest password-based system, with immunity to off-line dictionary attacks and password capture by server spoofing, faces increased exposure as time passes. Over a long period of time, a slow, ongoing, online guessing attack could be successful. Also, the likelihood of inadvertent disclosure by surreptitious observation, or exposure on a Trojan horse-infected computer, increases with time. Nevertheless, a good password, carefully chosen by the user to be safe from dictionary attacks and well memorized by the user, should not be changed casually. A change every six months may be appropriate for well-chosen, brute-force-resistant passwords.

Enforcing password changes is a complicated business, and one where the security community has not really done a good job. It is not difficult to keep track of the age of a password and to force a change when an appropriate time has passed. The difficulty is in forcing the new password to be independent of the old one. In fact, the likelihood is that the new password will be a slight variation of the old one. For example, appending the numeral designating a month to a fixed string enables users to have almost the same password, even if they are forced to change it every month. Some systems will keep a history of recently used passwords to prevent their reuse. A system that keeps a history of, say, five passwords can be fooled by rapidly changing the password six times. To prevent this, there are systems that will not allow a user to change password more than once a day. This has the unfortunate effect of actually increasing risk of password exposure, since a user who realizes that the current password may have been inadvertently exposed cannot change it, exactly when the need to do so is greatest.

TOKEN-BASED AUTHENTICATION 28 · 13

28.3.10 Authentication Using Recognition of Symbols. An interesting new approach to user authentication is recognition of particular faces from among a large selection of random photographs. Passfaces software works in this way:

[A] user sets up an array of photographs and puts some familiar ones into the pool to use as keys—the faces of people the user recognizes; then the software can produce a 3×3 grid of random selections including one of the key pictures. The user picks out the familiar picture and then repeats the exercise twice more with new sets of eight strangers and one friend to authenticate the user.⁵

A white paper explains how human beings are particularly good at recognizing faces; indeed, it seems that we have special circuits that have evolved for rapid and accurate perception of faces. According to the paper, advantages of “using Passfaces over passwords” are that Passfaces:

- Can’t be written down or copied
- Can’t be given to another person
- Can’t be guessed
- Involve cognitive not memory skills
- Can be used as a single or part of a dual form of authentication

The power of the system is enhanced by setting parameters to interfere with misuse of the faces. For example:

In some high-security applications the grids of faces may be displayed only for a very short time. A half second is long enough for practiced users to recognize their Passfaces. Combined with masking (faces in a grid are overwritten with a common mask face) it is extremely difficult for “shoulder surfers” to learn the Passfaces as the user clicks on them.

28.4 TOKEN-BASED AUTHENTICATION. Token-based authentication relies on something that the user possesses that no other user of the identifier is supposed to possess or be able to access. This authentication can be achieved in many ways, including:

- Access cards
- Proximity and Touch Cards
- Smart cards and dongles
- One-time password generators
- Soft tokens

For details of authentication in physical-security systems, see Chapter 23 in this *Handbook*.

28.4.1 Card Entry Systems.⁶ The means of encoding identification data on the cards include optical bar code, magnetic stripe, smart cards with embedded chips that store biometric data, and cards with embedded bits of metal. Most bar codes are not secure; the cards are easily duplicated. Although the newer, two-dimensional bar codes are nearly tamperproof, they cannot store much information. Magnetic-striped cards also have many drawbacks: They cannot store much data, and they are easily

28 · 14 IDENTIFICATION AND AUTHENTICATION

altered, copied, or erased (often by accident). In time, the magnetic data decays and must be reprogrammed. Magnetic card readers are not practical outdoors, because of weather and vandalism. Heavily used magnetic card readers and the cards themselves wear out quickly.

Cards with embedded metal bits are effective. The encoding cannot be seen except by X-ray, and it is durable and permanent. The cards must be held against, or inserted into, a reader that scans the card for the position of each bit. They hold very limited data, the coding is factory installed and cannot be changed, and spare cards must be inventoried. In addition, with the right equipment, these cards can be copied.

Methods of using entry cards include proximity, touch, insertion into a slot that returns the card when it is authenticated, or swiping the card through a narrow channel or in and out of a slot. Swiping is fast, but the card must be hand held. Inserting a card has the same shortcoming as swiping and is slower. Wear, weather, or vandalism can damage the card and card reader.

28.4.2 Proximity and Touch Cards.⁷ The best of the new card access control systems use proximity or touch cards. These cards communicate with readers using infrared or microwave transmissions. The reading device powers some types of cards, while others contain miniature batteries. Physically, the cards and card readers are weatherproof, vandal resistant, and do not wear out. Proximity card readers can be surface mounted, recessed flush into a wall, or entirely concealed within a partition so that they do not call attention to a security door.

Touch cards are functionally similar to proximity cards, but they must be held briefly against a reader that is usually visible. Touch cards cost a little less than proximity cards and are good only for entrances with little traffic. The touch-card system is slower, and the cards more easily lost, stolen, or forcefully taken.

Proximity cards (which include RFID cards as well) may be used while concealed inside a pocket, handbag, or wallet. Some are worn concealed or hung on neck lanyards. A proximity card that is also an ID badge that everyone in the workplace wears at all times can access both doors and workstations without being touched. Temporary badges customarily are issued to all visitors, even when escorted, and can be used for access control and to monitor areas entered. Temporary badges are quickly activated with specific privileges and can be revoked automatically and immediately when necessary. Increasingly, the visitor cards are created quickly on site with the visitor's picture and perhaps biometric data as well printed on the card. Longer-term temporary badges can be issued to vendors, contractors, and external employees, although it is best that security personnel store visitor badges safely while the person is off the premises.

The new systems provide many useful functions. They are usually laminated and sealed to prevent wear, damage, or alteration. Individual cards are quickly prepared, activated, and canceled—all on site. The system can restrict entry to specific places, days, and times, and holiday restrictions also can be programmed. Any card can be locked out immediately if lost or stolen, or when the owner leaves.

The newer 13.56-MHz proximity cards function up to three feet away from the card reader; older cards were limited to a range of about four inches. The newer cards are also faster, hold more data, and offer more functionality. Many of the card readers also can write data to the card. There is a trade-off, however, between useful operating range and the amount of data stored. The farther the range, the less the data stored. Most proximity systems are adjustable to optimize distance and speed. For example, on outer perimeter doors, where quick, convenient access is more important than tight security, the systems are set for maximum range. Inner doors that need higher security

TOKEN-BASED AUTHENTICATION 28 · 15

are adjusted to utilize more information and to function at a shorter distance, which is still far greater than the older systems allowed. It is not easy without very high-tech equipment, but proximity and touch cards can be compromised.

There are also self-expiring visitor badges that noticeably change color or prominently display the word “expired” after an elapsed period. Self-expiring badges are reusable and come with a fixed expiration period that is usually from two to 24 hours following each activation. These badges cannot be reactivated, except with very sophisticated equipment.

Cards are not the only proximity or touch devices. Keys or patches also are used. The keys can be small, rugged, and easily attached to a key ring or to a small handheld wand that a security guard might use. The patches work in place of touch cards, or with separate access control systems, to upgrade existing legacy systems. The patches are about the size and thickness of a quarter and are easily attached to anything a person normally carries, such as an ID card or badge, a pager or cell phone, or the inside of a wallet. The newer RFID devices will be even smaller.

Card access often is used for all equipment rooms containing servers, network components, or telephone gear; for off-hour access to information systems by users, technicians, and administrators; and for any areas where high-value items are stored. Card systems usually are integrated with premises security to control access to and egress from the building, elevators, service areas, parking, and restrooms, and other parts of the information infrastructure that can also take advantage of the access controls. Plan ahead and consider where additional access control points may someday be needed. Piecemeal additions at a later date can be costly.

Each entry into a controlled area should be logged in a way that cannot be compromised. Logs should provide an accurate audit trail of everyone who sought entry, when, and whether access was denied. Where stronger security is needed, each egress should be logged in the same way. The logging system is best monitored by software that can review all system data in real time, flag trouble quickly, issue periodic summary reports, and quickly search and summarize unusual events. Reviewing logs manually is a cumbersome, time-consuming task. If only manual auditing is possible, there must be a firm policy to do this every few days.

Token-based entry systems by themselves do not provide strong protection. Therefore, some degree of authentication is required.

28.4.3 Smart Cards and Dongles. Another form of token is a smart card. These cards can go into a PC card reader or can be read by a specialized reader. *Dongles* are smart cards that fit into input-output ports such as Universal Serial Bus (USB). A smart card has its own processing capability and typically stores a private key associated with the user. Often, a password or PIN is required to access the card, thereby providing two-factor authentication capability. The smart card enables user authentication by signing some challenge presented to it with the user’s private key. The signature is verified by means of the user’s public key. A complete discussion of such smart cards involves consideration of public key cryptography, public key certificates or so-called digital certificates, and supporting infrastructure or PKI (see Chapter 37 in this *Handbook*). Suffice it to say that smart cards have long been considered essential for widespread use of public key technology but so far have not been widely deployed.

Hardware tokens offer the potential for stronger authentication than passwords but have seen only limited use due to their perceived costs and their infrastructure requirements. Whether they can be deployed in a scale of millions of users remains to be seen. Authentication by tokens is really authentication by something that the

28 · 16 IDENTIFICATION AND AUTHENTICATION

token knows. Since tokens can be programmed to remember and use secrets much more effectively than humans can, they offer the potential of strong cryptographic authentication, but tamper-proof tokens are not easy to produce. In recent years, attacks based on differential power analysis have proven effective in determining the secret keys and PINs stored on smart cards. These attacks require physical access to a card whose loss probably would be known, so although they may not always be feasible, they certainly call into question the presumed tamper-proof nature of smart cards. As smart cards are more widely deployed, other ingenious attacks are likely to be pursued. Smart cards are more susceptible to secret extraction than tokens because their computations leak such information in the form of electromagnetic radiation.

In comparing tokens with passwords, one can argue that undetected theft is easier with passwords. The token is a physical object whose absence is noticeable. However, tokens create their own problems. Like password generators, they are vulnerable to physical damage and compromise. Smart cards typically are built from thin plastic with a chip embedded in them, making the entire unit susceptible to failure from damage. In addition, users typically may store cards with other credit cards in a wallet or pocket, either of which presents additional environmental hazards. However, for certain industries, the cost of replacing lost or damaged tokens may be worth the reduction in undetected sharing, since a token can be used by only one person at a time.

28.4.4 Soft Tokens. The idea of *soft tokens*, or *software tokens*, has been proposed as a low-cost alternative to hardware tokens. Early soft tokens consisted of a user's private key encrypted with a password and stored on some transportable medium, such as a floppy disk. Such a scheme is extremely vulnerable to dictionary attacks because a guessed password can be verified easily (by testing a putative private key to see if it decrypts a message encrypted using the user's known public key). Moreover, the physical transport of floppy disks and the possible lack of floppy disk drives have led people to store these soft tokens on network servers so they are accessible as needed. Unfortunately, this location also makes them easily accessible to attackers. Protecting access to soft tokens on a network server by means of a password simply returns to the problems of password-based authentication.

It has been suggested that a user's public key could be kept secret and known only to trusted servers to avoid dictionary attacks on the encrypted private key. This approach comes at the severe cost of a closed PKI rather than an open PKI. Schemes for retrieving the private key by means of a secure password-based protocol have been published and are being implemented by some vendors. These schemes ultimately revert to password-based authentication as their foundation. Schemes based on splitting the private key into two parts have been developed. One part of the private key is computed from the password; the other part is stored on an online server, which functions as a network-based virtual smart card. Both parts of the private key are needed for user authentication but are never brought together in one place.

An alternative scheme would be to store the user's entire private key on an online server and make its use contingent on a secure password-based protocol. This approach allows the server to impersonate any user at will and may not be suitable in all environments. However, in all cases, the security of the token relies on the integrity of the computer that uses it. Newer implementations of soft tokens rely on asymmetric cryptography to eliminate the security concerns with storing private keys in a file or other transportable location. The soft token can generate its own key pair and exchange public keys with the authentication server. Although this increases the security of the soft tokens, the concept is inherently weak and prone to attack.

TOKEN-BASED AUTHENTICATION 28 · 17

The U.S. Federal Government established a standard for the identity cards it has mandated for all federal employees and contractors. The original Homeland Security Presidential Directive 12 (HSPD 12), dated August 27, 2004, was entitled “Policy for a Common Identification Standard for Federal Employees and Contractors.” It defined these requirements for secure and reliable identification that—

- Is issued based on sound criteria for verifying an individual employee’s identity
- Is strongly resistant to identity fraud, tampering, counterfeiting, and terrorist exploitation
- Can be rapidly authenticated electronically
- Is issued only by providers whose reliability has been established by an official accreditation process⁸

28.4.5 One-Time Password Generators. A popular form of token, from vendors such as RSA Data Security Inc. and CryptoCard, displays a one-time password, typically a six- or eight-digit numeral, which changes each time an access button is pushed or when a given time has elapsed since the password was last used. The user authenticates by entering the user ID and current value displayed by the token.

The password is called one-time because it expires at the end of its allowable period for use. The token is typically contactless, in that it does not need electrical contact with the computer where the user is presenting authentication data. The user transfers the necessary information from the token via a keyboard or other input device. To make this a two-factor authentication, a fixed user password is also required in addition to the changing one-time password displayed by the token. These tokens are based on shared secret keys, so both the token and the server have a shared secret. The server and the token need to be initialized and then kept synchronized for this scheme to work. In the case of RSA’s SecurID, if the time discrepancy between a specific token and the authenticating system exceeds a specified limit, the authenticating software adjusts a value in an internal table to compensate for the time slippage. Vendors such as CryptoCard have developed event-based authentication algorithms to solve the “slippage” problem. These tokens use a seed with the algorithm to generate a unique value for each button push or other activating action. After the next login, that value is now “known” to the authenticating system, and a new value must be provided. The token’s value increments based on the seed, without requiring synchronization with the authenticating system.

Password generators must be protected against physical tampering. These devices typically include several measures to cause destruction of the electronic circuits if the outer case is opened; for example, in addition to epoxy-resin glue, tokens may include light-sensitive components that are destroyed immediately by exposure to light, rendering the unit unusable. Some password generators and smart cards cannot even be opened to replace batteries, so the entire card must be replaced on a predictable schedule. Tokens of this kind are available that are guaranteed to last one year, two years, or three years without having the batteries wear out.

28.4.6 Authentication Using Mobile Devices. Another form of token-based authentication increasingly available for access to restricted Websites such as bank account pages registers a mobile device such as a cell phone or a tablet that is capable of receiving Short Message Service (SMS) text messages. During login, the authorized user receives a one-time code to enter in the Web page. This approach puts

28 · 18 IDENTIFICATION AND AUTHENTICATION

additional value on security measures to protect the mobile device; for example, short timeouts requiring entry of a PIN, encryption of data in the device's memory, and methods for remotely inactivating or wiping the device if it is lost or stolen.

28.5 BIOMETRIC AUTHENTICATION. Biometric authentication looks like an excellent solution to the problem of authentication in cyberspace. However, there are several challenges in implementing biometrics, including drawbacks in accuracy (false-positive and false-negative results), loss of biometric identifiers, security of templates, and privacy concerns. For a full treatment of biometric authentication, refer to Chapter 29 in this *Handbook*.

28.6 CROSS-DOMAIN AUTHENTICATION. As more systems and processes become Internet enabled, people come to expect a seamless experience between organizations and applications across the Internet. Furthermore, the ever-increasing risk of identity theft has made individuals and organizations more careful about sending too much identity information across an untrusted network. Over the past several years, efforts have increased to enable easy but secure sharing of authentication and authorization information between organizations. Using the Security Assertion Markup Language (SAML), researchers are developing methods enabling one organization to share just enough information between systems to enable a transaction, without compromising privacy. Shibboleth, a project started in 2000 as an Internet2 middleware initiative, is an open-source project using SAML.⁹ Organizations using Shibboleth as a basis for designing new federated authentication and authorization implementations include InCommon¹⁰ and the UK Access Management Federation for Education and Research.¹¹ To facilitate this kind of information sharing, the participating organizations need to share details about their security policies and procedures, so each may decide in advance whether it will trust the authentication assertions of the other. This kind of information sharing is accomplished most effectively via the policies and practice statements used in a PKI. For details on PKI and certificate requirements for cross-domain authentication, refer to Chapter 37 in this *Handbook*.

28.7 RELATIVE COSTS OF AUTHENTICATION TECHNOLOGIES. One of the frequent responses from security professionals in discussions of identification and authentication using anything other than passwords is that the expense of buying new equipment is prohibitive. However, passwords are not free. In an analysis of the costs of managing passwords, RSA Data Security, makers of the SecurID token, estimated deployment costs (initializing user accounts) at about \$12 per user over three years and management costs (replacing forgotten passwords and resetting locked accounts) at about \$660 per user over three years. Worse yet, the well-established failure of most users to select strong passwords (i.e., those resistant to guessing, dictionary attacks, and brute-force attacks) makes passwords a weak authentication mode in practice. In comparison, token-based and biometric authentication are more readily affordable and more effective than passwords.¹²

28.8 CONCLUSIONS. Identification and authentication are the foundations for almost all other security objectives in cyberspace. In the past, these problems often were viewed as simple ones whose solution was assumed, before the real security issues came into play. Important standards for computer security were published and practiced without much attention to identification and authentication. In cyberspace,

SUMMARY 28 · 19

the problems associated with I&A are severe, and they have barely begun to be solved effectively. A robust, scalable identification and authentication infrastructure is vital to achieving security. Technologies such as tokens and biometrics hold out considerable promise, but their deployment requires infrastructure costs dominated by the cost of hardware for readers and the like. Meanwhile, passwords continue to be strengthened, as we better understand the real risks in using them and as we develop technical means to mitigate those risks. The fact that modern operating systems continue to provide simple password-based authentication, vulnerable to dictionary attacks, reflects poorly on the pace at which security technology is adopted in the marketplace. Looking to the future, one can predict that we will see a mix of passwords, biometrics, and tokens in use, perhaps in two- or three-factor configurations. Biometrics and tokens are likely to dominate the high-assurance end, while passwords will dominate the lower end.

One of the misconceptions that security specialists should seek to dispel is that identification and authentication are *sufficient* for improving public safety. However, assigning a reliable and nonrepudiable identity to someone is in no way equivalent to asserting the trustworthiness of that individual. In closed populations such as employee pools, employers can check the background of potential employees and monitor the performance of existing employees; under those circumstances, knowing someone's identity at the entrance gate may indeed improve security. However, when dealing with a large number of unscreened people, such as potential air passengers, confidently being able to name them tells us nothing about their trustworthiness. Having a clerk at a government office glance at fuel-oil invoices and a birth certificate before granting someone a photo ID is no basis for assuming that the carrier of the valid ID is an inoffensive traveler. As several writers have noted, unambiguously knowing the name of the suicide bomber sitting next to you in a plane is not a reasonable basis for complacency. Timothy McVeigh, the Oklahoma City bomber, was a perfectly identifiable citizen of the United States, but he committed his atrocity nonetheless.¹³

Since September 11, 2001, the air transport industry has been a very public example of both the difficulty of strong identification and authentication and the use of identification and authentication as a public relations substitute for substantive security involving thorough passenger screening. The inherent difficulties of authentication are becoming evident to the lay public and to political and corporate leaders. It is very difficult, perhaps even impossible, to guarantee foolproof identification and authentication in free societies. As technologists, we realize that absolute guarantees cannot be achieved in cyberspace. Too many security professionals seek absolute goals, and too many security technologies are marketed as being stronger than they really are. Our profession will benefit greatly if we address practical problems with practical cost-effective techniques and develop a sound security discipline that contains, bounds, and mitigates the inevitable residual risk that we must face in any large-scale human situation.

28.9 SUMMARY. Passwords are widely used in practice and will continue to be a dominant form of user authentication. There are many risks in deploying passwords, and a number of widely used password systems have serious vulnerabilities. Nonetheless, technical measures can mitigate the inherent vulnerabilities of passwords. Although it takes great skill and care, with our current understanding it is technically possible to build and deploy strong password-based authentication systems using commercial products. The truly inherent risks of undetected theft and undetected sharing can be largely mitigated by new technologies, such as intrusion detection systems. Undetected

28 · 20 IDENTIFICATION AND AUTHENTICATION

sharing may be deterred further by a system that couples high-value secret data, such as credit card account numbers, with passwords. Tokens are available to generate one-time passwords or to communicate directly with authentication systems. Although costs have been dropping, tokens are still not as widely deployed as early predictions suggested they would be. Biometric authentication has been implemented only infrequently and on a small scale but offers great potential, especially for high-security applications. Interesting new research and applications are extending the use of authentication (and authorization) over untrusted networks between federated organizations.

28.10 FURTHER READING

- Anderson, J., and R. Vaughn. *A Guide to Understanding Identification and Authentication in Trusted Systems* (1991). “Light Blue Book” in the Rainbow Series, NCSC-TG-017; www.fas.org/irp/nsa/rainbow/tg017.htm
- Birch, D. G. W., ed. *Digital Identity Management: Technological, Business and Social Implications*. Aldershot, Hampshire: Gower Publishing, 2007
- Integrity Sciences. “Bizcard ZKPP: A Zero-Knowledge Password Proof with Pencil and Paper” (2001). www.integritysciences.com/zkppcard.html
- Jablon, D., et al. “Publications on Strong Password Authentication” (2001). www.integritysciences.com/links.html
- Jain, L. C., et al., eds. *Intelligent Biometric Techniques in Fingerprint and Face Recognition*. Boca Raton, FL: CRC Press, 1999
- Kabay, M. E. “Identification, Authentication and Authorization on the World Wide Web,” 1998. www.mekabay.com/infosecmgmt/iaawww.pdf
- National Institute of Standards and Technology (NIST). Personal Identity Verification (PIV) of Federal Employees and Contractors, REVISED DRAFT, July 2012. csrc.nist.gov/publications/drafts/fips201-2/draft_nist-fips-201-2.pdf
- “One-Time Passwords.” FreeBSD Handbook, Chapter 14.5. www.freebsd.org/doc/en_US.ISO8859-1/books/handbook/one-time-passwords.html
- Ouf, M. O., R A. El-kammar, and A. E. S. Ahmed. *Implementation and Analysis Requirements of Ultra-Lightweight Radio Frequency Identification Authentication Protocols*. Lap Lambert Academic Publishing, 2012
- Radhakrishnan, R. *Identity & Security: A Common Architecture & Framework for SOA and Network Convergence*. London: futuretext, 2007
- Smith, R. E. *Authentication: From Passwords to Public Keys*. Reading, MA: Addison-Wesley, 2001
- Todorov, D. *Mechanics of User Identification and Authentication: Fundamentals of Identity Management*. Boca Raton, FL: Auerbach, 2007
- Tung, B. *Kerberos: A Network Authentication System*. Reading, MA: Addison-Wesley, 1999
- Vance, J., “Beyond Passwords: 5 New Ways to Authenticate Users,” *Network World*, May 30, 2007. www.networkworld.com/research/2007/060407-multifactor-authentication.html
- Wayman, J. L. “Biometric Technology: Testing, Evaluation, Results” (1999); www.engr.sjsu.edu/biometrics/publications_technology.html
- Windley, P. *Digital Identity*. Sebastopol, CA: O'Reilly, 2005
- Wu, T. “A Real-World Analysis of Kerberos Password Security.” Proceedings of the 1999 Network and Distributed System Security Symposium; www.isoc.org/isoc/conferences/ndss/99/proceedings/papers/wu.pdf

NOTES 28 · 21

28.11 NOTES

1. “No Token Resistance: Citi rolls out Digipass authentication devices to biz clients,” *Bank Systems & Technology*, May 25, 2006, www.banktech.com/features/showArticle.jhtml?articleID=188103117
2. J.-J. Quisquater, L. Guillou, and T. Berson, “How to Explain Zero-Knowledge Protocols to Your Children,” Proceedings of CRYPTO ‘89, *Advances in Cryptology*, vol. 435 (1989) pp. 628–631. www.cs.wisc.edu/~mkowalcz/628.pdf
3. R. Wright, “Secret Communication Using a Deck of Cards.” Abstract of presentation from DIMACS Research and Education Institute Cryptography and Network Security, July 28–August 15, 1997 (Abstracts of Talks Presented); <ftp://dimacs.rutgers.edu/pub/dimacs/TechnicalReports/TechReports/1997/97-80.ps.gz>
4. M. E. Kabay, “SiteKey Tries to Counter Phishing,” *Network World Security Strategies*, April 3, 2007. www.networkworld.com/newsletters/sec/2007/0402sec1.html
5. M. E. Kabay, “Password Management: Facing the Problem,” *Network World Security Strategies*, October 11, 2007. www.networkworld.com/newsletters/sec/2007/1008sec2.html
6. This section (28.4.1) was originally part of Chapter 23 on “Protecting the Physical Information Infrastructure,” by Franklin Platt, from the 5th edition of this *Handbook*.
7. This section (28.4.2) was originally part of Chapter 23 on “Protecting the Physical Information Infrastructure,” by Franklin Platt, from the 5th edition of this *Handbook*.
8. FIPS 201-1, “Personal Identity Verification (PIV) of Federal Employees and Contractors.” Gaithersburg, MD: NIST, 2006, p. iv, <http://csrc.nist.gov/publications/fips/fips201-1/FIPS-201-1-chng1.pdf>
9. Shibboleth, <http://shibboleth.net>
10. InCommon, www.incommonfederation.org
11. UK Access Management Federation for Education and Research, <http://ukfederation.org/>
12. RSA Data Security, “Are Passwords Really Free? A Closer Look at the Hidden Costs of Password Security,” 2006. www.rsa.com/go/wpt/wpindex.asp?WPID=3384
13. M. E. Kabay, “Airport Safety,” 2005. www.mekabay.com/opinion/airport_safety.pdf

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 29

BIOMETRIC AUTHENTICATION

Eric Salveggio, Steven Lovaas, David R. Lease, and Robert Guess

29.1 INTRODUCTION	29·2		
29.2 IMPORTANCE OF IDENTIFICATION AND VERIFICATION	29·2	29.6.1 General Considerations	29·17
		29.6.2 Health and Disability Considerations	29·18
		29.6.3 Environmental and Cultural Considerations	29·18
29.3 FUNDAMENTALS AND APPLICATIONS	29·2	29.6.4 Cost Considerations	29·19
29.3.1 Overview and History	29·2	29.6.5 Attacks on Biometric Systems	29·19
29.3.2 Properties of Biometrics	29·4	29.6.6 Privacy Concerns	29·20
29.3.3 Identification, Verification, and Authentication	29·5	29.6.7 Legal Issues	29·22
29.3.4 Application Areas	29·7		
29.3.5 Data Acquisition and Presentation	29·8	29.7 RECENT TRENDS IN BIOMETRIC AUTHENTICATION	29·22
		29.7.1 Government Advances in Biometric Authentication	29·22
29.4 TYPES OF BIOMETRIC TECHNOLOGIES	29·9	29.7.2 Face Scanning at Airports and Casinos	29·22
29.4.1 Finger Scan	29·9	29.7.3 Increased Deployment in the Financial Industry	29·23
29.4.2 Facial Scan/Recognition	29·11	29.7.4 Biometrics in the Healthcare Industry	29·23
29.4.3 Hand Geometry Scan	29·12	29.7.5 Increased Deployment of Time and Attendance Systems	29·24
29.4.4 Iris Scan	29·13		
29.4.5 Voice Recognition	29·14		
29.4.6 Other Biometric Technologies	29·15		
29.5 TYPES OF ERRORS AND SYSTEM METRICS	29·16		
29.5.1 False Accept	29·16	29.8 SUMMARY AND RECOMMENDATIONS	29·24
29.5.2 False Reject	29·16		
29.5.3 Crossover Error Rate	29·16	29.9 FURTHER READING	29·26
29.5.4 Failure to Enroll	29·17		
29.6 DISADVANTAGES AND PROBLEMS	29·17	29.10 NOTES	29·26

29 · 2 BIOMETRIC AUTHENTICATION

29.1 INTRODUCTION. Once exclusively the purview of law enforcement, intelligence, and national security agencies, biometrics—the automated recognition of people based on their physiological or behavioral characteristics—has finally entered the business mainstream as a method of identification and authentication for access to physical and logical infrastructure. Once a pipe dream, biometric authentication technologies have substantially improved security, convenience, and portability over other commonly used methods of authentication. Adoption rates are still slower than some security professionals have desired; but falling costs, improvements in technologies, increased security needs, and changing government regulations are encouraging the adoption of biometrics.

29.2 IMPORTANCE OF IDENTIFICATION AND VERIFICATION. Ensuring the identity and authenticity of persons is a prerequisite to security and efficiency in present-day organizational operations. Intruders can damage physical and logical infrastructure, steal proprietary information, compromise competitive assets, and threaten organizational sustainability. Traditional methods of recognition and identification, wherein one individual identifies another based on his or her voice, physical appearance, or gait, are impractical, inefficient, and inaccurate in the scope of contemporary organizational operations. To address the need for rapid, efficient, and cost-effective authentication, organizations today primarily rely on the two methods of “something you know” and “something you have” (either applied individually or in combination) to verify the identity of persons accessing their physical and/or logical infrastructure. The most robust authentication systems use multiple factors of authentication. These forms of authentication are described in Chapter 28 in this *Handbook*.

29.3 FUNDAMENTALS AND APPLICATIONS. Biometrics are based on the measurement and matching of distinctive physiological and/or behavioral characteristics. The former are based on direct measurement of a physiological characteristic of some part of the human body. Examples of physiological biometrics include finger, hand, retina, face, and iris scans. The latter indirectly measure characteristics of the human body based on measurements and data derived from an action. Commonly used behavioral biometrics include voice and signature scan and keystroke pattern.

29.3.1 Overview and History. The use of nonautomated biometrics dates back to the beginning of human civilization, when individuals first began identifying other individuals based on certain physical or behavioral characteristics. The concept of biometrics as a means of authentication dates back more than 2,000 years. As early as 300 BC, Assyrian potters used their thumbprint as an early form of brand identity for their merchandise. The use of handwritten signatures (*chops*) in classical China is another example of an early biometric. In the first instance of a formal, legal biometric authentication system, fingerprints were used to sign contracts during the Tang dynasty (AD 618–906).

The development of contemporary biometric systems can be viewed as an outgrowth of the efforts of forensic scientists and law enforcement agencies to identify and classify criminals in the late nineteenth and early twentieth centuries. In 1882, Alphonse Bertillon introduced a system of body measurements called *anthropometry* to identify criminals. This system was proven unreliable because the measurements taken were not globally unique. A student of Bertillon, Edmund Locard, later proposed a fingerprint system based on the work of Sir Edmond Galton to identify people by analyzing unique points in fingerprint ridges and pores. Locard’s system, which used 12 Galton points,

FUNDAMENTALS AND APPLICATIONS 29 · 3

is considered reliable to this day. This methodology underlies fully automated modern biometrics, including the Integrated Automated Fingerprint Identification Systems (IAFIS) used by law enforcement agencies. Commercial biometric systems (typically relying on hand geometry) designed for use in physical access to buildings, emerged in the 1960s and 1970s.

Biometric methods of identification and identity verification, including automatic fingerprint analysis and facial recognition technologies, have been available and used by some government/public agencies (e.g., law enforcement, intelligence, and national security) and a few private industries (e.g., facial recognition scans in casinos) since the 1960s and 1970s. Notwithstanding the potential benefits and advantages over other authentication methods, biometrics have not been widely applied, particularly in the corporate world. Analysts cite high costs of equipment and implementation, technological problems, vulnerabilities of specific biometrics, lack of standards, and user resistance (notably, concerns over privacy) as reasons for the lack of implementation.

However, beginning in the 1990s, significant improvements in biometric technologies, a movement toward standardization, regulatory changes requiring organizations to adopt stringent security and privacy controls, and significantly reduced costs have encouraged wider adoption. A number of U.S. Government agencies (e.g., Department of Homeland Security, Department of Transportation, Department of Defense, Customs and Border Protection, Department of Justice, National Library of Medicine) and businesses in certain industries (e.g., health care and finance) have significantly increased their use of biometrics during the past few years—a factor that is likely to encourage other organizations to adopt biometrics as well. HP and Lenovo offer fingerprint scanners on their notebook computers, enabling users to increase security by requiring a finger swipe and a password (or just a finger swipe) to access files. Other computer manufacturers and peripheral vendors have added fingerprint scanning to computer keyboards and/or developed standalone fingerprint scanners that connect to computers via a Universal Serial Bus (USB) port.

Some analysts see the 2001 attacks on the World Trade Center and the Pentagon also as a key impetus behind the increased usage of biometrics for authentication and identification. The terrorist attacks have been critical in encouraging adoption not only because they have heightened the security concerns of companies and agencies but also because the impact and implications of the terrorists' attacks seem to have lowered users' resistance to the use of biometrics by employers and government. In other words, in the same way that the threat of global terrorism reduced public objections to possible infringements on civil liberties as a result of implementation of the USA PATRIOT Act and other security-focused measures, these attacks also appear to have rendered many people less sensitive to the potential privacy-invading implications of biometrics. Boroshok reported in 2007 that a survey sponsored by AuthenTec found that 71 percent of U.S. consumers would pay more for biometric security options in their cell phones and 63 percent of consumers would pay an additional cost for these options to be added to their personal computers.¹

The changing security environment has prompted forecasts of rapid growth in biometrics. In 2004, the International Biometric Group (IBG) had predicted rapid growth for the biometrics industry over the next several years from revenues totaling under \$50 million in 2004 to revenues of almost \$200 million in 2008.² In late 2003, analysts at the San Jose, California-based market research firm Frost and Sullivan predicted that biometric applications from commercial applications (not including the Federal Bureau of Investigation's IAFIS) would jump from \$93.4 million in 2001 to \$2.05 billion by 2006—up from the \$700 million (in 2006) that these analysts predicted prior

29 · 4 BIOMETRIC AUTHENTICATION

to the September 11, 2001, attacks. In January 2009, IBG released their most recent projections of market growth for the biometrics industry. They estimated that biometric industry revenues will grow from \$3.01 million in 2007 to over \$7.4 million by 2012.³ This is a far cry from their initial projected growth rate published in 2004, and may be attributed to a lack of standards, regulatory issues, and privacy groups that decried the invasion of their privacy by these devices (IBG hasn't publicly released a more recent report so we can't tell if these projections were met).

Despite rosy projections from industry analysts, adoption of biometric authentication systems continues to lag. Some firms continue to cite cost issues (though this is dwindling) and privacy concerns, while others point to problems surrounding biometric implementation in airports and among government agencies. Overall, surveys of companies indicate that forecasts of dramatic and rapid growth in biometrics implementation may be overstated.

This has been proven to be true by the recent 2013 Spending Priorities Survey put out by *InformationWeek*. Of the 513 businesses that responded to the survey, 13 percent had budget decreases from their 2012 budgets, with 43 percent having no change in their budgets—for many, it's a “do what you can with what you have” mentality. However, in spite of this, 58 percent of those surveyed wanted to improve their security overall.⁴ In another 2013 survey, however, the *InformationWeek* survey “Strategic Security Survey” revealed 52 percent of the more than 1,000 respondents wanted to increase their identity and password management practices.⁵ The seeming disconnect may be attributed to management wanting to keep costs down, and the ever-changing world of regulations—many are waiting to see what the landscape will look like when the dust settles.

A 2003 Forrester Research survey found that only 1 percent of companies had implemented biometric systems, just 3 percent had a biometric system rollout in progress, only 15 percent were testing biometrics, and 58 percent of those surveyed had *no plans* to try biometrics. Ten years later, we still have yet to see any substantial adoption of biometric authentication.

29.3.2 Properties of Biometrics. The contemporary meaning of biometrics emphasizes its automated aspects, which allow for deployment on a large scale. The most widely cited definition of biometrics is some variation of “the automatic identification of a person based on his or her physiological or behavioral characteristics.”⁶ The term “biometrics” generally is used as a noun to refer to the automatic recognition of persons based on their physical or behavioral characteristics. The term “biometric” can be used as a noun in reference to a single technology or measure (e.g., finger scan is a commonly used biometric) or as an adjective, as in “a biometric system uses integrated hardware and software to conduct identification or verification.”⁷

Biometrics have long been touted as a possible solution to the problems and vulnerabilities of other commonly used methods of authentication and identification. They represent sophisticated versions of the traditional means of identification, such as a guard allowing access to a user whom the guard recognizes by sight. Biometrics commonly are defined as automated methods of recognition/verification/identification of individuals based on some measurable physiological or behavioral characteristics, such as fingerprints, hand geometry, facial shape, iris pattern, voice, signature, and the like.

Whereas identification (ID) badges and keys authenticate the user based on something the user possesses, and passwords/personal identification numbers (PINs) authenticate the user based on what the user knows, biometrics allows authentication and identity verification based on who the user *is*. Because biometric methodologies of authentication actually base identification on physiological or behavioral “pieces” of the

FUNDAMENTALS AND APPLICATIONS 29 · 5

user, biometrics represents the only form of authentication that *directly authenticates the user*.

Biometrics have a number of other obvious advantages over other commonly used authentication methods. Unlike an ID badge or a USB key, one cannot easily lose or misplace a fingerprint or other biometric measures. Likewise, unlike the case with passwords and PINs, one does not need to remember and one is not subject to forgetting a physiological or behavioral characteristic. Although biometric measures *can* be compromised, in general, a biometric is much more difficult to manipulate by stealing, forging, sharing, or destroying than other commonly used authentication tools. Biometrics also provide considerable convenience, as opposed to the hassle of memorizing dozens of passwords.

Although the initial costs are quite high, the implementation of biometric systems typically results in much lower administrative costs than other access methodologies due to fewer calls to the helpdesk for technical support to reset passwords, no need to issue replacement ID badges, and so on. For these and other reasons, biometrics are viewed as providing better security, increased efficiency, and more reliable identity assurance than other commonly used methods of authentication/identification based on what a user possesses or what a user knows.

In theory, almost any human physiological and/or behavioral characteristic can be used as a biometric measure. However, to fit within a viable, potentially accurate, and practical biometric system, the biometric used should also satisfy four other requirements offered by Jain⁸ and Bolle, Connell, Pankanti, Ratha, and Senior⁹:

- 1. Universality.** Every person should have the biometric characteristic.
- 2. Uniqueness.** No two persons should be the same in terms of the biometric characteristic. Jain proposed the somewhat lower standard of distinctiveness, defined as “any two persons would be sufficiently different in terms of the characteristic.”¹⁰
- 3. Permanence.** The biometric should be relatively invariant over a significant period of time.
- 4. Collectability.** The biometric characteristic should lend itself to quantitative measurement in a practical manner.

Bolle, Connell, Pankanti, Ratha, and Senior argued that the biometric should also have a fifth attribute: acceptability, defined as “the particular user population and the public in general should have no strong objections to the measuring/collection of the biometric.”¹¹ Jain argued that a practical biometric system should consider two other attributes: (1) performance, which is “the achievable recognition accuracy and speed, the resources required to achieve the desired performance, as well as the operational and environmental factors that affect the performance,” and (2) circumvention, which “reflects how easily the system can be fooled using fraudulent methods.”¹² Along with this line of thinking, we find that the BYOD explosion among the workplace is opening an avenue for identity and access management. With the acquisition of AuthenTec by Apple in 2012, suppositions were running amok with the proposition of a new mobile-enabled biometric device for use inside and outside the enterprise. While the BYOD era has brought about the “acceptability” attribute, Jain’s idea of the proper operational and environmental factors are still there—how largely this will be accepted will be determined by how well these can be managed and centralized using IAM systems.¹³

29.3.3 Identification, Verification, and Authentication. Biometrics can be used to fill several roles related to identification and authentication (for more

29 · 6 BIOMETRIC AUTHENTICATION

background on those concepts, see Chapter 28 in this *Handbook*). These roles are often conflated in the use of the term *biometric authentication*, when in fact the system using the biometric data may be performing identification, verification, or authentication... or some combination of these. To avoid confusion in the use of terms, it is important to distinguish the various ways in which biometrics are used.

Identification systems answer the question, *Who do you claim to be?* They allow for selecting a single subject from a possible group of subjects based on the presentation of the identifying information. Identification systems often are referred to as $I:N$ (one-to- N or one-to-many) systems because the subject's biometric information is compared against multiple (N) records in the attempt to determine which one is a match. This kind of system has application in searching for a positive match (such as the facial recognition one might want to do in airports or stadiums for finding terrorists), as well as a negative match (designed to ensure that a person's biometric information is *not* present in the database, such as preventing people from enrolling more than once in large-scale benefits programs). But biometric identification also can serve as a convenient substitute to entering a username or swiping a card to proffer a claim of identity to an automated entry or logical access control system.

Verification systems answer the question, *Are you who you claim to be?* They attempt to match proffered information (passwords, token identifiers, or biometric data) with previously registered ("enrolled") credentials stored in the system; the match to be attempted is selected by a specific user identifier (ID) presented by the subject. Biometric verification systems are referred to as $1:1$ (one-to-one) systems because, while they may contain thousands or even millions of biometric records, they are "always predicated on a user's biometric data being matched against only his or her own enrolled biometric data."¹⁴ Thus, the presentation of biometric information is intended to confirm (or authenticate) the veracity of the subject's claim to the presented identity.

Whether one of these biometric systems qualifies as an *authentication* system depends on the implementation. A biometric verification system, which pairs a user ID with biometric information, always performs both identification and authentication. In that case, the biometric data is being used in place of a known password or the output of an authentication token, and may increase security by avoiding some of the vulnerabilities that plague other authentication credentials. A biometric identification system, on the other hand, does not always perform authentication. It *may* also fill the role of an authentication system if additional authentication factors (such as a PIN) are presented to confirm that the biometric information is in fact being presented by the subject the system identifies as being linked to the information. Without a second piece of information to confirm the linking, the use of a biometric identifier substitutes for the presentation of other user identifiers (such as typing in a user ID), and in this sense is not performing authentication.

It may be argued that because biometric data refers directly to some part of a subject's body or to an observable characteristic of the subject's behavior, presentation of biometric credentials should be considered a form of direct authentication. On this line of argument, in the same way that a security guard recognizes authorized employees at the gate in a combination identification/authentication action, a biometric identification system is in a sense authenticating the subject. This argument, however, glosses over the vulnerabilities that can plague both approaches: the guard can be fooled through disguises, can misidentify siblings (especially twins), or can simply miss a sufficiently similar face. Similarly, the automated biometric identification system can be fooled by spoofed credentials, as well as failing to distinguish between very similar biometric data. Relying on the single biometric presentation for both identification and

FUNDAMENTALS AND APPLICATIONS 29 · 7

implicit authentication is certainly quicker and more convenient, but may introduce risk that should be mitigated by a second authentication factor.

Biometric identification systems are more difficult to design and implement than verification systems because of the extensive biometric database search capabilities needed. Additionally, identification systems are more subject to error than verification systems, because many more matches must be conducted, matches that increase the opportunity for error. Verification systems are overall much faster (often rendering a match/no match decision within less than a second) and more accurate than identification systems. Verification systems, as opposed to identification systems, predominate in private sector applications, particularly for computer and network security applications. Verification systems also predominate in applications designed to authenticate rights-to-access to buildings and rooms, although sometimes identification systems are also deployed in high-security environments. Identification systems are often found in public sector applications, such as law enforcement (e.g., parole and prison administration, forensics, etc.), large-scale public benefits programs, intelligence, and national security applications.

29.3.4 Application Areas. Although there are many potential applications for biometrics, the primary ones can be divided into four categories: systems security (logical access systems), facilities access (physical access systems), ensuring the uniqueness of individuals, and public identification systems. The common thread among these four applications is that they all rely on individuals enrolled in the systems. The significance of whether individuals are enrolled or not enrolled will be explained in Section 29.3.5.

29.3.4.1 Security (Logical Access Systems). Logical access systems “monitor, restrict, or grant access to data or information.”¹⁵ Examples include accessing a computer or network or accessing an account. In these systems, biometrics replace or complement PINs, passwords, and tokens. The volume and value of electronic commerce plus the value of sensitive and personal information transported and/or stored on networks and computers make the use of biometrics to secure logical access a much more robust industry segment than physical security.

The use of biometric technologies for logical access control is still very much in its infancy. The most common biometric approach is to use fingerprint readers, either with a stand-alone USB reader or with a reader embedded in a laptop. Manufacturers are beginning to incorporate the Trusted Platform Module (TPM) chip in new laptops to support a variety of cryptographic applications. In combination with biometric devices like fingerprint readers, the TPM chip can allow applications like Microsoft’s BitLocker to apply biometric access control to encrypted volumes on a hard drive. Such technologies are still new enough to suffer problems with backward compatibility and inter-vendor support, but they offer the promise of much more secure logical access.

29.3.4.2 Facilities Access (Physical Access Systems). Facilities access systems “monitor, restrict, or grant movement of a person or object into or out of a specific area.”¹⁶ In these systems, biometrics replace or complement keys, access cards, or security cards, allowing authorized users access to rooms, vaults, and other secure areas. Physical access systems often are deployed in major public infrastructure settings, such as airports, security checkpoints, and border facilities, in order to monitor and restrict movements of unauthorized or suspicious persons. In addition to entry to secure rooms, physical access systems, when applied in business settings, include

29 · 8 BIOMETRIC AUTHENTICATION

time-and-attendance systems by combining access to a location with an audit of when the authentication occurred.

Biometric technologies have been in use for physical access control for some time but still represent a range of possible implementations, from stand-alone fingerprint-reading door locks to complete systems with central storage of biometric templates, logging, and power failure protection. Selecting a system from within such a wide product range must fit into the overall security stance of the facility, particularly in regard to the storage of templates. In a simple stand-alone door system, for example, fingerprint templates would be stored in or near the locking device, and the system might not have any logging features. Unless coupled with surveillance and intrusion alarms anticipating physical compromise, these devices are more appropriate for lower security applications, such as storage facilities where physical key access might not be preventing theft. More centrally connected systems, however, offer greater integration into overall monitoring and control of access but suffer from communication and power issues as they push templates across a network. Selecting a product that interacts well with other parts of the overall access control system is crucial.

29.3.4.3 Ensuring the Uniqueness of Individuals. Uniqueness biometric identification systems typically focus on preventing double enrollment in programs or applications, such as a social benefits program. The main use of this application occurs in the public sector, although similar systems could be implemented to prevent double enrollment in employee benefits programs.

29.3.4.4 Public Identification Systems. A final biometric application of note is its use to identify criminals and/or terrorists. Criminals and terrorists can wear disguises, acquire fake documents, and change their names, but biometric data are fairly difficult to forge. In 2004, the United States Department of Defense started the Automated Biometric Identification System (ABIS), which collects biometric data on Iraqi insurgents in a manner compatible with IAFIS. This allows for identification of known repeat offenders and wanted persons. Soldiers also use the Biometric Automated Toolset (BAT) developed by the Army's Language Technology Office to identify persons on the scene of bombing attacks. Anyone present in the area can be cross-referenced with an existing database of insurgents. The BAT is also used to enroll and identify members of the Iraqi army. Although biometric systems hold much promise, it is also important to understand there are still limitations of current technology, but with the advent of nanotechnology, many of these are quickly being overcome.

29.3.5 Data Acquisition and Presentation. As Jain explains, “A biometric system is essentially a pattern recognition system that operates by acquiring biometric data from an individual, extracting a feature set from the acquired data, and comparing this feature set against the template set in the database.”¹⁷ The starting point for the biometric system is *enrollment*: A user’s biometric data are initially collected and processed into a template, the form in which they are then stored for ongoing use. As Woodward, Orlans, and Higgins explain, “Templates are not raw data or the scanned images of a biometric sample, but rather they are an accumulation of the distinctive features extracted by the biometric system.”¹⁸ Liu and Silverman describe the template as “a mathematical representation of biometric data. A template can vary in size from 9 bytes for hand geometry to several thousand bytes for facial recognition.”¹⁹ Templates are proprietary to each vendor and technology with little or no interoperability between systems. This lack of interoperability is attractive from a privacy perspective

TYPES OF BIOMETRIC TECHNOLOGIES 29 · 9

but unattractive from the perspective of cost effectiveness and the prospective implementer who is concerned about committing significant investment to a single nonstandardized technology.

The term “presentation” refers to the process by which a user provides biometric data to an acquisition device by looking in the direction of a camera, placing a finger on a pad or sensor, or some other specified physiological exam. For purposes of verification or identification, the user presents biometric data, which are then processed and converted to a template. This template is an extraction of distinctive features and is not adequate for the reconstruction of the original biometric data. The scanned template is then matched against the stored enrollment template(s). Each time a user makes a presentation, a new template is created and matched. It is important to note, especially from the perspective of privacy concerns, that biometric systems do not store raw biometric data; instead they use the data for template creation and, in most cases, discard the biometric data. The biometric system’s match/no-match decisions are based on a score, which is “a number indicating the degree of similarity or correlation resulting from the comparison of enrollment and verification templates.”²⁰ Like the templates, the scoring system is based on proprietary algorithms; there is no standard system.

29.4 TYPES OF BIOMETRIC TECHNOLOGIES. As previously noted, biometrics generally can be grouped into two categories: physiological and behavioral. The International Biometric Group provides data on comparative market share of various biometric technologies. The IBG data focus on market share from commercial and government applications. The top eight biometric technologies are all from the physiological category (see Exhibit 29.1).

Five of these biometrics are discussed in detail in the next sections.

29.4.1 Finger Scan. Finger scan or fingerprint technology is by far the most widely deployed biometric technology. Finger scan’s number-one status as a biometric is maintained even if the extensive use of fingerprinting by law enforcement agencies is excluded. The type of fingerprinting employed in commercial biometric systems differs from the one used in law enforcement. In most commercially available biometric applications, the station provides only for the scan of a single finger on one hand, whereas

EXHIBIT 29.1 Comparative Market Share of Biometric Technologies by 2015

Biometric Technology	Market Share	
		2015 (projected)
Finger scan		15%
Facial scan		15%
Hand scan		1%
Iris scan		16%
Vein		10%
Voice		13%
Signature		10%
AFIS/Livescan		16%

Sources: (1) Market share data from PRWeb, Retrieved May 29, 2013, from <http://www.prweb.com/releases/2007/04/prweb516626.htm>
(2) J. McHale, “Biometrics: The Body’s Keys,” *Military & Aerospace Electronics* 14, No. 12 (2003): 17–23.

29 · 10 BIOMETRIC AUTHENTICATION

law enforcement agencies often rely on full sets of fingerprints. In addition to being the most widely used biometric, fingerprinting is also one of the oldest and most well-researched biometric technologies. Because it is a widely used, well-documented, and mature technology, costs for the deployment of finger-scan-based technologies are relatively low. Single-quantity pricing for a workstation version with associated software can be as low as \$150; server versions are currently priced as low as \$50 per unit.

The strengths of finger scan are one of the principal reasons for its popularity and include:

- Wide use.
- Mature technology.
- Low cost.
- High ease of use. (Very little training is required to place a finger on a finger pad.)
- Ergonomic design. (Comfortable to use for most users.)
- Low error incidence. (False match rates are extremely low; crossover error rate is lower than voice scan and facial recognition, higher than hand geometry and iris scan.)
- Fast transaction times. (In most systems, authentication takes less than a second.)
- Capacity to be deployed in a wide range of environments (e.g., on workstations, doorways, indoors/outdoors).
- Ability to increase accuracy levels by enrolling multiple fingers.
- Can provide identification with a high level of accuracy (if properly configured to include multiple enrolled fingers) in addition to verification.

Despite its multiple strengths, finger scan is not without significant weaknesses. As Chirillo and Blaul note, some of this technology's weaknesses stem from the same factors that lend it its strengths. "Because fingerprint technology is one of the oldest and most well-known technologies, a good amount of information is publicly available on how to defeat it."²¹ A number of ways exist to foil finger scans and produce a false match (false accept), including the use of a dummy finger constructed of latex or other material, manipulation of the scanner so as to raise the latent print of the person who used the scanner previously, and even use of an actual finger that is no longer attached to a body. (Most finger scanners cannot discriminate between live and dead tissue.) Because of these factors, the security levels of finger scans are not actually as impressive as the low error rates seem to indicate. It should be noted that countermeasures could be taken to overcome the vulnerability of finger scans to fraud. For example, enrolling additional fingers makes fraud more difficult. To reduce the chance that the system will be foiled by synthetic or dismembered fingers, thermal and/or moisture scanners can be added to the sensors to detect finger temperature and moisture levels that would indicate the vitality of the finger.

Other weaknesses include:

- A scanner requires frequent maintenance because screens/sensors tend to retain an obstructing buildup of user skin oil and residue.
- Performance can deteriorate over time, both because of aging of the users (and wearing away of fingertips) and because of the need for system maintenance.
- Finger-scan biometrics are obviously not appropriate for users with missing hands or hand disabilities.

TYPES OF BIOMETRIC TECHNOLOGIES 29 · 11

- Performance levels deteriorate among users who have hand tremors because the presentation of biometric data will be distorted.
- Performance levels also deteriorate when users' fingers are either overly dry (a certain amount of normal skin moisture is needed for an accurate reading) or overly moist or oily (as from too much hand lotion).
- There is a small but significant failure to enroll (FTE) rate even among a population with hands and without disabilities. The FTE rate for finger scans is estimated at 2 to 10 percent and is attributed to persons with genetically indistinct prints, scarred fingers, dry skin, and fingerprints worn down by age and/or manual labor.

Perhaps the biggest weakness of finger scan, however, has nothing to do with the accuracy and reliability of the technology. Instead, it relates to user acceptance. Because of the association of finger scans with law enforcement and criminality, often such scans are not readily accepted by users who dislike the technology's "taint" with forensic applications and who may worry that finger-scan biometric data will be used for other purposes.

According to Chirillo and Blaul, "Another reason fingerprint technology is not highly accepted is that it may require individuals to share or touch the same device that others touch."²²

29.4.2 Facial Scan/Recognition. Bolle, Connell, Pankanti, Ratha, and Senior note that "face appearance is a particularly compelling biometric because it is one used every day by nearly everyone as the primary means for recognizing other humans. Because of its naturalness, face recognition is more acceptable than other biometrics."²³ However, user acceptance of facial scans drops significantly when users discover that it has been used covertly. As Imparato observes, "Of all the biometric technologies currently in use, face recognition is arguably the most controversial."^{23,24} Like finger scans, face recognition relies on the identification of unknown face images by comparison to a database of (known) face templates. Face recognition is used overtly in access control where it is used for one-to-one identification. This application yields relatively high performance because the environment is highly controlled and input data are predictable. Face recognition may also be used covertly in surveillance where it is used to locate people in crowds (one to many), albeit with mixed results. For example, the city of Virginia Beach, Virginia, has employed a face recognition system to identify known (pre-enrolled) felons for over five years, but has yet to identify a single criminal. A variety of facial recognition technologies, ranging from single image, video sequence, three-dimensional image, near infrared, to facial thermograms, are available.

Facial recognition offers these benefits:

- It has the capacity to leverage existing image acquisition equipment, such as digital cameras, Web, video, and the like.
- Because facial recognition is a software-based technology, it is often unnecessary to purchase new hardware, especially given the number of Closed Circuit Television (CCTV) and surveillance cameras in broad use.
- The lack of need for specialized hardware can help keep the cost of this technology down, assuming that high software costs do not counterbalance the savings from the hardware.
- It is the only biometric capable of identification at a distance without the subject's cooperation or even awareness.

29 · 12 BIOMETRIC AUTHENTICATION

- It is easy to use. All that is required is that the user (or target) look at the camera.
- It does *not* require the user to touch any device (a major objection for some users with finger scans and hand scans).
- When deployed in verification situations, facial scans have extremely low failure-to-enroll rates. (Unlike fingerprints, human faces are almost always distinctive.)
- They are capable of enrolling static images (e.g., photographs on driver's licenses), a factor that makes it possible to implement very large scale enrollments at a relatively low cost and in a brief amount of time.

Facial recognition systems have a number of serious weaknesses too. The predominant weakness (which derives from a combination of the technology's other weaknesses) is the low accuracy and high error rate of this biometric. Whether deployed covertly or overtly, facial recognition has the lowest accuracy rate among all five top biometrics. Tilted heads and low camera resolution still confounds the programs, as evidenced in the Boston bombings of 2013.²⁵

Evidence of the technology's low accuracy rate comes from a study at Palm Beach (Florida) International Airport that showed that the system failed more than 50 percent of the time to match the 15 employees who had enrolled in the database for a trial run. Out of 958 pass-throughs, the system matched the employees' faces just 455 times. Some studies suggest that accuracy improvements can be made in facial recognition systems, but these improvements will come at a very high cost. For example, a facial recognition software package from Visionics FaceIt resulted in impressively low error rates, as long as lighting conditions were perfect. The software costs \$30,000 for a three-camera system.

Other weaknesses include:

- False matches (false accepts) routinely occur in the case of twins, and most systems are insensitive enough for someone skillful at disguise and impersonation to trick the system into a false match.
- More likely than false matches, however, are false nonmatches (false rejects), which can occur as a result of facial expressions; changes in hairstyle, makeup, facial hair, significant changes in body weight, eyeglasses, and age-related facial changes.
- The acquisition environment can have a dramatic impact on facial recognition system accuracy. In particular, lighting, either too bright or too dim, can dramatically increase the error rate.
- The perceived threat to privacy. Overtly deployed facial recognition technologies (e.g., used for identification and access) are generally judged relatively unobtrusive and meet with a high level of user acceptance. However, covertly deployed systems, such as those used for surveillance, pose significant threats to privacy. This threat is generally viewed as much more serious than that posed by the other top biometrics.

29.4.3 Hand Geometry Scan. Hand geometry scans refer not to handprints or to any analogy of fingerprints but rather to the geometric structure (or geometric invariants) of the human hand. Nanavati, Thieme, and Nanavati explain that "hand-scan technology utilizes the distinctive aspects of the hand—in particular, the height and width of the back of the hand and fingers—to verify the identity of individuals."²⁶ The

TYPES OF BIOMETRIC TECHNOLOGIES 29 · 13

leading hardware maker for this technology, Recognition Systems, Inc. (RSI), has a basic hand scanner that takes upward of 90 measurements from three to four enrollments to create a user template that includes length, width, and thickness, plus surface area of the hand and fingers. Newer systems include temperature-sensing mechanisms to ensure “live” subjects. All the components of a hand scan system (acquisition hardware, matching software, storage components) reside within a stand-alone device. Hand scans are a well-established biometric technology (they have been in widespread use since the 1970s), but compared to other leading biometrics, hand scans tend to be much more limited in their range of applications. Hand scans are used exclusively for verification rather than for identification because the hand measurements are not distinctive or specific enough to allow for identification applications. For this reason, hand scans are used mostly for physical access and time-and-attendance applications. In the latter case, they are used as a way to eliminate the problem of “buddy-punching” whereby one employee punches in or out for a coworker who is not present.

Hand scan technology has changed very little since it was first introduced over 30 years ago, so its strengths and weaknesses are well established. The principal strengths of the hand scan include:

- Operates in very challenging environments. (The equipment is typically unaffected by light, dust, moisture, or temperature.)
- Established and reliable technology.
- Ease of use. (Users simply stick their hand in the unit; placement matters little.)
- Resistance to fraud compared to other biometrics. (It would be difficult and time consuming to substitute a fake sample.)
- Small template size (as low as 9 bytes; much smaller than other biometrics, allowing for storage of thousands of templates in a single unit).
- Based on a relatively stable physiological characteristic.
- High level of user acceptance and lack of attached stigma.

Problems reported in using hand scans include:

- Limited accuracy (which in turn limits its use to verification not identification). The relatively low accuracy of hand scan (higher than facial recognition and behavioral biometrics but lower than finger and iris scans) is a result of the general lack of physical variety expressed in the hand as well as the relatively small number of features measured by hand scan.
- Comparatively large form factor (This limits the technology’s deployment in computer-oriented applications that require hardware with a smaller footprint.)
- Some people resent forced contact with possibly unclean surfaces.
- Ergonomic design limits its use by some populations (e.g., the disabled).
- Comparatively high cost. While the cost of palm scanners has come down to around \$500 for a USB model, enterprise level scanners are still a rather pricey unit, even with the increase in companies producing them.

29.4.4 Iris Scan. Iris scan technology uses the unique pattern formed by the iris—the colored part of the eye bounded by the pupil and the sclera—to identify or verify the identity of individuals. The iris pattern is unique, for even in the same

29 · 14 BIOMETRIC AUTHENTICATION

individual, no two irises are alike. The uniqueness of iris patterns has been likened to that of multilayered snowflakes. To put this into perspective, the chance of mistaking one iris pattern for another is 1 in 10^{78} ! Iris scans accomplish this by allowing for more than 200 points of reference for comparison, as opposed to the 60 to 70 used in fingerprints. The unique aspects of the iris make it an ideal biometric for high-security applications; enrolling both irises from the same individual can enhance the level of security. In addition to high-security physical access applications, iris-scan technology has been used in automated teller machines (ATMs) and banking kiosks. It is also now being used among police departments in handheld units to identify suspects, and is being introduced for such jobs as border crossings, and is being regularly used in India, Iraq, and Dubai.²⁷

The most important strength of iris biometrics is its accuracy, the most critical weakness of facial scanning. Of all the leading biometrics, iris technology has the lowest error rate and the highest level of overall accuracy. Other strengths of this biometric include:

- Ability to be used both for verification and for identification.
- Stability of its biometric characteristics over a lifetime.
- Relative difficulty to fake or spoof because it is an internal biometric.
- The fact that the iris is minimally subject to outside influences when compared to biometrics like fingerprints and faces.

The major weaknesses of the iris biometric concern user perceptions and problems in the user-technology interface. Other weaknesses include:

- Acquisition of the image requires moderate training and attentiveness: Users must stand still and look straight into the scanner with eyes open and unblinking.
- Users often report some physical discomfort with the use of eye-based technology, although less so than with retina scanning technology.
- Anecdotal reports also suggest a fairly high level of user psychological resistance to iris-scanning technology, with some users believing that the scanner will lead to eye damage.
- Can be adversely affected by lighting and other environmental conditions (although not to the extent of facial scanning).
- In some cases eyewear adversely affects performance (although many iris devices can scan people wearing glasses or contact lenses).
- Although the iris is a relatively stable biometric, it is affected by aging and disease.
- Relies on proprietary hardware and software technologies.
- Costs tend to be high compared to finger scanning, hand scanning, and many facial recognition systems.

On the other hand, the per unit cost of the leading hardware/software combination technology has dropped to as low as \$300 per seat, still higher than finger scans but significantly lower than the over \$5,000-per-seat price seen a few years ago.

29.4.5 Voice Recognition. Voice recognition biometrics “utilizes the distinctive aspects of the voice to verify the identity of individuals.”²⁸ Voice recognition

TYPES OF BIOMETRIC TECHNOLOGIES 29 · 15

generally is classified as a behavioral biometric, although it actually combines elements of behavioral and physiological biometrics: “The shape of the vocal tract determines to a large degree how a voice sounds, a user’s behavior determines what is spoken and in what fashion.”²⁹ Stated somewhat differently, “voice is a behavioral biometric but is dependent on underlying physical traits, which govern the type of speech signals we are able and likely to utter.”³⁰ Because of comparatively low levels of accuracy and considerable user variability in voice dynamics, this biometric generally is used only for verification, not identification. Commonly deployed voice recognition systems can be divided into two types: text-dependent systems (the speaker is prompted to say a specific thing) and text-independent systems (the authentication system processes any utterances of the speaker), which provide a higher level of security because they are more difficult to spoof and provide better accuracy than text-dependent systems.

Strengths of voice recognition include:

- Capacity to leverage existing telephony infrastructure (as well as built-in computer microphones).
- Low cost when existing infrastructure is used.
- Ease of use.
- Interface with speech recognition and verbal passwords.
- High level of user acceptance. (This biometric does not suffer from the negative perceptions associated with all of the other leading biometrics.)

Weaknesses of voice recognition include:

- More susceptible to replay attacks than other biometrics.
- Accuracy levels are low compared to iris scanning, finger scans, and hand scans.
- Accuracy levels are negatively affected by ambient noise and low-quality capture devices.
- Accuracy, security, and reliability are challenged by individual variations in voice, such as speaking softly or loudly, hoarseness or nasality because of a cold, and so on.
- The stability of the biometric is affected by illness, aging, and other user behaviors including smoking.

29.4.6 Other Biometric Technologies. The five major biometric technologies just discussed collectively comprise the vast majority of biometric technology under deployment. Other biometric technologies that register on market share breakdowns are two of the behavioral type: signature scan and keystroke scan. Although both of these behavioral biometrics are well accepted (signature scanning more so than keystroke scanning), their usefulness is limited by their lack of accuracy.

Other behavioral biometrics under investigation include gait and lip motion. One physiological biometric that has received considerable attention because of its high accuracy and security rates is retinal scanning. However, most analysts believe that the problems associated with retinal scanning (lack of user acceptance, high cost, difficult and painful acquisition process) outweigh any advantages to this biometric. The consensus seems to be that iris scanning has replaced retinal scanning as the eye scanning biometric of choice.

29 · 16 BIOMETRIC AUTHENTICATION

The use of DNA as a biometric identifier has also been investigated, although it has significant weaknesses including the fact that DNA in body tissues (e.g., epithelial cells) can be obtained surreptitiously and transferred easily for nefarious purposes whereas the official methods of collection (e.g., taking blood samples) are relatively intrusive.

Other physiological biometrics that may prove useful in the future include body odor, skin reflectance, and ear shape.

29.5 TYPES OF ERRORS AND SYSTEM METRICS. All types of identification and authentication systems suffer from two types of errors: false accepts and false rejects.

29.5.1 False Accept. Also known as false match, false positive, or type 1 error, false accept is the likelihood, expressed as a percentage, that an imposter will be matched to a valid user's biometric. In some systems, such as those that attempt to secure entry to a weapons facility, a bank vault, or a high-level system administrator account, the false match/false accept rate is the most important metric to watch. In other systems, such as a facial recognition system deployed by a casino in an effort to spot card counters, a high level of false matches may be tolerated.

29.5.2 False Reject. Also known as false nonmatch, false negative, or type 2 error, false reject is the probability that "a user's template will be incorrectly judged to *not* match his or her enrollment template."³¹ False nonmatches typically result in the user being locked out of the system. These false nonmatches can occur because of changes in a user's biometric data, changes in how the biometric data is presented, and/or changes in the environment. Biometric systems are generally more susceptible to false rejects than they are to false accepts.

29.5.3 Crossover Error Rate. An important metric in biometric systems is the *crossover error rate* (CER), also known as the equal error rate (EER). This useful metric is the intersection of the false accept and false reject rates. In general, a lower CER indicates the biometric device is more accurate and reliable than another biometric device with a higher CER. Exhibit 29.2 provides a summary of benchmark test-based accuracy/error rates for the five most prevalent biometric technologies. Each biometric technology is rank-ordered from most accurate to least accurate based on CER.

EXHIBIT 29.2 Accuracy/Error Rates of Leading Biometric Technologies

Biometric	False Match Rate	False No-Match Rate
Iris scan	0.0001%	1.1–1.4%
Finger scan	0.02%	.6%
Hand scan	0.3%	3.0%
Voice (text independent)	2–5%	5–10%
Voice (text dependent)	2.0%	0.03%
Face scan	.1%	1–2.5%

Source: Based on data contained in Anil K. Jain (2008), "Biometric Authentication," *Scholarpedia*, 3(6); 3716.

DISADVANTAGES AND PROBLEMS 29 · 17

29.5.4 Failure to Enroll. Another critical metric in biometric systems is FTE. As Ashbourn explains, FTE refers to “a situation whereupon an individual is unable to enroll their biometric in order to create a template of suitable quality for subsequent automated operation.”³² Common reasons for failure to enroll include physical disability and a user whose physiological/behavioral characteristics are less distinctive than average. Nanavati, Thieme, and Nanavati observe that failure to enroll can be a major problem in “internal, employee-facing deployments” in which “high FTE rates are directly linked to increased security risks and increased system costs.”³³ A final important metric is the “transaction time.” Transaction time refers to “a theoretical time taken to match the live template against a reference sample.”³⁴

29.6 DISADVANTAGES AND PROBLEMS

29.6.1 General Considerations. Despite the many advantages over other commonly used authentication systems, the implementation of biometric authentication controls carries a number of risks and disadvantages. Even the most accurate biometric system is not perfect, and errors will occur. The error rates and the types of errors will vary with specific biometrics deployed and the circumstances of deployment. Certain types of errors, such as false matches, may pose fundamental, critical risks to organizational security. Other types of errors—failure to enroll, false nonmatch—may reduce organizational productivity and efficiency and increase costs. Organizations planning biometrics implementation will need to consider the acceptable error threshold. In any event, organizations deploying biometric authentication systems must not be lulled into a belief that they are invulnerable to errors and/or fraud. Certain biometric systems (e.g., iris scanning) are fairly impervious to fraud, while others (especially behavior-based systems) are much more susceptible to it. Facial scanning systems can be foiled with clothing, makeup, eyeglasses, and/or changes in hairstyle. Even relatively stable physiology-based biometrics like fingerprint scans can be defrauded with the use of rubber or gelatin fingers. Matsumoto outlines a gummy finger approach designed to fool even those countermeasures mentioned in Section 29.4.1.³⁵ The protein used has a similar galvanic response to flesh and, since it is very thin and attached to a live finger, has the correct temperature. In some cases, blowing warm air over the scanner may even raise the latent print of the intruder’s predecessor.

The deployment of commonly used authentication systems (i.e., ID badges, passwords, etc.) requires relatively little training, although one could argue that better training on the development and use of passwords would improve security. This limited need for training is not the case with most of the most commonly used biometric systems. Both systems administrators and users need instruction and training to ensure smooth operation of the system. Some biometric systems are exquisitely sensitive to intra- and inter-user variation in presentation and performance. Their effectiveness becomes substantially compromised and error rates substantially increase in cases of significant variation and/or irregular presentation. A related problem concerns user acceptance of the biometric system. Some users may object to the deployment of biometrics due to concerns over privacy and intrusiveness. In other instances, users may object to the deployment of biometrics and avoid optimal interface with the system because of safety and/or health concerns, general fears, and/or cultural and religious beliefs. For example, some individuals may be concerned that biometric systems that require them to touch a finger pad or hand pad will unnecessarily expose them to germs and place them at risk for illness. Some users may fear that eye scans will damage their

29 · 18 BIOMETRIC AUTHENTICATION

eyes. Other users may object to eye scans on the basis that the eyes are the window to the soul. Anderson notes that some persons may object to the use of biometrics due to a personal interpretation of religious doctrine.³⁶ Notwithstanding users' beliefs and perceptions about the biometric system, in many cases features or elements related to the users and/or the operating environment will influence the successful implementation and effectiveness of the system.

29.6.2 Health and Disability Considerations. Individuals with arthritis and/or certain other disabilities and physical limitations may be unable to enroll in systems and subsequently, to align themselves physically in an optimal position with respect to biometric sensors. For example, users with severe hand arthritis may be unable to place their hand firmly as required on the hand geometry sensor, and users with migraines and associated photophobia may find it physically too uncomfortable to look straight into the light sensor for the iris scan. Some disabled people may have to be excluded from biometric systems altogether. Some relatively minor disabilities, such as a slight tremor, may compromise a legitimate user's ability to gain access through certain biometric systems. Variations in physical size can also influence system accuracy. An iris scanner positioned for a standard height range may fail to capture images of either very short or very tall individuals, or in some cases an individual's hands or fingers may be either too large or too small to be read accurately in a hand or finger scanner. Likewise, individuals with neck and back problems may find it difficult to use some biometric devices, depending on the kind of positioning required. Systems that rely on behavioral biometrics such as voice or signature are particularly vulnerable to variations and irregularities in user characteristics. For example, users who speak too softly, too loudly, or too rapidly may cause system errors. Minor changes in users' health can affect some biometric readings. Excessive skin moisture or lack of skin moisture can impact finger scans.

Although one of the ideal properties of a biometric is its universality, in reality not everyone has the characteristic or has it to the same degree. For example, some people are born without distinct fingerprints. In other cases, users may have lost the distinctiveness of their fingerprints because of years of manual labor, use of certain chemicals, scarring, or the aging process. Anderson notes that "people with dark-colored eyes and large pupils give poorer iris codes."³⁷ Certain eye diseases and metabolic conditions may also reduce or negate the efficacy of eye scan authentication. Age has a significant impact on the user-biometric-system interface. Definite physiological changes are associated with the aging process and can result in poor template matching with the live biometric. In this case, reenrollment may be needed. Fingerprints are affected by the aging process as the skin becomes drier and more brittle; voice patterns change in tonal quality over time; and facial shape or appearance may shift with age. Overall, the acceptability of a biometric system will be lessened if there is the impression that implementation of the system discriminates against, or has an otherwise adverse impact on, the disabled, the ill, ethnic minorities, the elderly, and other protected or traditionally disadvantaged groups of users. Organizations must ensure compliance with the Americans with Disabilities Act when implementing biometric authentication systems. Compliance may involve providing alternative methods of authentication to those affected.

29.6.3 Environmental and Cultural Considerations. A broad range of factors in the operating environment can also impact the effectiveness and acceptability of biometric systems. User-related cultural, social, and behavioral factors can influence

DISADVANTAGES AND PROBLEMS 29 · 19

system performance. For instance, the accuracy of facial scans can be compromised by users' changes in hairstyle, facial hair, and headwear as well as by changes in an individual's physical appearance because of significant weight gain or loss. The accuracy of voice/speech recognition systems is affected by the distance between the scanner and the user as well as by the volume of speech. Fingerprint recognition is impeded in cases when users' skin is too dry, whether the condition arises as a result of aging, skin disease, environmental factors, or occupation-related factors, such as frequent hand washing among healthcare professionals. Factors in the surrounding ambient environment may also affect the accuracy of the biometric system. Ambient lighting will influence accuracy and error rate in facial scans and, to a lesser extent, in iris scans. Noise levels can impede the effectiveness of voice recognition systems. Humidity and air temperature can affect the accuracy of fingerprint and hand scans.

29.6.4 Cost Considerations. Although the cost of biometric system implementation has fallen dramatically in the past few years, it is still a major barrier for many organizations. Costs vary significantly depending on the type of system. Recent reports suggest that newer fingerprint scanners can be purchased for as little as \$50 per unit; voice recognition systems can cost in excess of \$50,000. However, even the least expensive biometrics systems are likely to cost more than simpler versions of traditional authentication systems. Experts estimate minimum costs, including hardware and software, at \$200 or so per user and upward of \$150,000 for corporate-wide protection in a medium-sized business. Compounding the cost issues are problems related to the lack of clear standards and the lack of clear interoperability between various biometric authentication systems.

Many of the problems and difficulties with biometrics systems are likely to be corrected or significantly mitigated with technological improvements, better user and administrator training, and good control of environmental conditions. In other cases, problems can be overcome or ameliorated with the use of countermeasures, such as combining different types of biometrics, combining biometrics with traditional authentication systems, and so on. Two major concerns that will continue to loom large and deserve closer examination are biometric identity theft and user privacy.

29.6.5 Attacks on Biometric Systems. Although biometrics are much less vulnerable to attack than other authentication controls, they are not immune to fraud. Moreover, when a biometric identity is stolen or spoofed, it creates a much bigger problem than that created by the theft of an ID badge, USB key, or password because a biometric cannot be simply canceled and replaced. One of the principal advantages to using biometrics for authentication is their invariability over time. Consequently, when an imposter or intruder defrauds a biometric authentication system and creates a false match error, the entire biometric security system is defrauded and the individual authorized user's biometric integrity is compromised. Likewise, Prabhakar, Pankanti, and Jain note, "One disadvantage of biometrics is that they cannot be easily revoked. If a biometric is ever compromised, it is compromised forever."³⁸

A number of analysts believe that the ultimate solution to the problem of biometric identity theft lies in the development of "cancelable biometrics." Researchers at IBM have developed a prototype for the cancelable biometric that incorporates a repeatable distortion of the biometric. Similar in theory to the use of public and private keys for encryption, a unique distortion of the biometric is introduced at each enrollment. Therefore, if a user's biometric is compromised, only the one system is defrauded, not every system in which the user is enrolled.

29 · 20 BIOMETRIC AUTHENTICATION

29.6.6 Privacy Concerns. The use of biometric authentication controls raises significant privacy concerns, particularly in comparison to conventional authentication methods like passwords and ID badges. User objections to biometrics are often based on privacy concerns, sometimes articulated in terms of the user's sense of the intrusiveness of the biometric system. Anecdotal reports suggest that public perceptions of intrusiveness vary among different biometrics and in how biometrics are implemented. With regard to the latter, Nanavati, Thieme, and Nanavati³⁹ report that there is a greater risk of privacy invasiveness when:

- Deployment is covert (users are not aware of the system's operation) versus overt.
- The system is mandatory versus opt-in.
- The system is used for identification rather than verification.
- It is deployed for an indefinite duration versus fixed duration.
- It is deployed in the public versus the private sector.
- The user is interfacing with the system as an employee/citizen versus an individual/customer.
- An institution, not the user, owns the biometric information.
- The biometric data are stored in a template database versus the user's personal storage.
- The system stores identifiable biometric data versus templates.

A vivid example of the public's lack of acceptance of the covert use of biometric systems comes from the 2001 Super Bowl and the uproar that ensued after the Tampa Police Department deployed facial scanning technology for the purpose of picking out criminal suspects from the audience. In contrast, in the aftermath of the 2001 attacks on the World Trade Center and the Pentagon, there has been fairly widespread public acceptance of the use of facial scanning at airports in the United States.

Users generally view behavior-based biometrics, such as voice recognition and signature verification, as less intrusive and less privacy-threatening than physiology-based biometrics. Facial scanning is viewed as having a high potential for privacy invasion because of the capacity to deploy it without the user's knowledge and participation. Finger scans may be viewed as intrusive and privacy-invasive because of their association with law enforcement functions. The level of intrusiveness of the scanning technique appears to affect users' perception of privacy invasion, with iris scanning provoking more privacy objections than hand scanning. Civil libertarians and users also raise privacy objections over biometric systems that have the potential to uncover additional information about the user beyond the biometric identity. For example, finger scans, because of their capacity to be linked to large law enforcement databases of fingerprints, could be used to reveal information about the user's criminal background. Iris scans have the capacity to reveal confidential medical/health information about the user. Probably one of the most troubling privacy-related aspects of biometrics is the potential for large-scale linkage between biometric systems and the use of biometric data to facilitate large-scale national ID programs. Even though employers may design a biometric system for purely in-house use in order to facilitate verification of employee identities on corporate networks, federal regulations and laws such as the USA PATRIOT Act may eventually compel employers to surrender employees' private biometric data to government authorities.

DISADVANTAGES AND PROBLEMS 29 · 21

In summary, the major privacy concerns associated with biometric deployments include:

- Users' loss of anonymity and autonomy
- Risk of unauthorized use of biometric information and/or unauthorized collection of biometric information
- Unnecessary collection of biometric information
- Unauthorized disclosure of biometric information to others
- Systematic reduction of users' reasonable expectation of privacy
- Potential for misuse on the part of overzealous or corrupt government agents

Many of these concerns can be generally lumped under the heading of "function creep" or "mission creep," wherein biometric systems designed for user authentication may, over time, be used for purposes not originally intended. An example of "mission creep" is the use of Social Security numbers (SSNs) for identification. The original Social Security cards were stamped "Not for Identification." However, many organizations (including the Internal Revenue Service) use SSNs for identification purposes.

Notwithstanding the privacy risks, supporters of biometric authentication systems argue that, properly deployed and with adequate best practice controls, biometric systems actually can function to enhance and protect privacy. Woodward, Orlans, and Higgins point out that "several newly developed biometric technologies use an individual's physical characteristics to construct a digital code for the individual without storing the actual physical characteristics," thus creating a sort of *biometric encryption* that can be used to protect the privacy of an individual's financial, medical, or other data.⁴⁰ Nanavati, Thieme, and Nanavati argue that "privacy-sympathetic" biometric systems can be designed.⁴¹ Such systems would:

- Have limited system scope.
- Eschew use of biometrics as a unique identifier.
- Limit retention of biometric information.
- Limit storage of identifiable biometric data.
- Limit collection and storage of extraneous information, while including "opt-out" provisions for users.
- Enable anonymous enrollment and verification.
- Provide means of correcting and accessing biometric-related information.
- Limit system access.
- Use security tools and access policies to protect biometric information.
- Make provisions for third-party audits.
- Disclose the system purpose and objective.
- Disclose enrollment, verification, and identification processes.
- Disclose policies and protections in place to ensure privacy of biometric information.
- Disclose provisions for system termination.⁴²

29 · 22 BIOMETRIC AUTHENTICATION

In contrast to this view, Alterman argues that the deployment of biometric systems and the use of biometric data for identification and verification are ethically questionable because they always entail a violation of privacy and autonomy. Alterman finds “something disturbing about the generalized use of biometric identification apart from the standard data privacy issue.”⁴³ He maintains that biometric data “has inherent moral value”⁴⁴ but does not go so far as to argue against *any* deployment of biometric identification or verification systems. Rather, he maintains that they must be judiciously implemented and deployed only with due consideration to users’ privacy concerns.

29.6.7 Legal Issues. Within the legal system, as mentioned, the question arises as to the validity of using biometrics. As many consider this to be an invasion of privacy on many levels, the actual definition of this should be briefly touched on.

According to current law, there are four separate categories of “invasion of privacy”:

1. An intrusion into your private life—this is usually done by someone attempting to secretly learn something about you.
2. A public disclosure of private facts about yourself by someone other than yourself.
3. The use of your name, likeness, or both for monetary gain without your consent.
4. Any publication placing you in a false light.

With biometrics, the primary concern seems to stem from the first, and sometimes the second categories. Possibly these fears (and valid concerns, in light of recent government mandates and actions) may be alleviated by only storing a digital template, or a marker similar to those used in algorithms, rather than the actual image. TSA has been forced to use scans that show only a resemblance of the person in the scanner, and not an actual naked picture. The company they were using up to this time could not comply with their software, and thus lost the contract. With this, the possibility of the government being kept on a leash and not infringing on our Constitutional rights becomes a bit brighter.

29.7 RECENT TRENDS IN BIOMETRIC AUTHENTICATION

29.7.1 Government Advances in Biometric Authentication. Although private-sector organizations are increasingly adopting biometric technologies for their authentication needs, the government (public) sector has led investment in biometrics. The 2001 terrorist attacks on the World Trade Center and the Pentagon, and the ensuing USA PATRIOT Act, have encouraged increasing government commitment to biometric technologies. The Department of Defense (DoD), the Department of Homeland Security (DHS), the Immigration and Naturalization Service, and the Department of Transportation are the government agencies most involved in the deployment of biometrics technologies. The DoD’s Common Access Card program involves putting biometric technology on a smart ID card.⁴⁵ The US-VISIT program under the DHS is another government program that incorporates biometrics (including face and fingerprint) into a smart ID card. Another DHS program, the Transportation Worker Identity Credential, incorporates biometric information in an ID card.

29.7.2 Face Scanning at Airports and Casinos. After the 2001 terrorist attacks on the World Trade Center and the Pentagon, most of the nation’s airports

RECENT TRENDS IN BIOMETRIC AUTHENTICATION 29 · 23

moved to incorporate face-scanning technologies into their security systems. Most studies of the effectiveness of these systems, however, have revealed their high error rates and low accuracy rates.

Casinos utilize facial scanning systems to identify professional “advantage players” and cheats. Although this is largely unregulated in the United States, Canadian casinos must notify players regarding the use of such systems. Casinos share data on professionals and cheats. One firm has networked 125 casino surveillance operators in the United States, Canada, Puerto Rico, Aruba, and the Bahamas and provides real-time alerts and other information useful in identifying suspicious players. However, it is unclear as to how such systems may be affected by international law. Article 12 of the United Nations Universal Declaration of Human Rights guarantees that “[n]o one shall be subjected to arbitrary interference with his privacy, family, home or correspondence, nor to attacks upon his honor and reputation. Everyone has the right to the protection of the law against such interference or attacks.”⁴⁶ Whether this system amounts to arbitrary interference or an attack on one’s honor has not been addressed by the courts, but a case clearly could be made that covert use of such systems does so. A best practice for any surveillance system is informed consent. Organizations should clearly post notifications about the use of surveillance systems in order to protect themselves from legal challenges.

29.7.3 Increased Deployment in the Financial Industry. Usually slow to embrace new technologies, the financial industry is one of the leaders in the adoption of biometric authentication controls. Current deployments range from fingerprint scanners securing computer networks for brokers, to facial recognition systems at ATMs, to iris scanning for high-security access points. International Biometric Group projected that U.S. financial services firms would spend \$672 million in 2007 for various biometric deployments. One of the biggest deployments to date has been United Bankers’ Bancorporation (UBB) adoption of U.are.U, a fingerprint recognition system that allows UBB customers to automatically log onto UBB’s Website with finger scans versus passwords. UBB also adopted a fingerprint authentication system for its employees. Wells Fargo, Bloomberg Financial, and Janus Capital Management are other well-known financial firms that have adopted biometric authentication systems for employees and/or customers. Although some financial institutions have selected voice, iris, or facial-scan-based systems, most seem to be choosing finger-scan systems.

29.7.4 Biometrics in the Healthcare Industry. Spurred in part by new regulations that require healthcare institutions to ensure the privacy and security of patient records, healthcare companies have also been at the forefront in the adoption of biometric authentication. Among the major healthcare organizations that have moved to biometric authentication is the Mayo Clinic, which adopted a fingerprint ID system in 2002. The majority of healthcare institutions that have adopted biometric authentication systems have selected finger-scan ID systems. However, deployment of these systems in healthcare organizations has not met with the same success as seen in the financial services industry. Issues involving the potential transmission of illness via physical contact with the fingerprint scanner are not trivial. Additionally, error rates have been higher and accuracy rates much lower than expected. The major reason behind the high incidence of errors appears to be the particulars of the healthcare environment, especially the characteristics of the hands of the doctors, nurses, and other healthcare workers using these systems. Specifically, system performance appears to be undermined by the chronically dry hands of these workers, a condition resulting from

29 · 24 BIOMETRIC AUTHENTICATION

frequent hand washing and the use of alcohol-based hand sanitizers. Another problem has been the resistance to using the fingerprint technology by both nurses and doctors, who feel that it involves a privacy intrusion. Add to this the costs that are associated with tying a single sign on type of authentication to the core hospital system, and the problems escalate even higher. Cerner and EPIC seem to have the top programs in the healthcare industry, and both have been certified for all three Meaningful Use levels required by HIPAA and HITECH. Yet, these programs come at a hefty cost. Even barebones, a recent change for a hospital in Wyoming cost them \$4.5 million. While it was a change that needed to be done, it stretched their budgets to the very maximum, and several cuts were made to accommodate this.

29.7.5 Increased Deployment of Time and Attendance Systems. An increasing number of companies across many different industries are deploying biometric-based time-and-attendance systems. A shift from the past practice is in the increased use of biometric attendance and tracking systems for white-collar workers. Previously, the focus was on blue-collar factory workers. Although some employers are using the traditional hand-scanning systems, there appears to be a shift toward the use of finger-scanning time-and-attendance systems. This shift seems to be related to the more competitive pricing structure for the finger-scanning systems.

29.8 SUMMARY AND RECOMMENDATIONS. There is no universal “best” biometric authentication system. Each of the five leading biometric technologies carries specific advantages and disadvantages. Some biometric technologies are more appropriate for certain applications and environments than their counterparts. An organization in the midst of evaluating potential biometrics authentication implementation must recognize that there will be trade-offs in any selection, such as cost for accuracy, privacy versus user acceptance, and so on, and there are not yet any universal decision factors for selecting a particular biometric technology for a specific application. There is, however, substantial research into many of the advantages and disadvantages of biometrics. Exhibit 29.3 provides a summary comparison of the features of the five leading biometric technologies discussed in this chapter. The features, shown in the extreme left column, were excerpted from various researcher efforts, and the rankings represent an amalgam of the rankings found in the literature.

Although biometric authentication systems promise cost savings and higher levels of security for organizations, they are not a panacea. Many factors affect how well or poorly biometric authentication controls will perform in a given organizational environment. Included among these factors are the users, the administration, the environment, the infrastructure, the budget, the communication system, and the existing security needs. Although many biometric technologies are capable of operating as stand-alone systems, in reality their accuracy and performance levels would be greatly improved by combining them with more conventional authentication methods, such as passwords and keys. Such multifactor systems offer greater security and reliability.

In selecting a biometric authentication system and preparing for an implementation, organizations should focus on the user-technology interface and the conditions in the operational environment that may influence the technology’s performance. For example, the healthcare industry’s unreflective embrace of finger scan technology illustrates the dangers of failing to heed environmental realities. It is important that organizations consider not only the practical impediments to effective implementation but also the potential psychological impediments, such as user fears about the technology. Ethically, the organization also has the obligation to consider carefully the extent to which the

EXHIBIT 29.3 Comparison of Leading Biometric Technologies

	Finger Scan	Facial Scan	Hand Scan	Iris Scan	Voice Recognition
Accuracy	High	Low	Medium	Very high	Low to medium
Ease of Use	High	Medium	High	Low to medium	High
User Acceptance	Medium	High (overt) (covert)	High	Low to medium	High
Privacy Concerns	High	Very high (overt)	Medium	High	Very low
Cost	Low to medium	Low to medium	Medium	High	Low
Performance	High	Low	Medium	High	Low
Potential for Circumvention	Medium	High	Low to medium	High	High
Distinctiveness	High	Low	Medium	Very high	Low
Barriers to Universality	Worn ridges; hand or finger impairment	None	Hand impairment	Visual impairment	Speech impairment
Susceptibility to Changes in Biometric	Low to medium	Medium to high	Medium	Low	Low to medium
Susceptibility to Changes in the Environment	Low	High	Very low	Low	Medium to high
Error-Causing Factors	Age, trauma, degradation of prints	Lighting, contrast, pose, movement, expression	Hand injury or trauma, inability to place correctly	Positioning, eye angle, glasses, disease	Illness, age, quality of communication system, ambient noise
Mitigations for Potential Errors	Periodic re-enrollment, enrollment of multiple fingers	Frequent re-enrollment, multiple scans, controlled environment	Periodic re-enrollment, enrollment of both hands	Periodic re-enrollment, user training, enroll both irises	Periodic re-enrollment, control ambient noise

29 · 26 BIOMETRIC AUTHENTICATION

implementation of biometric authentication compromises the privacy rights of users. In making this assessment, management must take into account the possibility that the organization may be compelled to release employees' biometric-related information to government authorities.

Recent studies in biometrics include banks and other financial institutions having identified the use of voice biometrics over fingerprint identification as one of the best means to secure its client accounts and financial information. Voice biometrics compares various characteristics drawn from a client's voice such as inflection, pitch, dialect, accent and others, matching it with data previously captured, and securely stored. For this technology to work, however, it will require banks and other financial institutions to register their clients' voice patterns, correlate them to personal data, and place this into a database.⁴⁷ Other studies have reached agreements with major corporations to help within the healthcare industry, such as BIO-Key having executed an OEM agreement with Caradigm for use of BIO-Key's Identity and Access Management suite. Caradigm will offer this suite within its Single Sign On solution to help hospitals and other caregivers writing electronic prescriptions of controlled substances the ability to meet federal and state regulations.⁴⁸ Other areas of study within this field have covered such areas as a wearable biometric tattoo for identification purposes, vein recognition, and mobile platforms being used to promote its acceptance. In a 2011 report released by the Unisys Corporation, acceptance of these and other biometric technologies are gaining, but more slowly than anticipated.⁴⁹ Future use and acceptance may come about as technology advances, people become more accustomed to their use in everyday life, and legislation forces the use for security purposes.

29.9 FURTHER READING

- Kaine, A. K. "The Impact of Facial Recognition Systems on Business Practices within an Operational Setting." 25th International Conference Information Technology Interfaces, June 16–19, 2003, Cavtat, Croatia, pp. 315–320.
- Hamilton, D. P. "Workplace Security; Read My Lips: Are Biometric Systems the Security Solution of the Future?" *Wall Street Journal*, September 29, 2003.
- Jain, A. K., P. Flynn, and A. A. Ross, eds. *Handbook of Biometrics*. New York: Springer, 2007.
- NIST Biometrics Resource Center Website, www.nist.gov/biometrics-portal.cfm
- Ratha, N. K., Connell, J. H., & Bolle, R. M. "Enhanced Security and Privacy in Biometrics-Based Authentication Systems." *IBM Systems Journal* 40, No. 3 (2001): 614–634.
- Ratha, N., and R. Bolle, eds. *Automatic Fingerprint Recognition Systems*. New York: Springer, 2003.
- Ross, A. A., K. Nandakumar, and A. K. Jain. *Handbook of Multibiometrics*. New York: Springer, 2006.
- Vacca, J. R. *Biometric Technologies and Verification Systems*. Boston: Butterworth-Heinemann, 2007.

29.10 NOTES

1. J. Boroshok, "Pointing the Finger at Biometrics." SearchSecurity Website, January 14, 2005, http://searchsecurity.techtarget.com/originalContent/0289142,sid14_gci1044805,00.html
2. International Biometric Group, "Biometrics Market and Industry Report 2004–2008," 2005, accessed June 4, 2005, www.biometricgroup.com/reports/public/market_report.html (URL inactive).

NOTES 29 · 27

3. International Biometric Group, “Biometrics Market and Industry Report 2007–2012,” January 2007, accessed February 4, 2007, www.biometricgroup.com/reports/public/market_report.html (URL inactive).
4. J. Feldman, “2013 IT Spending Priorities Survey,” *Information Week*, May 6, 2013, <http://reports.informationweek.com/abstract/83/10497/IT-Business-Strategy/Research:-2013-IT-Spending-Priorities-Survey.html>
5. M. A. Davis, “2013 Strategic Security Survey,” *Information Week*, May 27, 2013, www.informationweek.com/security/management/2013-strategic-security-survey/240155501
6. J. Chirillo and S. Blaul, *Implementing Biometric Security* (Indianapolis, IN: John Wiley & Sons, 2003), p. 2.
7. S. Nanavati, M. Thieme, and R. Nanavati, *Biometrics: Identity Verification in a Networked World* (Hoboken, NJ: John Wiley & Sons, 2002), p. 11.
8. A. K. Jain, “Biometric Recognition: How Do I Know Who You Are,” *IEEE Symposia* (2004): 3–5.
9. R. M. Bolle, J. H. Connell, S. Pankanti, N. K. Ratha, and A. W. Senior, *Guide to Biometrics* (New York: Springer-Verlag, 2004).
10. Jain, “Biometric Recognition,” p. 3.
11. Bolle et al., *Guide to Biometrics*, p. 6.
12. Jain, “Biometric Recognition,” p. 3.
13. Chickowski, E. “Mobile Biometrics: Your Device Defines You”, October 3, 2012, www.informationweek.com/quickview/mobile-biometrics-your-device-defines-you/1675?queryText=biometrics&wc=4
14. Nanavati et al., *Biometrics*, p. 12.
15. Nanavati et al., *Biometrics*, p. 14
16. Nanavati et al., *Biometrics*, p. 14.
17. Jain, “Biometric Recognition,” p. 3.
18. J. D. Woodward, N. M. Orlans, and P. T. Higgins, *Biometrics* (New York: McGraw-Hill/Osborne, 2003), p. 37.
19. S. Liu and M. Silverman, “A Practical Guide to Biometric Security Technology,” *IT Pro* (January/February 2001): 20.
20. Nanavati et al., *Biometrics*, p. 20.
21. Chirillo and Blaul, *Implementing Biometric Security*, p. 21.
22. Chirillo and Blaul, *Implementing Biometric Security*, p. 24.
23. Bolle et al., *Guide to Biometrics*, p. 36.
24. N. Imparato, “Does Face Recognition Have a Future?” *Intelligent Enterprise* 5, No. 7 (April 2002): 20.
25. H. Bray, “Facial Recognition Software Precision May Be Years Away,” April 29, 2013, www.bostonglobe.com/business/2013/04/28/facial-recognition-technology-after-marathon-bombings/VN9gyIUqtps5EVfjI47P1J/story.html
26. Nanavati et al., *Biometrics*, p. 99.
27. Brandom, R. “We Know Who You Are: The Scary New Technology of Iris Scanners,” May 2, 2013, www.theverge.com/2013/5/2/4270352/theyre-already-watching-the-scary-new-technology-of-iris-scanners
28. Nanavati et al., *Biometrics*, p. 87.
29. Nanavati et al., *Biometrics*, p. 87.
30. Bolle et al., *Guide to Biometrics*, p. 40.

29 · 28 BIOMETRIC AUTHENTICATION

31. Nanavati et al., *Biometrics*, p. 27.
32. J. Ashbourn, *Practical Biometrics: From Aspiration to Implementation* (London: Springer-Verlag, 2003), p. 10.
33. Nanavati et al., *Biometrics*, p. 35.
34. Ashbourn, *Practical Biometrics*, p. 10.
35. T. Matsumoto, “Importance of Open Discussion on Adversarial Analyses for Mobile Security Technologies: A Case Study for User Identification,” May 14, 2002, <http://crypto.csail.mit.edu/classes/6.857/papers/gummy-slides.pdf>
36. R. J. Anderson, *Security Engineering: A Guide to Building Dependable Distributed Systems* (Hoboken, NJ: John Wiley & Sons, 2001).
37. Anderson, *Security Engineering*, p. 274.
38. S. Prabhakar, S. Pankanti, and A. K. Jain, “Biometric Recognition: Security and Privacy Concerns,” *IEEE Security & Privacy* (March/April 2003): 39.
39. Nanavati et al., *Biometrics*.
40. Woodward et al., *Biometrics*. p. 211.
41. Nanavati et al., *Biometrics*.
42. Nanavati et al., *Biometrics*.
43. A. Alterman, “‘A Piece of Yourself’: Ethical Issues in Biometric Identification,” *Ethics and Information Technology* 5, No. 3 (2003): 143.
44. Alterman, “‘A Piece of Yourself,’ ” p. 145.
45. www.digitalpersona.com/products/UPOS.php
46. United Nations, “Universal Declaration of Human Rights,” December 10, 1948, www.un.org/en/documents/udhr/index.shtml
47. R. King, “Biometric Research Note: Mobile Devices To Drive Bank Adoption Of Voice Biometrics,” *BiometricUpdate.com*, January 1, 2013, www.biometricupdate.com/201301/mobile-devices-to-drive-bank-adoption-of-voice-biometrics
48. A. Vrankuli, “Caradigm, BIO-key International Ink OEM Agreement,” *Biometric Update.com*, May 13, 2013, www.biometricupdate.com/201305/caradigm-bio-key-international-ink-oem-agreement
49. J. Trader, “Public Opinion on Biometric Technology Points to Wider Acceptance,” *M2SYS Blog on Biometric Technology*, November 8, 2011, <http://blog.m2sys.com/comments-on-recent-biometric-news-stories/public-opinion-on-biometric-technology-points-to-wider-acceptance>

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 30

E-COMMERCE AND WEB SERVER SAFEGUARDS

Robert Gezelter

30.1	INTRODUCTION	30·2	30.3.7	Avoiding Overreaction	30·21
			30.3.8	Appropriate Responses to Attacks	30·21
30.2	BUSINESS POLICIES AND STRATEGIES	30·3	30.3.9	Counter-Battery	30·22
30.2.1	Step 1: Define Information Security Concerns Specific to the Application	30·3	30.3.10	Hold Harmless	30·22
30.2.2	Step 2: Develop Security Service Options	30·5	30.4	RISK ANALYSIS	30·22
30.2.3	Step 3: Select Security Service Options Based on Requirements	30·7	30.4.1	Business Loss	30·23
30.2.4	Step 4: Ensures the Ongoing Attention to Changes in Technologies and Requirements	30·9	30.4.2	PR Image	30·23
30.2.5	Using the Security Services Framework	30·9	30.4.3	Loss of Customers/ Business	30·24
30.2.6	Framework Conclusion	30·17	30.4.4	Interruptions	30·24
			30.4.5	Proactive versus Reactive Threats	30·25
			30.4.6	Threat and Hazard Assessment	30·25
30.3	RULES OF ENGAGEMENT	30·18	30.5	OPERATIONAL REQUIREMENTS	30·25
30.3.1	Website-Specific Measures	30·18	30.5.1	Ubiquitous Internet Protocol Networking	30·26
30.3.2	Defining Attacks	30·19	30.5.2	Internal Partitions	30·27
30.3.3	Defining Protection	30·19	30.5.3	Critical Availability	30·27
30.3.4	Maintaining Privacy	30·20	30.5.4	Accessibility	30·27
30.3.5	Working with Law Enforcement	30·20	30.5.5	Applications Design	30·28
30.3.6	Accepting Losses	30·20	30.5.6	Provisioning	30·28
			30.5.7	Restrictions	30·29
			30.5.8	Multiple Security Domains	30·30
			30.5.9	What Needs to Be Exposed?	30·30
			30.5.10	Access Controls	30·31

30 · 2 E-COMMERCE AND WEB SERVER SAFEGUARDS

30.5.11	Site Maintenance	30·32	30.6.11	Emerging Technologies	30·41
30.5.12	Maintaining Site Integrity	30·32			
30.6	TECHNICAL ISSUES	30·33	30.7	ETHICAL AND LEGAL ISSUES	30·41
30.6.1	Inside/Outside	30·33	30.7.1	Liabilities	30·41
30.6.2	Hidden Subnets	30·34	30.7.2	Customer Monitoring, Privacy, and Disclosure	30·42
30.6.3	What Need Be Exposed?	30·34	30.7.3	Litigation	30·43
30.6.4	Multiple Security Domains	30·35	30.7.4	Application Service Providers	30·44
30.6.5	Compartmentalization	30·37			
30.6.6	Need to Access	30·38			
30.6.7	Accountability	30·39	30.8	SUMMARY	30·45
30.6.8	Read-Only File Security	30·39	30.9	FURTHER READING	30·46
30.6.9	Going Offline	30·40			
30.6.10	Auditing	30·40	30.10	NOTES	30·46

30.1 INTRODUCTION. Today, electronic commerce involves the entire enterprise. The most obvious e-commerce applications involve business transactions with outside customers on mobile devices and the World Wide Web (WWW or Web), yet they are merely the proverbial tip of the iceberg. The presence of e-commerce has become far more pervasive, often involving the entire logistical and financial supply chains that are the foundations of modern commerce. Even the smallest organizations now rely on the Web for access to services and information.

The explosive growth of mobile applications in the past few years has integrated electronic transactions deeply in the minute-to-minute fabric of daily life. Walk down the street in any urban area, mobile devices have become ubiquitous. Nations which lacked wired infrastructure have skipped a technological generation and gone directly to a telecommunications infrastructure based on mobile phone technologies.

The pervasive desire to improve efficiency often causes a convergence between the systems supporting conventional operations with those supporting an organization's online business. It is thus common for internal systems at bricks-and-mortar stores to utilize the same back-office systems as are used by Web customers. It is also common for kiosks and cash registers to use wireless local area networks to establish connections back to internal systems. These interconnections can provide intruders with access directly into the heart of the enterprise.

The TJX case, which came to public attention in the beginning of 2007, was one of a series of large-scale compromises of electronically stored information on back-office and e-commerce systems. Most notably, the TJX case appears to have started with an insufficiently secured corporate network and the associated back-office systems, not a Website penetration. This breach escalated into a security breach of corporate data systems. It has been reported that at least 94 million credit cards were compromised.¹ On November 30, 2007, it was reported that TJX, the parent organization of stores including TJ Maxx and Marshall's, agreed to settle bank claims related to VISA cards for US\$ 40.9M.²

E-commerce has now come of age, giving rise to fiduciary risks that are important to senior management and to the board of directors. The security of data networks, both those used by customers and those used internally, now has reached the level where it significantly affects the bottom line. TJX has suffered both monetarily and in public relations, with stories concerning the details of this case appearing in the *Wall Street*

BUSINESS POLICIES AND STRATEGIES 30 · 3

Journal, the *New York Times*, *Business Week*, and many industry trade publications. Data security is no longer an abstract issue of concern only to technology personnel. The legal settlements are far in excess of the costs directly associated with curing the technical problem.

Protecting e-commerce information requires a multifaceted approach, involving business policies and strategies as well as the technical issues more familiar to information security professionals.

Throughout the enterprise, people and information are physically safeguarded. Even the smallest organizations have a locked door and a receptionist to keep outsiders from entering the premises. The larger the organization, the more elaborate the precautions needed. Small businesses have simple locked doors; larger enterprises often have many levels of security, including electronic locks, security guards, and additional levels of receptionists. Companies also jealously guard the privacy of their executive conversations and research projects. Despite these norms, it is not unusual to find that information security practices are weaker than physical security measures. Connection to the Internet (and within the company, to the intranet) worsens the problem by greatly increasing the risk and decreasing the difficulty, of attacks. Ubiquitous mobile devices only serve to compound the difficulties.

30.2 BUSINESS POLICIES AND STRATEGIES. In the complex world of e-commerce security, best practices are constantly evolving. New protocols and products are announced regularly. Before the Internet explosion, most companies rarely shared their data and their propriety applications with any external entities, and information security was not a high priority. Now companies taking advantage of e-commerce need sound security architectures for virtually all applications. Effective information security has become a major business issue. This chapter provides a flexible framework for building secure e-commerce applications and assistance in identifying the appropriate and required security services. The theoretical examples shown are included to facilitate the reader's understanding of the framework in a business-to-customer (B2C) and business-to-business (B2B) environment.

A reasonable framework for e-commerce security is one that:

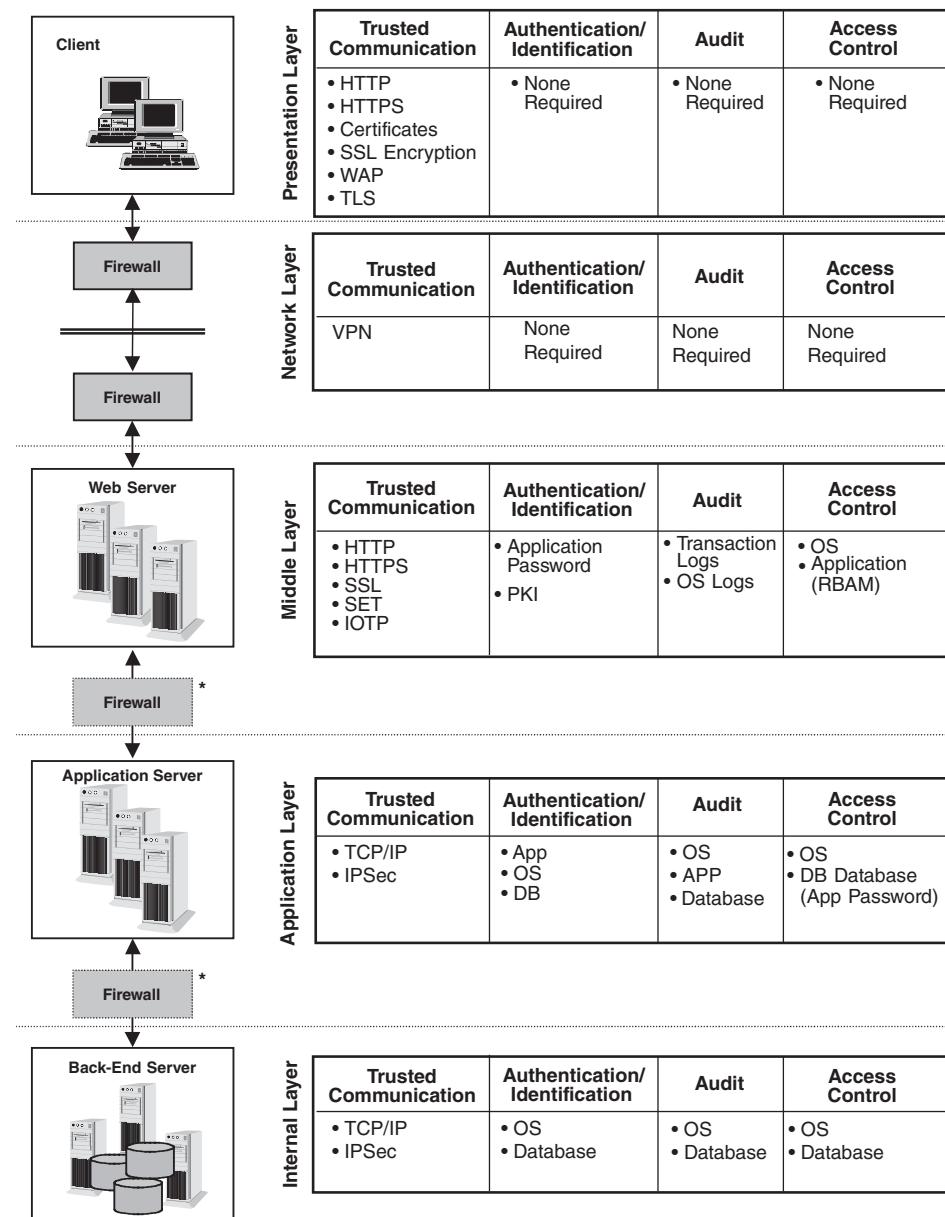
1. Defines information security concerns specific to the application.
2. Defines the security services needed to address the security concerns.
3. Selects security services based on a cost-benefit analysis and risk versus reward issues.
4. Ensures the ongoing attention to changes in technologies and requirements as both threats and application requirements change.

This four-step approach is recommended to define the security services selection and decision-making processes.

30.2.1 Step 1: Define Information Security Concerns Specific to the Application. The first step is to define or develop the application architecture and the data classification involved in each transaction. This step considers how the application will function. As a general rule, if security issues are defined in terms of the impact on the business, it will be easier to discuss with management and easier to define security requirements.

30 · 4 E-COMMERCE AND WEB SERVER SAFEGUARDS

The recommended approach is to develop a transactional follow-the-flow diagram that tracks transactions and data types through the various servers and networks. This should be a functional and logical view of how the application is going to work—that is, how transactions will occur, what systems will participate in the transaction management, and where these systems will support the business objectives and the organization's product value chain. Data sources and data interfaces need to be identified, and the information processed needs to be classified. In this way a complete transactional flow can be represented. (See Exhibit 30.1.)



*Firewall should be between the web layer/application layer or the application/back-end layers depending on the company's architecture.

EXHIBIT 30.1 Trust Levels for B2C Security Services

BUSINESS POLICIES AND STRATEGIES 30 · 5

Common tiered architecture points include:

- **Clients.** These may be PCs, thin clients (devices that use shared applications from a server and have small amounts of memory), personal digital assistants (PDAs), and wireless application protocol (WAP) telephones.
- **Servers.** These may include World Wide Web, application, database, and middleware processors, as well as back-end servers and legacy systems.
- **Network devices.** Switches, routers firewalls, NICs, codecs, modems, and internal and external hosting sites.
- **Network spaces.** Network demilitarized zones (DMZs), intranets, extranets, and the Internet.

It is important at this step of the process to identify the criticality of the application to the business and the overriding security concerns: transactional confidentiality, transactional integrity, or transactional availability. Defining these security issues will help justify the security services selected to protect the system. The more completely the architecture can be described, the more thoroughly the information can be protected via security services.

30.2.2 Step 2: Develop Security Service Options. The second step considers the security services alternatives for each architecture component and the data involved in each transaction. Each architectural component and data point should be analyzed and possible security services defined for each. Cost and feasibility should not be considered to any great degree at this stage. The objective is to form a complete list of security service options with all alternatives considered. The process should be comparable with, or use the same techniques as, brainstorming. All ideas, even if impractical or farfetched, should be included.

Decisions should not be made during this step; that process is reserved for Step 3.

The information security organization provides services to an enterprise. The services provided by information security organizations vary from company to company. Several factors will determine the required services, but the most significant considerations include:

- Industry factors
- The company's risk appetite
- Maturity of the security function
- Organizational approach (centralized or decentralized)
- Impact of past security incidents
- Internal organizational factors
- Political factors
- Regulatory factors
- Perceived strategic value of information security

Several factors contribute to the services that information security organizations provide. "Security services" are defined as safeguards and control measures to protect the confidentiality, integrity, and accountability of information and computing resources. Security services that are required to secure e-commerce transactions need to be based on the business requirements and on the willingness to assume or reduce the risk of the

30 · 6 E-COMMERCE AND WEB SERVER SAFEGUARDS

information being compromised. Information security professionals can be subject-matter experts, but they are rarely equipped to make the business decisions required to select the necessary services. Twelve security services that are critical for successful e-commerce security have been identified:

1. **Policy and procedures** are a security service that defines the amount of information security that the organization requires and how it will be implemented. Effective policy and procedures will dovetail with system strategy, development, implementation, and operation. Each organization will have different policies and procedures; best practice dictates that organizations have policies and procedures based on the risk the organization is willing to take with its information. At a minimum, organizations should have a high-level policy that dictates the proper use of information assets and the ramifications of misuse.
2. **Confidentiality and encryption** are a security service that secures data while they are stored or in transit from one machine to another. A number of encryption schemes and products exist; each organization needs to identify those products that best integrate with the application being deployed. For a discussion of cryptography, see Chapters 7 and 37 in this *Handbook*.
3. **Authentication and identification** are a security service that differentiates users and verifies that they are who they claim to be. Typically, passwords are used, but stronger methods include tokens, smart cards, and biometrics. These stronger methods verify what you have (e.g., token) or who you are (e.g., biometrics), not just what you know (password). Two-factor authentication combines two of these three methods and is referred to as strong authentication. For more on this subject, see Chapter 28 in this *Handbook*.
4. **Authorization** determines what access privileges a user requires within the system. Access includes data, operating system, transactional functions, and processes. Access should be approved by management who own or understand the system before access is granted. Authorized users should be able to access only the information they require for their jobs.
5. **Authenticity** is a security service that validates a transaction and binds the transaction to a single accountable person or entity. Also called nonrepudiation, authenticity ensures that a person cannot dispute the details of a transaction. This is especially useful for contract and legal purposes.
6. **Monitoring and audit** provide an electronic trail for a historical record of the transaction. Audit logs consist of operating system logs, application transaction logs, database logs, and network traffic logs. Monitoring these logs for unauthorized events is considered a best practice.
7. **Access controls and intrusion detection** are technical, physical, and administrative services that prevent unauthorized access to hardware, software, or information. Data are protected from alteration, theft, or destruction. Access controls are preventive—stopping unauthorized access from occurring. Intrusion detection catches unauthorized access after it has occurred, so that damage can be minimized and access cut off. These controls are especially necessary when confidential or critical information is being processed.
8. **Trusted communication** is a security service that assures that communication is secure. In most instances involving the Internet, this means that the communication will be encrypted. In the past, communication was trusted because it

BUSINESS POLICIES AND STRATEGIES 30 · 7

was contained within an organization's perimeter. Communication is currently ubiquitous and can come from almost anywhere, including extranets and the Internet.

9. **Antivirus** is a security service that prevents, detects, and cleans viruses, Trojan horse programs, and other malware.
10. **System integrity controls** are security services that help to assure that the system has not been altered or tampered with by unauthorized access.
11. **Data retention and disposal** are a security service that keeps required information archived, or deletes data when they are no longer required. Availability of retained data is critical when an emergency exists. This is true whether the problem is a systems outage or a legal process, whether caused by a natural disaster or by a terrorist attack (e.g., September 11, 2001).
12. **Data classification** is a security service that identifies the sensitivity and confidentiality of information. The service provides guides for information labeling, and for protection during the information's life.

Once an e-commerce application has been identified, the team must identify the security issues with that specific application and the necessary security services. Not all of the services will be relevant, but using a complete list and excluding those that are not required will assure a comprehensive assessment of requirements, with appropriate security built into the system's development. In fact, management can reconcile the services accepted with their level of risk acceptance.

30.2.3 Step 3: Select Security Service Options Based on Requirements. The third step uses classical cost-benefit and risk-management analysis techniques to make a final selection of security service options. However, we recommend that all options identified in Step 3 be distributed along a continuum, such as shown in Exhibit 30.2, so that they can be viewed together, and compared.

Gauging and comparing the level of security for each security service and the data within the transaction will facilitate the decision process. Feasible alternatives can then be identified and the best solution selected based on the requirements. The most significant element to consider is the relative reduction in risk of each option, compared with the other alternatives. The cost-benefit analysis is based on the risk versus reward issues. The effectiveness information is very useful in a cost-benefit model.

Four additional concepts drive the security service option selection:

1. Implementation risk or feasibility
2. Cost to implement and support
3. Effectiveness in increasing control, thereby reducing risk
4. Data classification

Implementation risk considers the feasibility of implementing the security service option. Some security systems are difficult to implement due to factors such as product maturity, scalability, complexity, and supportability. Other factors to consider include skills available, legal issues, integration required, capabilities, prior experience, and limitations of the technology.

Cost to implement and support measures the costs of hardware and software implementation, support, and administration. Consideration of administration issues is

30 · 8 E-COMMERCE AND WEB SERVER SAFEGUARDS

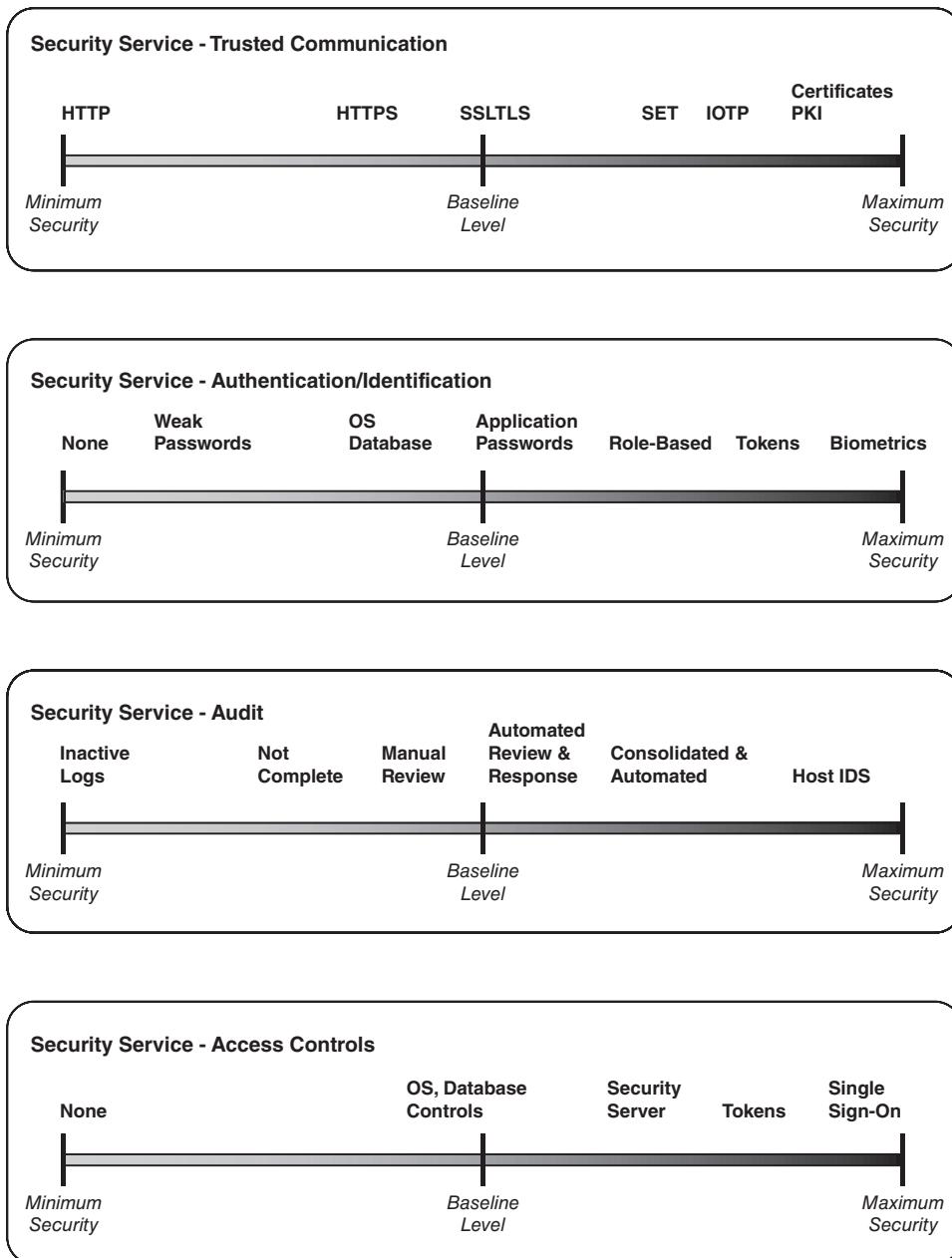


EXHIBIT 30.2 Continuum of Options

especially critical because high-level support of the security service is vital to an organization's success.

Effectiveness measures the reduction of risk proposed by a security service option once it is in production. Risk can be defined as the impact and likelihood of a negative event occurring after mitigating strategies have been implemented. An example of a negative event is the theft of credit card numbers from a business's database. Such an event causes not only possible losses to consumers but also negative public relations

BUSINESS POLICIES AND STRATEGIES 30 · 9

that may impact future business. Effective security service options reduce the risk of a negative event occurring.

Data classification measures the sensitivity and confidentiality of the information being processed. Data must be classified and protected from misuse, disclosure, theft, or destruction, regardless of storage format, throughout their life (from creation through destruction). Usually, the originator of the information is considered to be the owner of the information and is responsible for classification, identification, and labeling. The more sensitive and confidential the information, the more information security measures will be required to safeguard and protect it.

30.2.4 Step 4: Ensures the Ongoing Attention to Changes in Technologies and Requirements. The only constant in this analysis is the need to evolve and to address ever-increasing threats and technologies. Whatever security approaches the preceding steps identify, they must always be considered in the context of the continuing need to update the selected approaches. Changes will be inevitable, whether they arrive from compliance, regulation, technological advances, or new threats.

30.2.5 Using the Security Services Framework. The next two sections are examples to demonstrate the power of the security services methodology. The first example is a B2C model; the business could be any direct-selling application. The second example is a B2B model. Both businesses take advantage of the Internet to improve their product value chain. The examples are a short demonstration of the security services methodology, neither complete nor representative of any particular application or business.

30.2.5.1 Business-to-Customer Security Services. The B2C company desires to contact customers directly through the Internet, and allow them to enter their information into the application. These assumptions are made to prepare this B2C system example:

- Internet-facing business
- Major transactions supported
- External customer-based system, specifically excluding support, administration, and operations
- Business-critical application
- Highly sensitive data classification
- Three-tiered architecture
- Untrusted clients, because anyone on the Internet can be a customer

Five layers must be secured:

1. The *presentation layer* is the customer interface, what the client sees or hears using the Web device. The client is the customer's untrusted PC or other device. The security requirements at this level are minimal because the company will normally not dictate the security of the customer. The identification and authentication done at the presentation level are those controls associated with access to the device. The proliferation of appliances in the client space (e.g., traditional

30 · 10 E-COMMERCE AND WEB SERVER SAFEGUARDS

PCs, thin desktops, and PDAs) makes it difficult to establish a uniform access control procedure at this level. Although this layer is the terminal endpoint of the secure connection, it is not inherently trustworthy, as has been illustrated by all too many incidents involving public computers in cafes, hotels, and other establishments.

2. The *network layer* is the communication connection between the business and the customer. The client, or customer, uses an Internet connection to access the B2C Web applications. The security requirements are minimal, but sensitive, and confidential traffic will need to be encrypted.
3. The *middle layer* is the Web server that connects to the client's browser and can forward and receive information. The Web server supports the application by being an intermediary between the business and the customer. The Web server needs to be very secure. Theft, tampering, and fraudulent use of information needs to be prevented. Denial of service and Website defacement are also common risks that need to be prevented in the middle layer.
4. The *application layer* is where the information is processed. The application serves as an intermediary between the customer requests and the fulfillment systems internal to the business. In some examples, the application server and database server are the same because both the application and database reside on the same server. However, they could reside within the Web server in other cases.
5. The *internal layer* is comprised of the business's legacy systems and databases that support customer servicing. Back-end servers house the supporting application, including order processing, accounts receivable, inventory, distribution, and other systems.

For each of these five levels, we need four security services:

1. Trusted communications
2. Authentication/identification
3. Audit
4. Access controls

Step 1: Define Information Security Concerns Specific to the Application. Defining security issues will be particular to the system being implemented. To understand the risk of the system, the best starting place is with the business risk; then defining risks at each element of the architecture.

Business Risk

- The application needs high availability, because customers will demand contact at off-hours and on weekends.
- The Web pages need to be secure from tampering and cyber-vandalism, because the business is concerned about the loss of customer confidence as a result of negative publicity.
- Customers must be satisfied with the level of service.
- The system will process customer credit card information and will be subject to the privacy regulations of the Gramm-Leach-Bliley Act, 15 USC §§ 6801–09.

BUSINESS POLICIES AND STRATEGIES 30 · 11

Technology Concerns

Of the five architectural layers in this example, four will need to be secured:

1. The *presentation layer* will not be secured or trusted. Communications between the client and the Web server will be encrypted at the *network layer*.
2. The *network layer* will need to filter unwanted traffic, to prevent denial of service (DoS) attacks and to monitor for possible intrusions.
3. The *middle layer* will need to prevent unauthorized access, be tamper-proof, contain effective monitoring, and support efficient and effective processing.
4. The *application layer* will need to prevent unauthorized access, support timely processing of transactions, provide effective audit trails, and process confidential information.
5. The *internal layer* will need to prevent unauthorized access, especially through Internet connections, and to protect confidential information during transmission, processing, and storage.

Step 2: Develop Security Services Options. The four security services reviewed in this example are the most critical in an e-commerce environment. Other services such as nonrepudiation and data classification are important but not included in order to simplify the example. Services elected are:

- Trusted communication
- Authentication and identification
- Monitoring and auditing
- Access control

Many security services options are available for the B2C case, with more products and protocols on the horizon. There are five architectural layers for each of the services defined in Step 1.

1. **Presentation layer.** Several different options can be selected for trusted communication. While the most common application-level protocol is the Hypertext Transfer Protocol (HTTP); the ever increasing need for privacy and assurance of integrity argues for the use of its more secure encrypted sibling, Transport Layer Security,³ the successor to Secure Socket Layer (SSL) protocol,⁴ in conjunction with certificates in a Public Key Infrastructure (PKI), to authenticate and cryptographically secure the communications pathway. The use, or possible use, of wireless connections somewhere in the path only increases the importance of secured, encrypted paths. Because the client is untrusted, the client's authentication, audit, or access control methods cannot be relied on.
2. **Network layer.** Virtual private networks (VPNs) are considered best practice for secure network layer communication. Firewalls are effective devices for securing network communication. The client may have a personal firewall. If the client is in a large corporation, there is a significant likelihood that a firewall will intervene in communications. If the client is using a home or laptop computer, then a personal firewall may protect traffic. There will also be a firewall on the B2C company side of the Internet.

30 · 12 E-COMMERCE AND WEB SERVER SAFEGUARDS

3. **Middle layer.** The Web server and the application server security requirements are significant. Unauthorized access to the Web server can result in Website defacement by hackers who change the Web data. More important, access to the Web or application server can lead to theft, manipulation, or deletion of customer or proprietary data. Communication between the Web server and the client needs to be secure in e-commerce transactions. HTTP is the most common form of browser communication. In 2000, it was reported that over 33 percent of credit card transactions were using unsecured HTTP.⁵ In 2007, VISA reported that compliance with the Payment Card Industry Data Security Standard (PCI DSS) had improved to 77 percent of the largest merchants in the United States,⁶ still far from universal. This lack of encryption is a violation of both the merchant's agreement and the PCI, but episodes continue to occur. In the same vein, it is not unusual for organizations to misuse encryption, whether it involves self-signed, expired, or not generally recognized certificates. (See Chapter 37 in this *Handbook*.) HTTPS (HTTP over TLS) is the most common secure protocol solution, but the encryption key length and contents are critical: The larger the key, the more secure the transaction is from brute-force decryption (see Chapter 7). Digital certificates in a PKI and digital signatures are not as common, but they are more effective forms of security.
4. **Application layer.** The application server provides the main processing for the Web server. This layer may include transaction processing and database applications; it needs to be secure to prevent erroneous or fraudulent processing. Depending on the sensitivity of the information processed, the data may need to be encrypted. Interception and the unauthorized manipulation of data are the greatest risks in the application layer.
5. **Internal layer.** In a typical example, the internal layer is secured by a firewall protecting the external-facing system. This firewall helps the B2C company protect itself from Internet intrusions. Database and operating system passwords and access control measures are required.

Management can decide which security capability is required and at what level. This format can be repeated to discuss security services at all levels and all systems, not just e-commerce-related systems.

Step 3: Select Security Service Options Based on Requirements. In Step 3 the B2C company can analyze and select the security services that best meet its legal, business, and information security requirements. There are four stages required for this analysis:

1. Implementation risk or feasibility
2. Cost to implement and support
3. Effectiveness in increasing control, thereby reducing risk
4. Data classification

Implementation risk is a function of the organization's ability to effectively roll out the technology. In this example, we assume that implementation risk is low and resources are readily available to implement the technology.

Costs to implement and support are paramount to the decision-making process. Both costs need to be considered together. Unfortunately, the cost to support is difficult

BUSINESS POLICIES AND STRATEGIES 30 · 13

to quantify and easily overlooked. In this example, resources are available to both implement and support the technology.

Effectiveness in increasing control is an integral part of the benefit and risk-management decisions. Each component needs to be considered in order to determine cross-benefits where controls overlap and supplement other controls. In this case, the control increase is understood and supported by management.

Data classification is the foundation for requirements. It will help drive the cost-benefit discussions because it captures the value of the information to the underlying business. In this example, the data are considered significant enough to warrant additional security measures to safeguard the data against misuse, theft, and loss.

There are many technology decisions required to secure the example environment. Management can use this approach to plot the security levels required by the system. For example, for system audit services, in order of minimal service to maximum:

- The minimal level of security is to have systems with a limited number of system events logged. For example, the default level of logging from the manufacturer is used but does not contain all of the required information. The logs are not reviewed but are available for forensics in the event of a security incident.
- A higher level of security is afforded with a log that records more activities based on requirements and not, for example, the manufacturer's default level. As in the minimal level of security, the activities are logged and available to support forensics but are not reviewed. In this case, more types of information are recorded in the system log, but it still may not contain all that is required.
- A sufficient log is kept on each server and is manually monitored for anomalies and potential security events.
- The log is automatically reviewed by software on each server.
- System logs are consolidated onto a centralized security server. Data from the system logs are transmitted to the centralized security server, and software is then used to scan the logs for specific events that require attention. Events such as attempts to gain escalated privileges to root or administrative access can be flagged for manual review.
- The maximum service level is a host-based intrusion detection system (IDS) used to scan the system logs for anomalies and possible security events. Once detected, action needs to be taken to resolve the intrusion. The procedure should include processes such as notification, escalation, and automated defensive response.

30.2.5.2 Business-to-Business Security Services. The second case study uses the security services framework in a B2B example. Following is a theoretical discussion of how the framework can be applied to B2B e-commerce security. These assumptions may be made in this B2B system example:

- Internet-facing.
- Supports major transactions.
- Descriptions will be external and customer based (excluding support, administration, and operations security services).
- Trusted communication is required.

30 · 14 E-COMMERCE AND WEB SERVER SAFEGUARDS

- Three-tier architecture.
- Untrusted client.
- Business-critical application.
- Data are classified as highly sensitive.

There are five layers in this example that need to be secured:

1. The *presentation layer* is the customer interface and is what the client sees or hears using the Web device. The client is the customer's untrusted PC, but more security constraints can be applied because the business can dictate enhanced security. As noted previously, the security of the presentation layer is complicated by the wide range of potential client devices that may be employed.
2. The *application layer* is where the information is processed. The application serves as an intermediary between the business customer's requests and the fulfillment systems internal to the business (the back-end server). The application server is the supporting server and database.
3. The *customer internal layer* is the interface between the application server supporting the system at the customer's business location, and the customer's own internal legacy applications and systems.
4. The *network layer* is the communication connection between the business and another business. The Internet is used to connect the two businesses. Sensitive and confidential traffic will need to be encrypted. Best practice is to have the traffic further secured using a firewall.
5. The *internal layer* is the business's legacy systems that support customer servicing. The back-end server houses the supporting systems, including order processing, accounts receivable, inventory, distribution, and other systems.

The four security services are:

1. Trusted communications
2. Authentication/identification
3. Audit
4. Access controls

Step 1: Define Information Security Concerns Specific to the Application. Defining security issues will be particular to the system being implemented. To understand the risk of the system, it is best to start with the business risk; then define risk at each element of the architecture.

Business Risk

- Communication between application servers needs to be very secure. Data must not be tampered with, stolen, or misrouted.
- Availability is critical during normal business hours.
- Cost savings realized by switching from electronic data interchange (EDI) to the Internet is substantial, and will more than cover the costs of the system.

BUSINESS POLICIES AND STRATEGIES 30 · 15

Technology Concerns

There are six architectural layers in this example, five of which need to be secured:

1. The *presentation layer* will not be secured or trusted. The communication between the client and the customer application is trusted because it uses the customer's private network.
2. The *application server* will need to be secure. Traffic between the two application servers will need to be encrypted. The application server is inside the customer's network and demonstrates a potentially high, and perhaps unnecessary degree of trust between the two companies (see Section 30.6.5).
3. The *customer's internal layer* will be secured by the customer.
4. The *network layer* needs to filter out traffic that is not required, prevent DoS attacks, and monitor for possible intrusions. Two firewalls are shown: one to protect the client and the other to protect the B2B company.
5. The *application layer* will need to prevent unauthorized access, support timely processing of transactions, provide effective audit trails, and process confidential information.
6. The *internal layer* will need to prevent unauthorized access (especially through Internet connections) and protect confidential information during transmission, processing, and storage.

Step 2: Develop Security Services Options. There are four security services reviewed in this example. Others could have been included, such as authenticity, nonrepudiation, and confidentiality, but they have been excluded to simplify this example. Elected security services include:

1. Trusted communication
2. Authentication/identification
3. Audit
4. Access control

Many security services options are available for B2B environments.

- **Presentation layer.** Several different options can be selected for communication. HTTP is the most common. The communications between the presentation layer residing on the client device and the application server, in this example, are internal to the customer's trusted internal network and will be secured by the customer. Past episodes involving this scenario include the previously mentioned TJX affair and the episode involving Hannaford.⁷ Recently, Bashas', a chain of grocery stores in Arizona, was compromised by malware on their Point of Sale network.⁸ These compromises of internal networks argue for the use of HTTPS (HTTP over TLS) as a presentation layer, even on "internal" networks where information transmitted may include personally identifiable information (PII) in some form. As a matter of safety, when in doubt it is far safer to run an encrypted infrastructure as a guard against exposing PII than it is to discover at a later date that some data in certain circumstances is protected PII. Encrypting unsensitive data does accrue excessive costs; failing to encrypt even a small amount of PII could lead to a significant liability.

30 · 16 E-COMMERCE AND WEB SERVER SAFEGUARDS

- **Application layer.** Communication between the two application servers needs to be secure. The easiest and most secure method of peer-to-peer communication is via a VPN.
- **Customer internal.** Communications between the customer's application server and the customer's back-end server are internal to the customer's trusted internal network and will be secured by the customer.
- **Network layer.** It is common in a B2B environment that a trusted network is created via a VPN. The firewalls will probably participate in these communications, but hardware solutions are also possible.
- **Application layer.** The application server is at both the customer and B2B company sites. VPN is the most secure communication method. The application server also needs to communicate with the internal layer, and this traffic should be encrypted as well.
- **Internal layer.** The internal layer may be secured with another firewall from the external-facing system. This firewall helps the B2B company to protect itself from intrusions and unauthorized access. In this example, a firewall is not assumed, so the external firewall and DMZ need to be very secure.

Intrusion detection, log reading, and other devices can easily be added and discussed with management. This format can be repeated to discuss security services at all levels and all systems, not just e-commerce-related systems.

Step 3: Develop Security Service Options. In Step 3, the B2B company can analyze and select the security services that best meet its legal, business, and information security requirements. The biggest difference between B2C and B2B systems is that the B2C system assumes no level of trust. The B2B system assumes trust to some degree, but additional coordination and interface with the B2B customer or partner is required. Trustworthiness is not an either/or question and should not be presumed (see Section 30.6.4). This coordination and interoperability must not be underestimated, because they may prove difficult and expensive to resolve. There are four stages required for this analysis:

1. Implementation risk or feasibility
2. Cost to implement and support
3. Effectiveness in increasing control, thereby reducing risk
4. Data classification

Implementation risk is a function of the organization's ability to effectively roll out the technology. In this example, we assume that implementation risk is low and resources are readily available to implement the technology.

Cost to implement and support are paramount to the decision-making process. Both businesses' costs need to be considered. Unfortunately, the cost to support is difficult to quantify and easily overlooked. In this example, resources are available both to implement and to support the technology.

Effectiveness in increasing control is an integral part of the benefit and risk-management analysis. Each security component needs to be considered in order to determine cross-benefits where controls overlap and supplement others. In this example, increased levels of control are understood and supported by management.

BUSINESS POLICIES AND STRATEGIES 30 · 17

Data classification is the foundation for requirements and will help drive the cost-benefit discussions because it captures the value of the information to the underlying business. In this example, the data are considered significant enough to warrant additional security measures to safeguard the data against misuse, theft, and loss.

Each security service can be defined along a continuum, with implementation risk, cost, and data classification all considered. Management can use this chart to plot the security levels required by the system. This example outlines the effectiveness of security services options relative to other protocols or products. Each organization should develop its own continuums and provide guidance to Web developers and application programmers as to the correct uses and standard settings of the security services. For example, for authentication/identification services, in order of minimal service to maximum:

- The minimal level of security is to have no passwords.
- Weak passwords (e.g., easy to guess, shared, poor construction) are better than no passwords but still provide only a minimal level of security.
- Operating system or database level passwords usually allow too much access to the system but can be effectively managed.
- Application passwords are difficult to manage but can be used to restrict data access to a greater degree.
- Role-based access distinguishes users by their need to know to support their job function. Roles are established and users are grouped by their required function.
- Tokens are given to users and provide for two-part authentication. Passwords and tokens are combined for strong authentication.
- Biometrics are means to validate the person claiming to be the user via fingerprints, retina scans, or other unique body function.

For more information on identification and authentication, see Chapters 28 and 29 in this *Handbook*.

30.2.6 Framework Conclusion. Internet e-commerce has changed the way corporations conduct business with their customers, vendors, suppliers, and business units. The B2B and B2C sectors will only continue to grow. Despite security concerns, the acceleration toward increased use of the Internet as a sales, logistics, and marketing channel continues. The challenge for information security professionals' is to keep pace with this change from a security perspective, but not to impede progress. Another equal challenge is that the products that secure the Internet are new and not fully functional or mature. The products will improve, but meanwhile, existing products must be implemented, and later retrofitted, with improved and more secure security services. This changing environment, including the introduction of ever more serious and sophisticated threats, remains difficult to secure.

The processes described in this section will allow the security practitioner to provide business units with a powerful tool to communicate, select, and implement information security services. Three steps were described and demonstrated with two examples. The process supports decision making. Decisions can be made and readily documented to demonstrate cost effectiveness of the security selections. The risk of specific decisions can be discussed and accepted by management. The trade-offs between cost and benefit can be calculated and discussed. Therefore, it becomes critical that alternatives be

30 · 18 E-COMMERCE AND WEB SERVER SAFEGUARDS

reviewed and good decisions made. The processes supporting these decisions need to be efficient and quickly applied. The information security services approach will allow companies to implement security at a practical pace. Services not selected are easily seen. The risk of not selecting specific security services needs to be accepted by management.

30.3 RULES OF ENGAGEMENT. The Web is a rapidly evolving, complex environment. Dealing with customers electronically is a challenge. Web-related security matters raise many sensitive security and privacy issues. Attacks against Websites are not private incidents; they are often public and must always be taken seriously. Correctly differentiating “false alarms” from real attacks is an ongoing challenge with little hope of clean solutions. Distinguishing a popularity spike from a Distributed Denial of Service (DDoS) attack is imprecise and error prone. Indeed, a single traffic surge could be both simultaneously. The Victoria’s Secret online lingerie show in February 1998 was popular beyond even the most optimistic expectations of its creators, and the volume of visitors caused severe problems. Obviously, the thousands of people were not attacking the site; they were merely a virtual mob attempting to access the same site at the same time. Similar episodes have occurred when sites were described as interesting on Usenet newsgroups or on social networking sites (e.g., YouTube, Facebook, Twitter, and others). Request tsunamis often occur with little or no warning. Physical mobs are limited by transportation, timing, and costs; virtual mobs are solely limited by how many can attempt to access a resource simultaneously, from locations throughout the world. A recent example was the debut of a limited Missoni product set by Target on September 13, 2011, which generated a demand surge that rendered Target’s Website effectively inaccessible (93 percent of requests were refused) for hours.^{9,10}

30.3.1 Website-Specific Measures. Protecting a Website means ensuring that the site and its functions are available 24 hours a day, seven days a week, and 365 days a year. It also means ensuring that the information exchanged with the site is accurate and secure.

This section focuses on issues specific to Web interactions with customers as well as supply and distribution chains. Practically speaking, the Website is an important, if not the most important, component of an organization’s interface with the outside world.

Website protection lies at the intersection of technology, strategy, operations, customer relations, and business management. Website availability and integrity directly affect the main streams of cash flow and commerce: an organization’s customers, production chains, and supply chains.

Availability is the cornerstone of all Web-related strategies. Idle times have become progressively more sparse. Depending on the business and its markets, there may be some periods of lower activity. In the financial trading community, there remain only a few small windows during a 24-hour period when updates and maintenance can be performed. As global business becomes the norm, customers, suppliers, and distributors increasingly expect information, and the ability to effect transactions, at any time of the day or night, even from modest-size enterprises. On the Internet, “nobody knows that you are a dog” also means “nobody knows that you are *not* a large company.” The playing field has indeed been leveled, but it was not uniformly raised or lowered. Expectations have increased while capital and operating expenses have dramatically dropped.

Causation is unrelated to impact. The overwhelming majority of Web outages are caused by unglamorous problems. High-profile, deliberate attacks are much less

RULES OF ENGAGEMENT 30 · 19

frequent than equipment and personnel failures. The effect on the business organization is indistinguishable. Having a low profile is no defense against random scanning attack or someone looking for a springboard from which to attack a third party.

Additionally, while an individual firm's assets may not be an attractive target, the firm can still be a useful springboard for an attack against a third party. Such attacks may occur because of a relationship (e.g., vendor or supplier) or may merely be opportunistic attempts at recruitment (e.g., recruiting, or perhaps more appropriately, press-ganging desktop computers for botnets).

External events and their repercussions can also wreak havoc, both directly and indirectly. The September 11, 2001, terrorist attacks that destroyed New York City's World Trade Center complex had worldwide impact, not only on systems and firms located in the destroyed complex. Telecommunications infrastructure was damaged or destroyed, severing Internet links for many organizations. Parts supply and all travel was disrupted when North American airspace was closed for several days. Manhattan was sealed to exits and entries, while within the city itself, and throughout much of the world, normal operations were suspended. The September 11 attacks were extraordinarily disruptive, but security precautions similar to those described throughout this *Handbook* served to ameliorate damage to Web operations and other infrastructure elements of those concerns that had implemented them. Indeed, the existence of the Web and the resulting ability to organize groups without physical presence proved a means to ameliorate the damage from the attacks, even to firms that had a major presence in the World Trade Center. In the period following the attacks on the World Trade Center, Morgan Stanley and other firms with offices in the affected area implemented extensive telecommuting and Web-based interactions first to account for their staffs¹¹ and then to enable work to continue.¹²

Best practices and scale are important. Some practices, issues, and concerns at first glance appear relevant only to very large organizations, such as Fortune 500 companies. In fact, this is not so. Considering issues in the context of a large organization permits them to appear magnified and in full detail. Smaller organizations are subject to the same issues and concerns but may be able to implement less formal solutions. "Formal" does not necessarily imply written procedures. It may mean that certain computer-related practices, such as modifying production facilities in place, are inherently poor ideas and should be avoided. Very large enterprises might address the problems by having a separate group, with separate equipment, responsible for operating the development environment. Incremental staged deployments are an excellent technique to calibrate risk exposure.

30.3.2 Defining Attacks. Repeated, multiple attempts to connect to a server could be ominous, or they could be nothing more than a customer with a technical problem. Depending on the source, large numbers of failed connects or aborted operations coming from gateway nodes belonging to an organization could represent a problem somewhere in the network, an attack against the server, or a blend of both. It could also represent something no more ominous than a group of users within a locality accessing a Web resource through a firewall.

30.3.3 Defining Protection. There is a difference between protecting Internet-visible assets and protecting Websites. For the most part, Internet-visible assets are not intended for public use. Thus, it is often far easier to anticipate usage volumes and to account for traffic patterns. Websites are subject to the vicissitudes of worldwide public usage. A dramatic surge in traffic could be an attack, or it could be an unexpected

30 · 20 E-COMMERCE AND WEB SERVER SAFEGUARDS

display of the site's URL in a television program or in a relatively unrelated news story. Differentiating between belligerence and popularity is difficult, if not impossible.

Self-protective measures that do not impact customers are always permissible. However, care must be exercised to ensure that the measures are truly impact-free. As an example, some sites, particularly public FTP servers, often require that the Internet protocol (IP) address of the requesting computer have an entry in the inverse domain name system, which maps IP addresses to host names (e.g., node 192.168.0.1 has a PTR [pointer record] 1.0.168.192.in-addr.arpa) (RFC1034, RFC1035; Mockapetris 1987a, 1987b) as opposed to the more widely known domain name system database, which maps host names into IP addresses. It is true that many machines have such entries, but it is also true that many sites, including company networks and many ISPs, do not provide inverse DNS information. Whether this entire population should be excluded from the site is a policy and management decision, not a purely technical decision. Even a minuscule incident rate on a popular Website can be catastrophic, both for the provider and for the naïve end user who has no power to understand or resolve the situation.

30.3.4 Maintaining Privacy. Logging interactions between customers and the Website is also a serious issue. A Website's privacy policy is again a managerial, legal, and customer relations issue with serious overtones. Technical staff needs to be conscious that policies, laws, and other issues may dictate what information may be logged, where it can be stored, and how it may be used. For example, the 1998 Children's Online Privacy Protection Act (COPPA) (15 U.S.C. § 6501 et seq.) makes it illegal to obtain name and address information from children under the age of 13 in the United States. Many firms are party to agreements with third-party organizations such as TRUSTe,¹³ governing the use and disclosure of personal information. For more information on legal aspects of protecting privacy, see Chapter 69 in this *Handbook*.

30.3.5 Working with Law Enforcement. Dealing with legal authorities is similarly complicated. Attempts at fraudulent purchases and other similar issues can be addressed using virtually the same procedures that are used with conventional attempts at mail or phone order fraud. Dealing with attacks and similar misuses is more complicated and depends on the organization's policies and procedures, and the legal environment. The status of the Website is also a significant issue. If the server is located at a hosting facility, or is owned and operated by a third party, the situation becomes even more legally complicated. Involving law enforcement in a situation will likely require that investigators have access to the Web servers and supporting network, which may be difficult. Last, there is a question of what information is logged, and under what circumstances. For more information on working with law enforcement, see Chapter 61 in this *Handbook*.

30.3.6 Accepting Losses. No security scheme is foolproof. Incidents will happen. Some reassurance can be taken from the fact that the most common reasons for system compromises in 2001 appear to remain the same as when Clifford Stoll wrote *The Cuckoo's Egg* in 1989. Then and now, poorly secured systems have:

- Management accounts with obvious passwords
- Unprotected system files
- Unpatched known security holes

RULES OF ENGAGEMENT 30 · 21

However, eliminating the simplest and most common ways in which outsiders can compromise Websites does not resolve all problems. The increasing complexity of site content, and of the applications code supporting dynamic sites, means that there is an ongoing design, implementation, testing, and quality assurance challenge. Server-based and server-distributed software (e.g., dynamic www sites) is subject to the same development hazards as other forms of software. Security hazards will slip into a Website despite the best efforts of developers and testers. The acceptance of this reality is an important part of the planning necessary to deal with the inevitable incidents. When it is suspected that a Website, or an individual component, has been compromised, the reaction plans should be activated. The plans required are much the same as those discussed in Chapter 56 in this *Handbook*. The difference is that the reaction plan for a Website has to take into consideration that the group primarily impacted by the plan will be the firm's customers. The primary goal of the reaction plan is to contain the damage. For more information on computer security incident response, see Chapter 56.

30.3.7 Avoiding Overreaction. Severe reactions may create as much, if not more, damage than the actual attack. The reaction plan must identify the decision-making authority and the guidelines to allow effective decisions to be made. This is particularly true of major sites, where attacks are likely to occur on a regular basis. Methods to determine the point at which the Website must be taken offline to prevent further damage need to be determined in advance.

Offline clones should be prepared to enable fast responses. When under time pressure during an emergency, the last thing one wants to do is building production system images.

In summary, when protecting Websites and customers, defensive actions are almost always permissible and offensive actions of any kind are almost always impermissible. Defensive actions that are transparent to the customer are best of all.

30.3.8 Appropriate Responses to Attacks. Long before the advent of the computer, before the development of instant communications, international law recognized that firing upon a naval vessel was an act of war. Captains of naval vessels were given standing orders summarized as *fire if fired upon*. In areas without readily accessible police protection, the right of citizens to defend themselves is generally recognized by most legal authorities. Within the body of international law, such formal standards of conduct for military forces are known as *rules of engagement*, a concept with global utility.

In cyberspace, it is tempting to jettison the standards of the real world. It is easy to imagine oneself master of one's own piece of cyberspace, without connection to real-world laws and limitations on behavior. However, information technology (IT) personnel do not have the legal standing of ships' captains with no communications to the outside world. Some argue that "fire if fired upon" is an acceptable standard for online behavior. Such an approach does not take into account the legal and ethical issues surrounding response strategies and tactics. It may be unsatisfying, but the only response relatively free of legal consequences (criminal, civil, and regulatory) is "Defensive Action Only." Blocking address ranges is defensive; firing back crosses the line into offense.

Any particular security incident has a range of potential responses. Which response is appropriate depends on the enterprise and its political, legal, and business environment. Acceptability of response is also a management issue as well as potentially a political

30 · 22 E-COMMERCE AND WEB SERVER SAFEGUARDS

issue. Determining what responses are acceptable in different situations requires input from management on policy, from legal counsel on legality, from public relations on public perceptions, and from technical staff on technical feasibility. Depending on the organization, it also may be necessary to involve unions and other parties in the negotiation of what constitutes appropriate responses.

What is acceptable or appropriate in one area is not necessarily acceptable or appropriate in another. Often the national security arena has lower standards of proof than would be acceptable in normal business litigation. In U.S. civil courts, cases are decided upon a preponderance of evidence. Standards of proof acceptable in civil litigation are not conclusive when a criminal case is being tried, where guilt must be established beyond a reasonable doubt.

Gambits or responses that are perfectly legal in a national security environment may be completely illegal and recklessly irresponsible in the private sector, exposing the organization to significant legal liability.

Rules of etiquette and behavior are similarly complex. The rights of prison inmates in the United States remain significant, even though they are subject to rules and regulations substantially more restrictive than for the general population. Security measures, as well, must be appropriate for the persons and situations to which they are applied.

30.3.9 Counter-Battery. Some suggest that the correct response to a perceived attack is to implement the cyberspace equivalent of *counter-battery*, that is, targeting the artillery that has just fired upon you. However, counter-battery tactics, when used as a defensive measure against Internet attacks, will be perceived, technically and legally, as an attack like any other.

Counter-battery tactics may be emotionally satisfying but are prone to both error and collateral damage. Counter-battery can be effective only when the malefactor is correctly identified and the effects of the reciprocal attack are limited to the malefactor. If third parties are harmed in any way, then the retaliatory action becomes an attack in and of itself. One of the more celebrated counter-battery attacks gone awry was the 1994 case when two lawyers from Phoenix, Arizona, spammed over 5,000 Usenet newsgroups to give unsolicited information on a U.S. Immigration and Naturalization Service lottery for 55,000 green cards (immigration permits). The resulting retaliation—waves of email protests—against the malefactors flooded their Internet service provider (ISP) and caused at least one server to crash, resulting in a denial of service to all the other, innocent, customers of the ISP.¹⁴

30.3.10 Hold Harmless. Operational policies for Internet operations must adopt a hold harmless position vis-à-vis good faith actions within policy. Dealing with an Internet crisis often requires fast reactions. If employees act in good faith, in accordance with their responsibilities, and within documented procedures, those making the decisions should be immune from reprisal. Punishment for actions undertaken in good faith is counterproductive. If the procedures are wrong, improve the rules and procedures; do not blame the employees for following established policies. Disciplinary actions are manifestly inappropriate in such circumstances.

30.4 RISK ANALYSIS. As noted earlier in this chapter, protecting an organization's Websites depends on an accurate, rational assessment of the risks. Developing effective strategies and tactics to ensure site availability and integrity requires that all potential risks be examined in turn.

RISK ANALYSIS 30 · 23

Unrestricted commercial activity has been permitted on the Internet since 1991. Since then, enterprises large and small have increasingly integrated Web access into their second-to-second operations. The risks inherent in a particular configuration and strategy are dependent on many factors, including the scale of the enterprise and the relative importance of the Web-based entity within the enterprise. Virtually all high-visibility Websites (e.g., Yahoo!, America Online, cnn.com, Amazon.com, and eBay) have experienced significant outages at various times.

The more significant the organization's Web component, the more critical is availability and integrity. Large, traditional firms with relatively small Web components can tolerate major interruptions with little damage. Firms large or small that rely on the Web for much of their business must pay greater attention to their Web presences, because a serious outage can quickly escalate into financial or public relations catastrophe.

For more details of risk analysis and management, see Chapters 62 and 67 in this *Handbook*.

30.4.1 Business Loss. Business losses fall into several categories, any of which can occur in conjunction with an organization's Web presence. In the context of this chapter, customers are both outsiders accessing the Internet presence and insiders accessing intranet applications. In practice, insiders using intranet-hosted applications pose the same challenges as the outside users.

30.4.2 PR Image. The Website is the organization's public face 24/7/365. This ongoing presence is a benefit, making the firm visible at all times, but the site's high public profile also makes it a prime target.

Government sites in the United States and abroad have often been the targets of attacks. In January 2000, "Thomas," the Website of the U.S. Congress, was defaced. Earlier, in 1996, the Website of the U.S. Department of Justice was vandalized. Sites belonging to the Japanese, U.K., and Mexican governments also have been vandalized.

Attacks against Websites continue to be an ever-increasing problem. On March 21, 2013, a series of malware infections crippled broadcasters and financial institutions in South Korea during a time of increasing tensions over North Korean nuclear and missile testing, including a renunciation by North Korea of the 1953 armistice suspending the Korean War.¹⁵

These incidents continue, with hackers on both sides of various issues attacking the other side's Web presence. While no conclusive evidence of the involvement of nation states in such activities has become generally known, it is inevitable. Surges of such activity have coincided with major public events, including the 2001 U.S.–China incident involving an aerial collision between military aircraft, the Afghan and Iraqi operations, and the Israeli–Lebanese war of 2006, the Russian–Georgian conflict,¹⁶ and Estonia.¹⁷ Such activity has not been limited to the national security arena, however. Many sites were defaced during the incidents following publication of cartoons in a Danish newspaper that some viewed as defaming the prophet Muhammad.

In some cases, companies have suffered collateral damage from hacking contests, where hackers prove their mettle by defacing as many sites as they can.¹⁸ The fact that there is no direct motive or animus toward the company is not relevant; the damage has still been done.

In the corporate world, company Websites have been the target of attacks intended to defame the corporation for real or imagined slights. Some such episodes have been reported in the news media, whereas others have not been the subjects of extensive reporting. The scale and newsworthiness of the episode is unimportant; the damage

30 · 24 E-COMMERCE AND WEB SERVER SAFEGUARDS

done to the targeted organization is the true measure. An unreported incident that is the initiating event in a business failure is more damaging to the affected parties than a seemingly more significant outage with less severe consequences.

Other cybervandals (e.g., sadmind/IIS) have used address scanners to target randomly selected machines. Obscurity is not a defense against address-scanner attacks.

30.4.3 Loss of Customers/Business. Internet customers are highly mobile. Website problems quickly translate into permanently lost customers. The reason for the outage is immaterial; the fact that there was a problem is often sufficient to provoke customer flight.

In most areas, there is competitive overlap. Using the overnight shipping business as an example, in most U.S. metropolitan areas there are (in alphabetical order) Federal Express, United Parcel Service, and the United States Postal Service. All of the firms offer Web-based shipment tracking, a highly popular service. Problems or difficulties with shipment tracking will quickly lead to a loss of business in favor of a different company with easier tracking.

30.4.4 Interruptions. Increasingly, modern enterprises are being constructed around ubiquitous 24/7/365 information systems, most often with Websites playing a major role. In this environment, interruptions of any kind are catastrophic.

Production. The past 20 years have seen a streamlining of production processes in all areas of endeavor. Twenty years ago, it was common for facilities to have multiday supplies of components on-hand in inventory. Today, *zero latency* or *just-in-time* (JIT) environments are common, permitting large facilities to operate with minimal inventory. Fiscally, zero latency environments may be optimally cost efficient, yet the paradigm leaves little margin for disruptions of the supporting logistical chain. This chain is sometimes fragile, and subject to disruption by any number of hazards.

Supply Chain. Increasingly, it is common for Web-based sites to be an integral part of the supply chain. Firms may encourage their vendors to use a Web-based portal to gain access to the vendor side of the purchasing system. XML¹⁹-based gateways and *service-oriented architecture* (SOA) approaches, together with other Web technologies, are used to arrange for and manage the flow of raw materials and components required to support production processes. The same streamlining that speeds information between supplier and manufacturer also provides a potential for serious mischief and liability.

Delivery Chain. Web-based sites, both internal and external, human readable and intersystem, have become the vehicle of choice for tracking the status of orders and shipments. Increasingly system-system links, implemented as electronic-documents (e.g., XML or JSON) have become the backbone of many enterprises' delivery chain management and inquiry systems.

Information Delivery. Banks, brokerages, utilities, and municipalities are increasingly turning to the Web as a convenient, low-cost method for managing their relationships with consumers. Firms are also supporting downloading records of transactions and other relationship information in formats required by personal database programs and organizers. These outputs, in turn, are often used as inputs to other processes, which then generate other transactions. Not surprisingly, as time passes, more and more people and businesses depend on the availability of information on demand. Today's Web-based customers presume that information is accessible wherever they can use a smartphone or tablet. This is reminiscent of usage patterns of automatic teller machines in an earlier decade, which allowed people to base their plans on access

OPERATIONAL REQUIREMENTS 30 · 25

to teller machines, often making multiple \$20 withdrawals instead of cashing a \$200 check weekly.

30.4.5 Proactive versus Reactive Threats. Some threats and hazards can be addressed proactively, whereas others are inherently reactive. When strategies and tactics are developed to protect a Web presence, the strategies and tactics themselves can induce availability problems.

As an example, consider the common strategy of having multiple name servers responsible for providing the translation of domain names to IP addresses. It is required before a domain name (properly referred to as a Domain Name System [DNS] zone) can be entered into the root-level name servers, that at least two name servers be identified to process name resolution requests. Name servers are a prime example of resources that should be geographically diverse.

Updating DNS zones requires care. If an update is performed improperly, then the resources referenced via the symbolic DNS names will become unresolvable, regardless of the actual state of the Web server(s) and related infrastructure. The risk calculus involving DNS names is further complicated by the common, efficient, and appropriate practice of designating ISP name servers as the primary mechanism for the resolution of domain names. In short, name translation provides a good example of the possible risks that can affect a Web presence.

30.4.6 Threat and Hazard Assessment. Some threats are universal. Others are specific to an individual environment. The most devastating and severe threats are those that simultaneously affect large areas or populations, where efforts to repair damage and correct the problem are hampered by the scale of the problem. Differences in perspective can impede effective response. What seems a minor outage to a carrier may be a business-threatening outage to a customer.

On a basic level, threats can be divided into several categories. The first is between deliberate acts and accidents. Deliberate acts comprise actions done with the intent to damage the Website or its infrastructure. Accidents include natural phenomena (acts of God) and clumsiness, carelessness, and unconsidered consequences (acts of clod).

Broadly put, a deliberate act is one whose goal is to impair the system. Deliberate acts come in a broad spectrum of skill and intent. For the purpose of risk analysis and planning, deliberate acts against infrastructure providers can often appear to be acts of God. To an organization running a Website, an employee attack against a telephone carrier appears simply as one more service interruption of unknown origin.

No enterprise or agency should consider itself an unlikely target. Past high-profile incidents have targeted the FBI (May 26 and 27, 1999), major political parties, and interest groups in the United States. On the consumer level, numerous digital subscriber line (DSL)-connected home systems have been targeted for subversion as preparation for the launching of DDoS attacks. In 2007, several investigations resulted in the arrest of several individuals for running so-called botnets (ensembles of compromised computers). These networks numbered hundreds of thousands of machines.²⁰ Recently, Aramco was victimized by a major attack²¹; as were major Korean financial institutions.²² The potential for such networks to be used for mischief cannot be underestimated.

30.5 OPERATIONAL REQUIREMENTS. Internet-visible systems are those with any connection to the worldwide Internet. It is tempting to consider protecting

30 · 26 E-COMMERCE AND WEB SERVER SAFEGUARDS

Internet-visible systems as a purely technical issue. However, technical and business issues are inseparable in today's risk management. For example, as noted earlier in this chapter, the degree to which systems should be exposed to the Internet is fundamentally a business risk-management issue. Protection technologies and the policies behind the protection can be discussed only after the business risk questions have been considered and decided, setting the context for the technical discussions. In turn, business risk-management evaluation (see Chapter 62) must include a full awareness of all of the technical risks. Ironically, nontechnical business managers can accurately assess the degree of business risk only after the technical risks have been fully exposed.

Additional business and technical risks result from outsourcing. Today, many enterprises include equipment owned, maintained, and managed by third parties. Some of this equipment resides on the organization's own premises and other equipment resides offsite: for example, at application service provider facilities.

Protecting a Website begins with the initial selection and configuration of the equipment and its supporting elements, and continues throughout its life. In general, care and proactive consideration of the availability and security aspects of the site from the beginning will reduce costs and operational problems. Although virtually impossible to achieve, the goal is to design and implement an automatic system, with a configuration whose architecture and implementation operates even in the face of problems, with minimal customer impact.

That is not to say that a Website can operate without supervision. Ongoing, proactive monitoring is critical to ensuring the secure operation of the site. Redundancy only reduces the need for real-time response by bypassing a problem temporarily; it does not eliminate the underlying cause. The initial failure must be detected, isolated, and corrected as soon as possible, albeit on a more schedulable basis. Otherwise, the system will operate in its successively degraded redundancy modes until the last redundant component fails, at which time the system will fail completely.

30.5.1 Ubiquitous Internet Protocol Networking. Business has been dealing with the security of internets (i.e., interconnected networks) since the advent of internetworking in the late 1960s. However, the ubiquity of Transmission Control Protocol/Internet Protocol (TCP/IP) networks and of the public Internet has exposed much more equipment to attack than in the days of closed corporate networks. In addition, a much wider range of equipment, such as voice telephones based on voice-over IP (VoIP), fax machines, copiers, and even soft drink dispensers, are now network accessible.

IP connectivity has been a great boon to productivity and ease of use, but it has not been without a darker side. Network accessibility also has created unprecedented opportunities for improper, unauthorized access to networked resources and other mischief. It is not uncommon to experience probes and break-in attempts within hours or even minutes of unannounced connection to the global Internet.

Protecting Internet-visible assets is inherently a conflict between ease of access and security. The safest systems are those unconnected to the outside world. Similarly, the easiest to use systems are those that have no perceivable restrictions on use. Adjusting the limits on user activities and access must balance conflicting requirements. As an example, many networks are managed *in-band*, meaning that the switches, routers, firewalls, and other elements of the network infrastructure are managed using the actual network as the connection medium. If the management interfaces were not managed over properly encrypted connections, management passwords would be visible on the

OPERATIONAL REQUIREMENTS 30 · 27

network. An outsider or unauthorized insider may thus monitor the network, gaining sufficient information to paralyze the network. This can result in severe damage and liability to the organization.

30.5.2 Internal Partitions. Complex corporate environments can often be secured effectively by dividing the organization into a variety of interrelated and nested security domains, each with its own legal, technical, and cultural requirements. For example, there are specific legal requirements for medical records (see Chapter 71 in this *Handbook*) and for privacy protection (see Chapter 69). Partners and suppliers, as well as consultants, contractors, and customers, often need two-way access to their corporate data and facilities. These diverse requirements mean that a single corporate firewall is often insufficient. Different domains within the organization will often require their own firewalls and security policies. Keeping track of the multitude of data types, protection and access requirements, and different legal jurisdictions and regulations makes for previously unheard-of degrees of complexity. It is simply not possible to configure a complete set of rules to implement the multiple, different security requirements without creating security gaps or contradictions.

Much has been written about so-called *Advanced Persistent Threats*.²³ While the terminology is relatively new, the concept is not, having been mentioned in earlier editions of this handbook.²⁴ The fundamental underlying change with advanced threats is the degree of targeting and the attacker's intention to remain covert. In these respects a useful historic precedent is involvement in intelligence gathering by submarine forces. During World War II, submarines were first explicitly tasked with intelligence gathering. The orders given can be summarized as "get the information and/or photographs; avoid detection at all costs." The translation is straightforward: No matter how attractive, do not announce your presence by making an attack; your purpose is to stay covert and extract the desired information. This is precisely the model of today's targeted, information extraction attacks.

Damage control is another property of a network with internal partitions. A system compromised by an undetected malware component will be limited in its ability to spread the contagion beyond its own compartment.

30.5.3 Critical Availability. Networks are often critical for second-to-second operations. The side effects of ill-considered countermeasures may be worse than the damage from an actual attack. For example, shutting down the network or even part of it, for maintenance or repair can wreak more havoc than penetration by a malicious hacker. Damage control, often starting with compartmentalization, is a critical part of limiting the impact of an attack.

30.5.4 Accessibility. Users must be involved in the evolution of rules and procedures. Today, it is still not unheard of for a university faculty to take the position that any degree of security will undermine the very nature of their community, compromising their ability to perform research and inquiries. This extreme position persists despite the attention of the mass media, the justified beliefs of the technical community, and documented evidence that lack of protection of any Internet-connected system undermines the safety of the entire connected community.

Connecting a previously isolated computer system or network to the global Internet creates a communications pathway to every corner of the world. Customers, partners, and employees can obtain information, send messages, place orders, and otherwise interact 24 hours a day, seven days a week, 365 days a year, from literally anywhere on

30 · 28 E-COMMERCE AND WEB SERVER SAFEGUARDS

or near Earth. The Space Shuttle had Internet access; as does the International Space Station. Under these circumstances, the possibilities for attack or inadvertent misuse are limitless.

Despite the vast increase in connectivity, some businesses and individuals do not need extensive access to the global Internet for their day-to-day activities, although they may resent being excluded. The case for universal access is therefore a question of business policy and political considerations.

30.5.5 Applications Design. Protecting a Website begins with the most basic steps. First, a site processing confidential information should always support the Secure Hypertext Transfer Protocol (HTTPS), typically using Transmission Control Protocol (TCP) port 443. Properly supporting HTTPS requires the presence of an appropriate digital certificate (see Chapter 37).

When the security requirements are uncertain, the site design should err on the side of using HTTPS for communications. Eavesdropping on the Internet is an ever-present hazard. There have been cases of routing database manipulation and state-sponsored surveillance. The Hannaford and TJ Maxx episodes described earlier in this chapter both involved the application of what can only be described as classic signals intelligence (SIGINT) techniques. Even though details on such attacks are often not made public, the voluminous historical literature on communications intelligence (COMINT) and SIGINT from the World War II period makes one point abundantly clear: Encryption of all potentially sensitive traffic is the only way to protect information.

Encryption also should be used within an organization, possibly with a different digital certificate, for sensitive internal communications and transactions. Earlier, it was noted that organizations are not monolithic security domains. This is nowhere more accurate than when dealing with human resources, employee evaluations, compensation, benefits, and other sensitive employee information. There are positive requirements that this information be safeguarded, but few organizations have truly secured their internal networks against internal monitoring. It is far safer to route all such communications through securely encrypted channels. Such measures also demonstrate a good faith effort to ensure the privacy and confidentiality of sensitive information. Such efforts are an important element of an organization's defense when an incident occurs. Upfront care can significantly reduce legal liability.

Failure to heed these hazards has been at the root of several incidents, including the TJ Maxx and Hannaford affairs.

It is also important to avoid providing all of the authentication information on a single page or, for that matter, in a sequence of pages. When parts of information are suppressed, as, for example, portions of a credit card or account number, the division between suppressed and displayed portions should be maintained. Displaying all of the information, even if it is on different screens, is an invitation to a security breach.

30.5.6 Provisioning. Although today's hardware has unprecedented reliability, any failure of hardware between the customer and the data center will impair an enterprise's Web presence. For a highly available, front-line Website, the effective requirement is a minimum of two diversely located facilities, each with a minimum of two servers. This is not necessarily an expensive proposition. Fairly powerful Web servers can be purchased for less than \$5,000, so the total hardware expenditure for four servers is reasonable, substantially less than the annual cost of a single technician. Cloud-computing, the availability of data center-housed virtual servers as a service (e.g.,

OPERATIONAL REQUIREMENTS 30 · 29

Amazon, Verizon, RackSpace, Google) presents another option with low entry costs. In most cases, the cost of the extra hardware is more than offset by the business cost of downtime, which can sometimes exceed the total cost of the duplicative hardware by a factor as much as 100, in a single episode.²⁵

Duplicate hardware and geographic diversity ensure constant customer access to some degree of functionality. The degree of functionality that must be maintained depends on the market and the customers. Financial firms supporting online stock trading have different operational, regulatory, and legal requirements than supermarkets. The key is matching the support level to the activities. Some degree of planned degradation is generally acceptable. Total unavailability is not an option.

Remote hosting services generally referred to as “cloud computing” do not eliminate risks. While renting time and other resources from a provider does reduce capital and operating costs, there is a corresponding loss of visibility and control. A failure at a remote vendor operated facility can take down a critical service. Depending on your applications, industry, regulatory, and legal environments, the risks of cloud computing may be quantifiable, or they may be immeasurable. Difficulties in vendor relationships can give rise to “hostage data” situations. From a security and information assurance perspective, planning in advance to deal with changes in cloud vendors should be part of planning.

There are also unresolved legal issues surrounding the use of shared computing facilities operated by third parties. If a search warrant is served on an internal facility, the organization will clearly be aware of what has happened. If a shared facilities vendor is served with a warrant, it is an open question whether the vendor has the desire (or capability) to present the defenses clearly available with data stored within the company.

30.5.7 Restrictions. All Web servers should be located behind a firewall in a DMZ, as discussed in Section 30.6.4. Incoming and outgoing services should be restricted using protocols such as HTTP, HTTPS, and Internet Control Message Protocol (ICMP). For troubleshooting purposes, it is desirable to implement ICMP, which is used by PING, an echo requester, as a way to check connectivity. All unused ports should be disabled. Furthermore, the disabled ports should be blocked by the firewalls separating the DMZ from the outside world.

Customer information should, to the extent possible, be stored on systems separate from the systems actually providing Web serving. Many security episodes have exploited file protection errors on Web servers, in order to access databases directly. Segregating customer data on separate machines, and ensuring that the only way to access customer data is through the documented pathways, is likely severely to impede improper attempts to access and modify information.

These safeguards are especially important for high-security information such as credit card numbers. The number of incidents in which malefactors have downloaded credit card numbers directly from Website is an indication of the importance of such precautions. The systems actually storing the sensitive information should never be accessible from the public Internet. Storage and processing of payment card information is also subject to industry guidelines, particularly the PCI DSS.

The TJX case, which came to public attention in the beginning of 2007, was one of a series of large-scale compromises of electronically stored information on back-office and e-commerce systems. Most notably, the TJX case appears to have started with an insufficiently secured corporate network and the associated back-office systems, not

30 · 30 E-COMMERCE AND WEB SERVER SAFEGUARDS

a Website penetration. This breach escalated into a security breach of corporate data systems. It is reported on the TJX Website that at least 45.7 million credit cards were compromised (original reports in *USA Today* and other publications cite 94 million credit card numbers as being compromised). On November 30, 2007, it was reported that TJX, the parent organization of stores including TJ Maxx and Marshall's, had agreed to settle bank claims related to VISA cards for US\$ 40.9M.²⁶ Also, as of March 2008, a class action suit on behalf of customers was in process of being settled. It called for fees to be paid to each of the plaintiffs for credit monitoring and identity theft, and for new driver's licenses, as well as \$6.5 million to plaintiff's counsel. In spite of several such high-profile attacks, the lessons appear not to have been learned. As recently as March 2008, Hannaford Bros. Co. reported that 4.2 million credit card numbers had been exposed, with about 1,800 cases of fraudulent usage reported as of that date. An incident involving the network supporting Sony's PlayStation 3²⁷ disclosed personal information relating to an estimated 70 million subscribers. Clearly, more stringent controls, care, and safeguards are called for.

30.5.8 Multiple Security Domains. The front-line Web servers, and the database servers supporting their activities, comprise two different security domains.

The Web servers, as noted previously, need to be globally accessible via HTTP, HTTPS, and ICMP (while not needed for actual operations, ICMP is often needed for troubleshooting connectivity problems by both customers and providers). In turn, they need to access the application or database servers, and *only* those servers. In a production environment, it is preferable that application or database servers interact with the Web servers using a dedicated, restricted-use protocol. Properly implemented, such a restriction prevents a hijacked Web server from exploiting its access to the application or database server.

These second-tier servers should be in a security domain separated by restrictive firewalls from the externally accessible front-line Web servers. This seems like a significant expenditure, but it is often less expensive and lower risk than the costs associated with correcting a single significant incident.

30.5.9 What Needs to Be Exposed? Publicly accessible Websites need publicly accessible Web servers to perform their functions. The challenge is to provide the desired services to the public without simultaneously providing access that can be leveraged in unauthorized ways to subvert the site. A penetration incident may lead to significant financial losses, embarrassment, and financial and (in some cases) criminal liability.

Generally speaking, Web servers should not be connected to the public Internet (unless it is possible to verify that all other ports are closed to traffic). All connections to the public network should be made through a firewall system, with the firewall configured to pass only Web-related traffic to those hosts. A DMZ guarded by a firewall often creates a more easily controlled security gateway.

Many sites will opt to place externally visible Web servers in a security compartment of their own, on a separate port of the firewall (if not a totally separate firewall), using a separate DMZ from other publicly accessible resources. These precautions may seem excessive, but having improperly secured systems can lead to security breaches that are extremely difficult to correct and can lead to extended downtime while the problems are analyzed and remedied. In this case, an ounce of prevention is worth substantially more than a pound of cure.

OPERATIONAL REQUIREMENTS 30 · 31

30.5.9.1 Exposed Systems. Exposed systems are inherently a security hazard. Systems without access from or to the public network cannot be compromised from the public network. Only systems that absolutely need to be publicly accessible should be so configured. Minimizing the number of exposed systems is generally desirable, but this is best considered in terms of machine roles rather than actual counts of systems. Increasing the load on each publicly accessible server by increasing the size of the server, thus increasing the amount of capacity impacted by a single system outage, is not a benefit. However, this must be balanced against the new trend toward server virtualization, which does not increase the size of a server, but which increases its utilization, in order to lower costs and improve reliability.

The introduction of SSL-based VPN implementations and services supporting remote access via secure HTTPS connections (e.g., gotomypc.com) create an entirely new class of security hazard. SSL and other encrypted techniques have been used to exfiltrate data from security penetrations. While encrypted connections may be difficult or impossible to observe, connections to unexpected locations are a signal of possible security compromise.

30.5.9.2 Hidden Subnets. The servers directly supporting the Website need to be accessed by the outside world, and thus must generally have normal, externally visible Internet addresses. Alternatively, traffic directors, firewalls, or similar devices may map an external port number to an internal address and port number. However, in most cases the other systems supporting the Web servers generally have no legitimate need for unrestricted access from or to the public network.

The safest address assignments for supporting, non-outside-visible systems are the IPv4 addresses allocated for use in private internets²⁸ and their equivalent IPv6 addresses.²⁹ Needless to say, these systems should be in a security compartment separated from the publicly accessible Web servers, and that compartment should be isolated from the publicly accessible compartment with a very restrictive firewall.

30.5.10 Access Controls. Publicly accessible systems are both the focus of an organization's security efforts and the primary target of attempts to compromise that security. The number of individuals authorized to make changes to the systems and the ways in which changes may be made need to be carefully controlled, reported, and monitored. The cleared individuals should use individual accounts, and the access to sensitive functions through such accounts should be immediately invalidated if the individual's access is no longer authorized or if the employee ceases employment, is under investigation for impropriety, or other reason. The most efficient and effective way to implement these requirements is role-based access control, which leverages the underlying operating system protection model to enforce file security.³⁰ The key aspect of role-based protection is that file access is protected by role, not identifier. It is thus possible to remove a user's access to production files with a single command.

Access controls are also not merely a question of traditional role-based control. Increasingly, applications are directly accessed by customers. Such direct access requires more precautions and care than that provided by customer service personnel. It is often appropriate for customer service personnel to have access to most accounts; it is just as inappropriate to give corresponding unfettered access to end customers, who should only have access to their specific account. Applications-level attacks (e.g., SQL injection) become increasingly hazardous.

30 · 32 E-COMMERCE AND WEB SERVER SAFEGUARDS

The transition to Web browser and mobile device-based interfaces needs additional security precautions. Remote clients are subject to compromise and reverse-engineering. As an example, using plaintext, client-supplied information to control access (e.g., cookies or URL-encoded) can easily give rise to an SQL injection (or equivalent) attack, merely by editing the unencoded information. Switching to proper variable, session-based numbering scheme would make tampering far less successful (indeed, the greater number of failures resulting from using such an opaque number would be easily detectable by applications and security logging systems).

The challenge to such systems is that the safeguards must be implemented within the application. Such safeguards are needed to ensure applications security.

SQL injection attacks of this type have occurred to major organizations. In 2011, details of 200,000 credit card accounts with Citibank were compromised in this way.³¹

For more information on access controls, see Chapters 28 and 29 in this *Handbook*.

30.5.11 Site Maintenance. Maintaining and updating the Website requires great care. The immediate nature of the Web makes it possible for a single-character error in a major enterprise-level application to cause hundreds of thousands of dollars of damage within moments. Web servers need to be treated with respect; the entire enterprise is riding on the electronic image projected by the server. Cybervandalism, which most commonly consists of defacing the home page of a well-known site, requires unauthorized updating of the files comprising the site. In 2000 alone, well-known public and private entities including the FBI, OPEC, World Trade Organization, and NASA, as well as educational institutions including the University of Limerick, have been harmed in this way. These attacks continue to be a danger, both in terms of damage to the organization's image and covertly, as using the WWW site as a launching pad for Web-based exploits.

The Web is inherently a highly leveraged environment. Small changes in the content of a single page percolate throughout the Web in a matter of minutes. Information disseminates easily and quickly at low cost. Leverage helps tremendously when things go right; when things go badly, leverage dramatically compounds the damage. For more information on change control, see Chapters 38, 39, and 40 in this *Handbook*.

30.5.12 Maintaining Site Integrity. Every Website and its servers is a target. Antagonists can be students, activists, terrorists, disgruntled former employees, or unhappy customers as well as state actors. Because Websites are an enterprise's most public face, they represent extremely desirable targets.

Maintaining integrity requires that updates and changes to the site be done in a disciplined manner. Write access to the site must be restricted, and those authorized must use secure methods to access the Web servers. The majority of reported incidents appear to be the result of weak security in the update process. For example, unsecured FTP access from the general Internet is a poor practice. Safer mechanisms include:

- FTP from a specific node within the inner firewall
- KERMIT on a directly wired port or over a secured connection
- Logins and file transfers (e.g., SFTP) over secure authenticated connections via an underlying SSH transport
- Physical media transfers

TECHNICAL ISSUES 30 · 33

Most of the technologies do not inherently require that an onsite individual perform server updates, which would preclude remote maintenance. It does mean that in order to get to a machine from which an update can be performed, it is necessary to come through a virtual private network with point-to-point tunneling protocol or Layer2 tunneling protocol (VPN PPTP/L2TP) authenticated by at least one of the secure gateways. For more information on VPNs, see Chapter 32 in this *Handbook*.

30.6 TECHNICAL ISSUES. There are many technical issues involved in protecting Internet-accessible resources. The technologies used to protect network assets include routers, firewalls, proxy servers, redundancy, and dispersion. When properly designed and implemented, security measures produce a positive feedback loop, where improvements in network security and robustness are self-reinforcing. Each improvement makes other improvements possible and more effective.

30.6.1 Inside/Outside. Some visions of the future include a utopian world where everything is directly accessible from anywhere without effort, and with only beneficial results. The original Internet operated on this basis, until a number of incidents (including the 1988 Morris worm) caused people to rethink the perceived inherent peacefulness of the networked world. It is a trend that has only accelerated, as the total number of systems grows ever larger.

The architecture and design of protective measures for a network depend on differentiating inside trustable systems from outside untrustworthy systems. This is equally true for intranets, the Internet, or an *extranet*, so called to distinguish private interconnections of networks from the public Internet. Unfortunately, trust is not a black-and-white issue. A system may be trustworthy from one perspective and untrustworthy from another, thus complicating the security design. In addition, the vast majority of inappropriate computer use is thought to be done by those with legitimate access to some aspect of the system and its data. Alternatively, it is not possible to be strong everywhere. Trusted but compromised systems are a significant danger.

Basic connectivity configuration is one of those few areas that are purely technical, without a business risk element. One of the most obvious elements involves the tables implemented in routers connecting the enterprise to the public carrier-supplied IP connection. The table rules must prevent *IP spoofing*, which is the misrepresentation of IP packet origins. This is also true when originators are within the organization.

There are three basic rules for preventing IP spoofing applicable to all properly configured networks:

1. Packets entering the network from the outside should never have *originator* addresses within the target network.
2. Packets leaving a network and going to the public network must have *originator* addresses within the originating network.
3. Packets leaving a network and going to the public network must not have *destination* addresses within the originating network.

An exception to these rules is in the use of stealth internal networks, those whose internal addresses correspond to legal external addresses.³²

A corollary to these rules is that packets with originator or destination addresses in the most local intranet addresses range³³ or dynamic IP addresses,³⁴ should never be

30 · 34 E-COMMERCE AND WEB SERVER SAFEGUARDS

permitted to enter or leave an internal network. Nested address spaces may deliberately create aliased address spaces.^{35,36}

30.6.2 Hidden Subnets. Firewalls funnel network traffic through one or more choke points, concentrating the security task in a small number of systems. The reasoning behind such concentration is that the likelihood of security breaches of the entire network rises rapidly with the number of independent access points.

Firewalls and proxy servers (see Chapter 26 in this *Handbook*) are effective only in topologies where the firewall filters all traffic between the protected systems and the less-trusted world outside its perimeter. If the protected systems can in any way be accessed without going through the firewall, then the firewall itself has been rendered irrelevant. Security audits often uncover systems that violate this policy; relying on administrative sanctions to preclude such holes in the security perimeter generally does not work. The rule of international diplomacy, “Trust but verify,” applies.

The simplest solution to this problem is the use of RFC 1918³⁷ addresses within protected networks.³⁸ RFC1918 provides a range of IPv4 addresses guaranteed by the Internet Assigned Numbers Authority (IANA) never to occur in the normal, public Internet. The address ranges used for dynamic IP address assignment have similar properties.³⁹ IPv6 networks are similarly governed by RFC 4193.

Filtering these addresses on both inbound and outbound data streams is straightforward and highly effective at stopping a wide range of attacks.⁴⁰ Requiring the use of such addresses, and prohibiting the internal use of externally valid addresses, goes a long way toward preventing the use of unauthorized protocols and connections.

30.6.3 What Need Be Exposed? A security implementation starts with an analysis of the enterprise’s mission, needs, and requirements. All designs and implementations are a compromise between absolute security, achieved only in a powered-down or disconnected system, and total openness, in which a system is completely open to any and all access from the outside.

Although in most cases communications must be enabled between the enterprise and the outside world for normal business functions, total disconnection, known as an *air gap*, is sometimes both needed and appropriate between specific components. Industrial real-time control systems, life-critical systems, and systems with high requirements for confidentiality remain appropriate candidates for total air gaps.

Often, systems that do not need to receive information from outside sources must publish statistical or other information to less secure systems. This requirement can often be satisfied through the use of limited functionality links such as media exchange, or through tightly controlled one-way transfers. These mechanisms can be implemented with IP-related technologies or with more limited technologies, including KERMIT,⁴¹ UUCP, or vendor-specific solutions.

In other cases, restrictions reflect policies for permitted use and access rather than for protection against outside attack. For example, it is reasonable and appropriate for a public library to limit access to HTTP and related protocols while prohibiting access to such facilities as FTP and TELNET. Remote access using the Secured Shell Protocol (SSH) is extremely useful, as it allows for two-way authentication (both client and server can be authenticated against public/private keys). There is a collateral issue of what provisions should be made for systems within the perimeter to connect to outside networks using various protocols. PPTP and L2TP are the most well-known. Whether these methods of accessing outside networks should be permitted or not is

TECHNICAL ISSUES 30 · 35

a management question of no small import, equivalent to totally suppressing outside communications.

From a security standpoint, tunnels to the outside world are a perfect storm, potentially permitting unlimited, nonmonitored communications to the outside. From the standpoint of enabling business operations, such access may be a practical necessity. Vendors, contractors, suppliers, and customers all need access to internal systems at their respective organizations, and such access will invariably require the use of a VPN tunnel. Perhaps the best solution is to reverse the traditional inside/outside dichotomy, making the LAN an untrusted network. The LAN then becomes a universal dial tone. Access to internal corporate systems can then be secured separately using VPN technology.

The advent of SSL/HTTP-based tunnels presents another challenge. gotomypc.com offers such a service, and individuals have used it to circumvent firewalls and implemented ad hoc remote access. The cost is nominal, often a small fraction of the monthly out-of-pocket outlay for an individual's daily commute to the office. The solution of blocking connections to the IP addresses assigned to such a provider is an often recommended, yet inherently flawed, prophylactic. Blocking one or more such services does nothing to secure the network against an as-yet-unidentified service of this type, such as that hosted on a home server. Limiting the duration of SSL connections is also merely a speed bump to such schemes. These remain a challenge, as does the advent of directly usable SSL/HTTP tunnels communicating using the standard HTTP (TCP Port 80) and HTTPS (TCP Port 443) ports.

In these cases, firewalls are a reasonable solution, so long as their limitations are recognized. For example, firewalls do not prevent mobile code such as Active-X and JAVA from working around prohibited functions. Network attacks using code supplied to browsers for local execution within the local network context have been published. These attacks can persist across time, and are one of the attack classes within the so-called Advanced Persistent Threat spectrum. Such code, in combination with weak security practices on network infrastructure, can cause severe damage (see Chapters 16, 17, 20, and 21 in this *Handbook*).

30.6.4 Multiple Security Domains. Many networks implement security solely at the point of entry, where the organization's network connects to the public Internet. Such *monolithic firewalls* are less than effective choices for all but the simplest of small organizations. As a starting point, systems available to the general population should be outside of the internal security domain in a no-man's land between the public Internet and the private internal net such isolation is referred to as a *demilitarized zone* (DMZ). These systems should be afforded a degree of protection by sandwiching them between an outer firewall (protecting the DMZ from the public network) and an inner firewall (protecting the internal network from the public network and controlling the communications between the publicly accessible servers located in the DMZ to the internal network). Alternatively, the DMZ may be a cul-de-sac attached to a separate port on a single firewall, using a different set of traffic management rules. Such a topology permits implementation of the differing security restrictions applicable to the two networks. Systems located within the DMZ should also be suspect, as they are targets for compromise.

Although the industry as a whole agrees on the need for DMZ configurations, it is less appreciated that such restrictions also have a place within organizations. Different groups, departments, and functions within an organization have different security and access requirements. For example, in a financial services firm, the security

30 · 36 E-COMMERCE AND WEB SERVER SAFEGUARDS

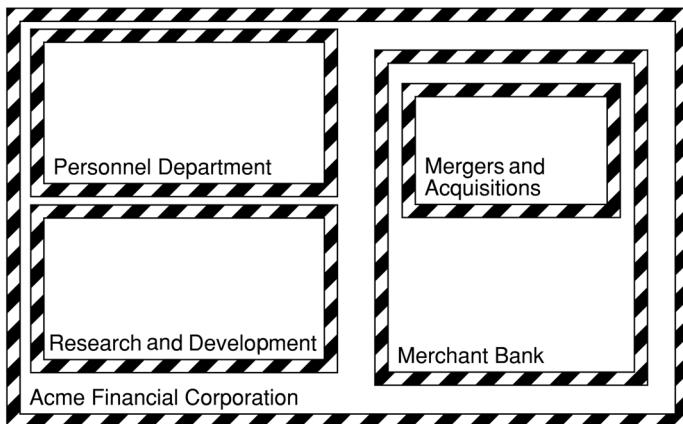


EXHIBIT 30.3 Sibling and Nested Security Domains

requirements differ dramatically among departments. Three obvious examples of departments with different requirements are personnel, mergers and acquisitions, and research and development (see Exhibit 30.3).

The personnel department is the custodian of a wide range of sensitive information about the firm, its employees, and often outsiders who are either regularly on company premises or work regularly with the company on projects. Some of this information, such as residence addresses, pay levels, and license plates, is sensitive for personal or cultural reasons. Other information subjects the organization to legal or regulatory sanctions if it is improperly disclosed or used. In the United States, examples of sensitive data include Social Security numbers, age, sexual orientation, and the presence of human immunodeficiency virus (HIV) or other medical details. Still other information may be subject to specific legal or contractual confidentiality provisions.

The mergers and acquisitions department handles sensitive information of a different sort. Information about business negotiations or future plans is subject to strict confidentiality requirements. Furthermore, the disclosure of such information is subject to a variety of regulations on governmental and securities industry levels. Within the mergers and acquisitions department, access to information often must be on a need-to-know basis, both to protect the deal and to protect the firm from exposure to civil and criminal liability.

Some information in the research and development department is completely open to the public, whereas other information is restricted to differing degrees.

A full implementation of an adequate security environment will require protections that are not only logically different on a departmental basis but also require that different departments be protected from each other. It is difficult, therefore, if not topologically impossible for a single firewall, located at the connection to the outside world, to implement the required security measures. It is difficult to provide assurances that the corporate backbone is free from surreptitious monitoring. In addition, even authorized monitoring can produce communications logs containing sensitive protected information. The network logs then become a hazard in and of themselves.

Securing systems in isolated logical areas is an example of necessary distrust, merely a matter of ensuring that the interactions between the third-party systems and the outside world are allowed to the extent that they are expected. As an example, consider the straightforward situation at Hypothetical Brokerage. Hypothetical Brokerage uses two

TECHNICAL ISSUES 30 · 37

trading networks, Omega and Gamma. At first glance, it would seem that that it would be acceptable to place Omega's and Gamma's network gateways on the usual DMZ, together with Hypothetical's Web servers.

However, this grants a high degree of trust to Omega and Gamma and all of their staff, suppliers, and contractors. The most important operative question is whether there is a credible hazard.

Either of the two gateways is well situated to:

- Monitor the communications traffic to and from Hypothetical's Web servers
- Monitor the traffic between Hypothetical and the other, competing network
- Attack the other gateway
- Disrupt communications to and from the other gateway
- Attack Hypothetical's network

Network providers also represent an attractive attack option. A single break-in to a network provider-supplied system component has the effect of compromising large numbers of end user sites. There is ample history of private (PBX) and public (carrier-owned) switches being preferred targets.⁴² Provider-supplied SOHO routers with remote management or Wi-Fi capability can also become an attack vector if the provider-supplied configurations are deficient, or if the provider-generated management passwords are compromised.⁴³

The solution (see Exhibit 30.4) is to isolate the third-party systems in separate DMZs, with the traffic between each of the DMZs and the rest of the network scrupulously checked as to transmission control protocol and user datagram protocol (TCP/UDP), port number, and source and destination addresses, to ensure that all traffic is authorized. One method is to use a single firewall, with multiple local area network (LAN) ports, each with different filtering rules, to recast Hypothetical's original single DMZ into disjoint, protected, DMZs.

30.6.5 Compartmentalization. Breaking the network into separate security compartments reduces the potential for total network meltdown. Limiting the potential damage of an incident is an important step in resolving the problems.

The same rule applies to the DMZs. For example, in a financial trading or manufacturing enterprise, it is not uncommon to have gateways representing access points to trading and partner networks. Where one places these friendly systems is problematic. Many organizations have chosen to place these systems within their regular DMZ.

Sites have belatedly discovered that such gateways have, on occasion, been found to have acted as routers, taking over that function from the intended routers. In other cases, the gateways have experienced malfunctions and impaired the functioning of the rest of the network (or DMZ). As always, the only solutions certain to work are shutdown or isolation.

Compartmentalization also prevents accidents from cascading. A failure in a single gateway is not likely to propagate throughout the network, because the unexpected traffic will be stopped by the firewall isolating the gateway from the network. Such an event can be made to trigger a firewall's attack alarms.

The network problem is not limited to externally provided nodes. An errant system operating within a noncompartmented network, located as it is within the outer security perimeter, can wreak havoc throughout the entire corporation. Constructing the network as a series of nested and peer-related security domains, each protected by appropriate

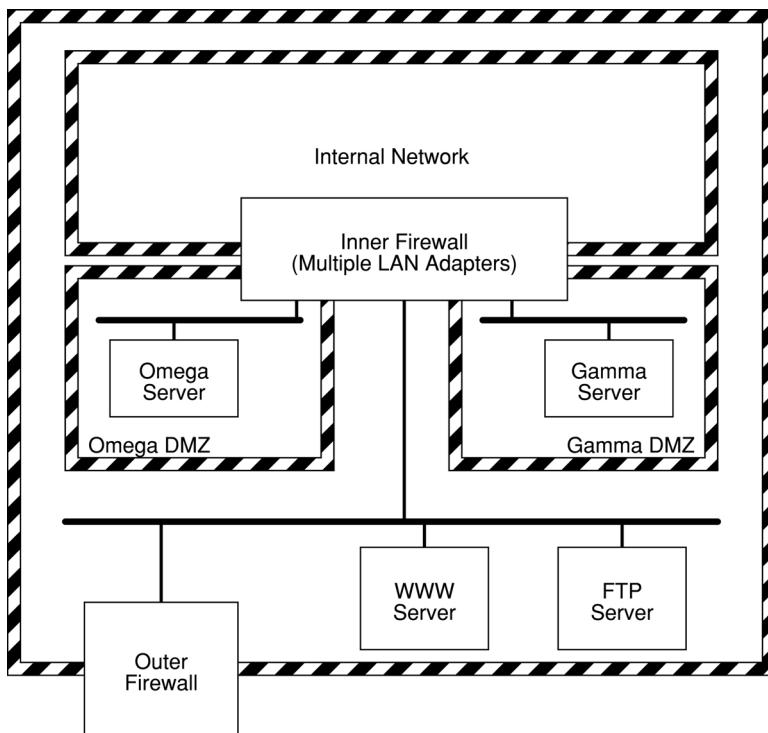
30 · 38 E-COMMERCE AND WEB SERVER SAFEGUARDS

EXHIBIT 30.4 Omega and Gamma Servers in Separate DMZs from Hypothetical's Server

firewalls, localizes the impact of the inevitable incidents, whereas an uncompartmented network permits the contagion to spread unchecked throughout the organization. Larger networks dramatically raise the costs of an incident. It is quite conceivable that the entire budget for compartmenting a corporate network will be less than the cost of a single hour of downtime resulting from the first errant system.

The ready availability of portable storage devices, including USB memory devices and external hard drives, makes compartmentalization an even more serious issue. Exploits have shown the desirability of disabling automatic execution of software stored on plug-in devices.

30.6.6 Need to Access. Sometimes it is easy to determine who requires access, to which resources, and to which information. However, the question of access control often involves painful choices, with many nuances and subtleties. Legal and contractual responsibilities further complicate the question. For example, lack of access may exculpate persons alleged to have misused information.

Physical and logical access controls (see Chapters 23, 28, 29, 32, and 35 in this *Handbook*) need to be implemented for Internet-accessible systems as for any other sensitive and critical system. Controls for such systems must be enforced and respected by all members of the organization. It is important that personnel with restricted access to the network and security infrastructure understand and comprehend the reasons for the security rules and that they not take measures that circumvent those rules. The integrity of an organization's firewalls and network infrastructure is only as good as the

TECHNICAL ISSUES 30 · 39

physical and logical security of the personnel, equipment, infrastructure, and systems comprising the firewall and network. Regular auditing of both physical and logical access to infrastructure assets is critical and necessary. Scanning for unauthorized (rogue) Wi-Fi access points is also prudent. As was the past with dial-in modems; rogue Wi-Fi access points can compromise an organization's entire security posture.

The need to maintain information security on communications within organizations argues for the extensive use of security technologies, even when the data packets are never expected to leave the premises. It is sound planning, even within the enterprise, to require applications with privacy or confidentiality requirements to make use of privacy infrastructure such as the Secure Sockets Layer (SSL) for Web-based applications or tunneling such as Layer 2 Tunneling Protocol⁴⁴ and Point-to-Point Tunneling Protocol.⁴⁵ This helps ensure that sensitive information is limited in distribution.

Needless to say, these approaches should employ high-grade encryption and properly signed X.509 certificates from well-accepted Certification Authorities. Self-signed, expired, and not-generally-accepted certificates should not be used.

30.6.7 Accountability. People often talk about impenetrable systems. Despite the castle-and-moat analogies used in many discussions of security, perimeters are imperfect. Given the likelihood of successful attacks, security personnel must use both technical and managerial measures for effective response.

When securing infrastructure, priority should be given to protective measures that ensure accountability for actions. Just as it is desirable to prevent inappropriate activity, it is even more important to ensure that activities can be accounted for.

For example, there are many ways to carry out denial-of-service attacks. The most troublesome are those which are completely legal. Although some of these attacks, such as distributed denial of service, are belligerent or politically and ideologically motivated, involving remote-control *zombie* programs and botnets (described in Chapter 16), many accidental DoSs can occur without malice in the course of software development. For example, two of the most famous worms that inadvertently led to DoSs, the Morris worm and the WANK worm, were detected inadvertently as a result of implementation errors in their replication mechanisms. These errors caused both worms to proliferate extremely rapidly, effectively producing unintended DoS attacks and subsequent detection.

When a compromised machine is analyzed forensically, the presence of the attack may remain undetected. This is particularly true when the attack involves designer malware not in general circulation (see Chapters 16, 17, and 20 in this *Handbook*).

It is important to analyze security breaches to distinguish among attacks, accidents, and experiments. It is better for weaknesses to be uncovered, even accidentally, within an organization than it is to deal with a truly belligerent security breach. Policies and practices should therefore encourage employees to report accidents rather than try to hide them. As for false alarms caused by overenthusiastic security neophytes, management should avoid punishing those who report illusory breaches or attacks. Accountability provides the raw material to determine what actually happened. The resulting information is the critical underpinning for analysis and education, thus enabling the enterprise to evolve to higher levels of security and integrity.

30.6.8 Read-Only File Security. Many sites allow downloads, usually via FTP and HTTP, of large numbers of files. These files may include copies of forms,

30 · 40 E-COMMERCE AND WEB SERVER SAFEGUARDS

manuals, instructions, maps, and service guides. If such file serving is provided, then it is critical to ensure that:

- The servers supporting the downloading service are secure.
- The contents of the publicly accessible file store are read-only and subject to change control.
- The entire contents of the public file store can be restored quickly in the event of a possible compromise.
- There is a clearly designated party with responsibility for maintaining and protecting the public file service.

30.6.9 Going Offline. Internet-connected systems require responsiveness that is directly related to the out-of-service costs. This may be the cost of lost business; in other organizations, the cost may be that of lost professional time, of damaged public relations, and of lowered morale. In any event, the larger the proportion of the organization (or its customers) affected by the problem, the higher the cost, and the greater the urgency to effect repairs.

Today's interconnected world makes Internet disconnection a truly painful option for a network or security manager. Although the cost of disconnection (or partial disconnection) can run into hundreds of thousands of dollars in lost business and productivity, disconnection in certain situations is both necessary and appropriate. At times, disconnection presents the lowest-cost, most effective way to protect users, systems, and the public. On May 4, 2000, during the epidemic of the "I Love You," Microsoft Outlook-exploiting virus attack, network managers at Ford Motor Company disconnected Ford's network from the outside world to limit the entry of contaminated email into Ford's systems and to prevent Ford's systems from spreading the contagion.⁴⁶ The response achieved its goals. Disconnection resulted in less pain than the alternatives. Segmented networks help provide compartmentalization, permitting a more nuanced response than simple disconnection.

The primary issue surrounding disconnection is: What can be disconnected? Who has the authority? Murphy's Law often intervenes. Such important incidents require short response times, inevitably on occasions when senior managers are not available. The best way to provide for this contingency is to furnish guidelines for personnel with authority to defend the systems and a guarantee that actions within the guidelines will be immune from reprisal. If an organization fails to authorize such actions, they also forswear the benefits accruing from the steps to contain damage.

30.6.10 Auditing. In any organization, facilities usage should be monitored and analyzed, including network activity. Because people interested in attacking or subverting the enterprise's networks will attack when they wish, such monitoring and analysis must be part of a continuing process. Ongoing observation and review should include:

- Physical communications infrastructure
- Firewalls, router tables, and filtering rules
- Host security
- File security
- Traffic patterns on backbones, DMZ, and other network segments
- Physical security of systems and communications infrastructure

ETHICAL AND LEGAL ISSUES 30 · 41

These areas are both synergistic and independent. They are synergistic in that the combination of all of them is mutually reinforcing in promoting a secure computing environment. They are independent because any one of them may represent a weak point in an otherwise secure environment. For more information on auditing, see Chapter 54 in this *Handbook*.

30.6.11 Emerging Technologies. Emerging technologies continue to recast the challenges in the Internet technology arena. One such challenge is emerging in the use of SSL/TLS-based firewalls.

Using HTTPS as a basis for building encrypted tunnels appears at first impression to be a tremendously enabling technology. Most firewalls permit the initiation of connections from within the protected zone on TCP port 443 for the purpose of allowing connection to the wide variety of Websites that need to secure information, as has been mentioned in other locations throughout this chapter. For privacy and security, such traffic should be opaque to monitoring.

The use of HTTPS as the basis for tunneling also provides a tailor-made technique for abuse. For example, such technologies make it possible for a compromised desktop system to monitor the network, extracting data of interest, and to transmit the resulting data securely to an outside location through the firewall. The opaque nature of the SSL/TLS-based connection makes it impossible to scan the outgoing data.

This technology may well be the death knell of shared-connectivity LANs. It is ever more important to secure the network infrastructure against attacks that allow a rogue machine to accumulate data by monitoring network traffic.

Intrusion detection systems may also need to be repurposed to identify HTTPS connections that have significant traffic volumes, for signs of compromise or abuse. Approaches that treat the entire network as inherently compromised, with the use of IPSec VPN technology from the desktop server, to secure communications infrastructure, may be needed to protect against such attacks and malfeasance. In this context, the use of telnet for console communications is strongly discouraged. SSH and other encrypted protocols should be used.

30.7 ETHICAL AND LEGAL ISSUES. Managing a Website poses a variety of ethical and legal issues, mostly surrounding the information that the site accumulates from processing transactions as well as from performance tracking, problem identification, and auditing. Policies need to be promulgated to ensure that information is not used for unauthorized purposes, and the staff running such systems needs to be aware of, and conform to, the policies. Many of these ethical issues have civil and criminal considerations as well. The topic of information privacy is more completely addressed in Chapter 69 of this *Handbook*. For more information on security policy standards, development and implementation, see Chapters 44, 45, 47, 48, 49, 50, and 66.

30.7.1 Liabilities. The liability environment surrounding Web servers is too new for litigation to have run its course. However, there is no reason to believe that the myriad laws governing the disclosure of personal information will not be fully enforced in the context of Websites.

Websites increasingly handle sensitive information. Financial industry sites routinely handle bank and securities transactions. Email services handle large volumes of consumer traffic, and more and more insurance companies, employee benefits departments, and others are using Websites to deal with extremely sensitive information covered by a variety of regulations.

30 · 42 E-COMMERCE AND WEB SERVER SAFEGUARDS

Part of the rationale for recommending careful attention to the control of Web servers and their supporting systems is the need to create an environment of due diligence, where an organization can show that it took reasonable steps to ensure the integrity, safety, and confidentiality of information.

30.7.2 Customer Monitoring, Privacy, and Disclosure. Customer monitoring is inherently a sensitive subject. The ability to accumulate detailed information about spending patterns, for example, is subject to abuse. A valid use of this information helps to pinpoint sales offerings that a customer would find relevant while eliminating contacts that would be of no interest. Used unethically or even illegally, such information could be used to assemble a dossier that could be the subject of embarrassing disclosure, of insurance refusal, and even of job termination. The overall problem predates the Web. In fact, more than 20 years ago a major network news segment reconstructed someone's life using nothing more than the information contained in their canceled checks, supplemented with publicly available information. The resulting analysis was surprisingly detailed. A similar experiment was reported in 1999, using the Web.⁴⁷

Organizations sometimes violate the most basic security practices for protecting online information, when all of the information required to access an account improperly is contained on the single page of a billing statement. There have been repeated incidents (e.g., CDUniverse, CreditCard.com) where extremely sensitive information has been stored unencrypted on Web-accessible systems. These incidents recur with regularity and are almost always the result of storing large amounts of sensitive client information on systems that are Internet accessible. There is little question that it is inappropriate to store customer credit card, and similarly sensitive data, on exposed systems.

The security and integrity of systems holding customer order information is critical. The disclosure of customer ordering information is a significant privacy hazard. Failure to protect customer banking information (e.g., credit card numbers and expiration dates) can be extremely costly, both in economic terms and in damaged customer relations. Firms processing payment cards are also subject to the Payment Card Industry Security Standards Council's PCI DSS; which can create additional liabilities and obligations on the enterprise.

A Website, by its monitoring of customer activity, will accumulate a collection of sensitive material. It may be presumed that the information is useful only for the site, and is inherently valid, but there are a variety of hazards here, most of which are not obvious:

- Information appearing to originate from a single source may indeed be a compilation of data from multiple sources. Shared computers, firewalls, and proxy servers can give rise to this phenomenon.
- Casual correlations may arise between otherwise unrelated items. For example, it is not an uncommon acceptable business practice for one member of a business group to pay for all expenses of a group of traveling colleagues. Failure to correctly interpret such an event could be misconstrued as proof of illicit or inappropriate behavior.
- Individuals often have a variety of societal roles (e.g., parent, domestic partner, professional, nonprofit involvement). Conclusions drawn from such data may not reflect the individual, but can easily be an unrelated consequence of some other involvement (e.g., an inquiry about the SAT or GRE could indicate an intent to attend a college; it could also mean that a child, niece, neighbor, or other acquaintance could have sought counsel).

ETHICAL AND LEGAL ISSUES 30 · 43

The problem with casual associations is the damage they can cause. In the national security area, the use of casual associations to gather intelligence is a useful tool, albeit one that is recognized to have serious limitations. In other situations, it is an extremely dangerous tool, with significant potential to damage individuals and businesses. An example:

A California-based married businessman flies to New York City. When he arrives, he checks into a major hotel. A short time later, he makes a telephone call, and shortly a young woman goes up to his room and is greeted warmly. Apparently, a compromising situation. The businessman is old enough to be the woman's father. In fact, he *is* her father. That single fact changes apparently inappropriate behavior into a harmless family get-together. Peter Lewis⁴⁸ of *The New York Times* correctly notes that this single fact, easily overlooked, dramatically changes the import of the information.

The danger with correlations and customer monitoring is that there is often no control on the expansive use of the conclusions generated. The information often has some degree of validity, but it is both easy to overstep the bounds of validity, and it is difficult, if not impossible, to later correct damage, once damage has been done.

30.7.3 Litigation. The increasing pervasiveness of the Web has led to increasing volumes of related litigation. In this chapter, the emphasis is on litigation or regulatory investigation involving commercial or consumer transactions, and the issues surrounding criminal prosecution for criminal activities involving a Website. More detailed information on this subject appears in Chapter 61 of this *Handbook*.

Civil. Website logs and records can become involved in litigation in many ways. In an increasing number of cases, neither the site owner nor the operator is a party to the action; the records merely document transactions and events involved in a dispute. The dispute may be a general commercial matter, a personnel matter, or even a domestic relations matter involving divorce.

It is important that records handling and retention policies be developed in concert with counsel. The firm's management and counsel also must determine what the policy is to be with regard to subpoenas and related requests. Counsel also will determine what materials are subject to which procedures and regulations. For example, in a case of an email provider (e.g., hotmail.com), material may be subject to the Electronic Communications Privacy Act of 1986 (18 U.S.C.A. § 2510 et seq.). Other material may be subject to different legal or contractual obligations.

Regulatory. A wide range of rules is enforced (and often promulgated) by various regulatory agencies. In the United States, such agencies exist at the federal, state, regional, and local level. (Outside the United States, many nations have agencies only at the national and provincial levels.) Many of these agencies have the authority to request records and conduct various types of investigations. For many organizations, such investigations are significantly more likely than civil or criminal investigations. Regulatory agencies also may impose record-keeping and retention requirements on companies within their jurisdiction.

Criminal. Criminal prosecutions receive more attention than the preceding two categories of investigation yet are much less frequent. Criminal matters are expensive to investigate and prosecute and must pass a higher standard of proof than regulatory or civil prosecutions. Relatively few computer-related incidents reach the stage of a criminal prosecution, although because of its seriousness, the process is the most visible.

30 · 44 E-COMMERCE AND WEB SERVER SAFEGUARDS

Logs, Evidence, and Recording What Happened. The key to dealing effectively with any legal proceeding relating to a Website is the maintenance of accurate, complete records in a secure manner. This is a complex topic, some details of which are covered in Chapter 53 of this *Handbook*.

In the context of protecting a Website, records and logs of activity should be offloaded to external media and preserved for possible later use, as determined by the site's policy and its legal obligations. Once offloaded, these records should be stored using the strict procedures suitable for evidence in a criminal matter. The advent of inexpensive CD-ROM and DVD writers greatly simplifies the physical issues of securely storing such media. A cautionary note about CD-ROM and DVD media is in order: When creating such records, care should be taken to use blank media manufactured to archival standards. Archival quality media is more expensive than the more readily available media of lesser quality; however, archival media has a rated life of decades, cheaper media has a lifespan of a few years.

Archival records not intended for evidentiary use also should be stored offline. This should take the form of physically offline, or at least on systems that are not accessible from the public Internet.

The media should be stored in signed, sealed containers in an inventoried, secure storage facility with controlled access. For this reason, the copies archived for records purposes should not be the copies normally used for system recovery. In the event of an investigation or other problem, these records will be carefully examined for possible modification or misuse.

30.7.4 Application Service Providers. In recent years, it has become increasingly common to outsource entire applications services and hosting. External organizations providing such services are known as applications service providers, more commonly referred to as ASPs. More recently, this has been rechristened as *software as a service* (SaaS). Both raise similar security and integrity concerns. In both cases, significant applications and data are stored outside of the organization, with the organization retaining responsibility for the integrity and confidentiality of these records.

A more bare-boned offering, virtual hosting, has experienced an explosion of popularity, often referred to as *Infrastructure as a Service* (IaaS). Many of the same issues with SaaS apply to environments purchased as IaaS. The dividing line for responsibility is different, but the same classes of problems occur.

Conceptually, ASPs are not new. Many organizations have historically outsourced payroll processing and other applications. Theoretically, the ASP is responsible for the entire application. Often, paying a package price for the application seems attractive: no maintenance charges, no depreciation costs, lower personnel costs, latest technology, and moderately priced upgrades. However, just as a ship's captain retains responsibility for the safety of his ship despite the presence of a harbor pilot, an enterprise must not forget that if something goes wrong, the enterprise, not the ASP, will likely bear the full consequences. In short, the ASP must be required to answer the same questions, and held to the same standards, as an inside IT organization regarding privacy, security, and integrity issues.

The security and integrity issues surrounding the use of ASPs are the same as those surrounding the use of an internal corporate service. Questions of privacy, integrity, and reliability remain relevant, but as with any form of outsourcing, there are additional questions. For example, is stored information commingled with that of other firms, perhaps competitors? Is information stored encrypted? What backup provisions exist?

SUMMARY 30 · 45

Is there offsite storage of backups? What connectivity does the ASP have? Where are the ASP's servers, and are they dispersed? What are the personnel practices of the ASP? Does the ASP itself own and operate its facilities, or does it in turn contract out to other providers?

The bottom line is that although outsourcing promises speedy implementation, lower personnel costs, and economies of scale, the customer organization will suffer considerable harm if there is a problem with the ASP, with availability, results, or confidentiality.

In the end analysis, the organization retains liability for its operations and its data. While it is comforting to consider that the legal system will accept "My provider did it" as an excuse for lost, compromised, or untrustable data, it remains an untested theory and is not likely to prevail. Recent experiences with major manufacturers, subcontractors, and tainted products would seem to indicate that outsourcing risk remains a serious potential liability.

Legal process requests, either criminal warrants or civil subpoenas, against ASPs for third-party data represent another hazard. If data belonging to third parties is commingled on the ASP's system, care is required to prevent the unauthorized and inappropriate disclosure of unrelated data belonging to third parties other than the one that is the subject of the request.

None of the items cited are sufficient to justify a negative finding on ASPs as a group. They should, however, serve as reminders that each of the issues related to keeping a Web presence secure, available, and effective apply no less to an ASP than they do to an in-house IT organization.

For information about outsourcing and security, see Chapter 68 in this *Handbook*.

30.8 SUMMARY. Availability is the cornerstone of all Web-related strategies. Throughout this chapter, it has been noted that redundant hosting and routing were necessary to ensure 24/7/365 availability. It was also noted that although some providers of services offer various guarantees, these guarantees almost never provide adequate compensation for consequential damage done to the enterprise. In the end, an enterprise's only protection is to take adequate measures to ensure their own security and integrity.

Operating guidelines and authority are also critical to ensuring the availability of Web resources on a 24/7/365 basis. Systems must be architected and implemented to enhance availability on an overall level. Operating personnel must have the freedom and authority to take actions that they perceive as necessary without fear of reprisal if the procedures do not produce the desired outcome.

Privacy and integrity of information exchanged with the Website is also important. The implementation and operation of the site and its components must be in compliance with the appropriate laws, regulations, and obligations of the site owner, in addition to being in conformance with the expectations of the user community.

Protecting Internet and Web assets is a multifaceted task encompassing many disciplines. It is an area that requires attention at all levels, beginning at the highest levels of business strategy and descending successively to ever more detailed implementation and technology issues.

It is also an area where the smallest detail can have catastrophic impact. Recent events have shown that the lessons of communications history apply. Even a small error in network architecture, key management, or implementation can snowball until it is an issue that can be felt in the corporate boardroom.

30 · 46 E-COMMERCE AND WEB SERVER SAFEGUARDS

30.9 FURTHER READING

- Anderson, R. *Security Engineering: A Guide to Building Dependable Distributed Systems*, 2nd ed. Wiley, 2008.
- Bhargav, A., and B. V. Kumar. *Secure Java: For Web Application Development*. CRC Press, 2010.
- Cobb, Stephen. *Privacy for Business: Web Sites and Email*. Dreva Hill, 2002.
- Elliott, E. *Programming JavaScript Applications: Robust Web Architecture with Node, HTML5, and Modern JS Libraries*. O'Reilly Media, 2013.
- Gabarro, S. A. *Web Application Design and Implementation: Apache 2, PHP5, MySQL, JavaScript, and Linux/UNIX*, 2nd ed. Wiley, 2013.
- Garfinkel, S. *Web Security, Privacy and Commerce*, 2nd ed. O'Reilly Media, 2002.
- Habiyaremye, J d'A. *E-Commerce Security Threats*. GRIN Verlag, 2011.
- Hassler, V. *Security Fundamentals for E-Commerce*. Artech House, 2000.
- Kissoon, T. *Securing Web Applications*. Auerbach, 2013.
- Long, F., D. Mohindra, R. C. Seacord, D. F. Sutherland, and D. Svoboda. *The CERT Oracle Secure Coding Standard for Java*. Addison-Wesley Professional, 2011.
- McGraw, G., and E. Felten. *Securing Java: Getting Down to Business with Mobile Code*, 2nd edition. Wiley, 1999.
- Musciano, C., B. Kennedy, and E. Weyl. *HTML5: The Definitive Guide*, 7th ed. O'Reilly Media, 2014.
- Re.think. *The Book on Professional Website Design*. The Bestselling Book.com, 2013.
- Scambray, J., Shema, M., and C. Sima. *Hacking Exposed Web Applications*, 2nd ed. McGraw-Hill Osborne Media, 2006.
- Sullivan, B. *Web Application Security, A Beginner's Guide*. McGraw-Hill Osborne Media, 2011.
- Welling, L., and L. Thomson. *PHP and MySQL Web Development*, 5th ed. Addison-Wesley Professional, 2013.
- Zalewski, M. *The Tangled Web: A Guide to Securing Modern Web Applications*. No Starch Press, 2011.

30.10 NOTES

1. J. Swartz, "TJX Data Breach May Involve 94 Million Credit Cards," *USA Today*, October 24, 2007, www.usatoday.com/money/industries/technology/2007-10-24-tjx-security-breach_N.htm
2. J. Vijayan, "TJX Agrees to Pay \$40.9M to Visa Card Issuers in Breach Case," *PCWorld*, November 30, 2007, www.pcworld.com/article/140174/article.html
3. T. Dierks and E. Rescoria. "RFC 5246: The Transport Layer Security (TLS) Protocol, v 1.2." *Internet Engineering Task Force, Network Working Group*. Updated by RFCs 5746, 5878, 6176. <http://tools.ietf.org/html/rfc5246>
4. A. Freier, P. Karlton, P. Kocher. "The Secure Sockets Layer (SSL) Protocol Version 3.0." *Internet Engineering Task Force, Network Working Group*. <http://tools.ietf.org/html/rfc6101>
5. E. Murray, "Survey on Information Security," *Information Security Magazine* (October 2000).
6. VISA, "PCI Compliance Continued to Grow in 2007," Press Release, January 22, 2008. <http://corporate.visa.com/newsroom/press-releases/press753.jsp>

NOTES 30 · 47

7. Dan Kaplan, "Hannaford Tells Regulators How Breach Happened," *SC Magazine*, April 1, 2008, www.scmagazine.com/hannaford-tells-regulators-how-breach-happened/article/108569/
8. Jennifer Thomas, "Cyber Attack Hits Bashas' Chain of Stores," *azfamily.com*, February 5, 2013, www.azfamily.com/news/Cyber-attack-hits-Bashas-chain-of-stores-189911491.html
9. Ann Bednarz, "Target's Website Woes Getting Fixed, CEO Says," *Network World*, November 18, 2011, www.networkworld.com/news/2011/111811-target-website-253296.html
10. Stephanie Clifford, "Demand at Target for Fashion Line Crashes Web Site," *The New York Times*, September 13, 2011, www.nytimes.com/2011/09/14/business/demand-at-target-for-fashion-line-crashes-web-site.html
11. S. Schiesel and R. Atlas, "By-the-Numbers Operation at Morgan Stanley Finds Its Human Side," *The New York Times*, September 16, 2001, www.nytimes.com/2001/09/16/business/by-the-numbers-operation-at-morgan-stanley-finds-its-human-side.html
12. A. Harmon (2001, October 29) "Breaking Up the Central Office; Staffs Make a Virtue of Necessity," *The New York Times*, October 29, 2001, corrected October 30, 2001. www.nytimes.com/2001/10/29/business/breaking-up-the-central-office-staffs-make-a-virtue-of-necessity.html
13. TRUSTe. www.truste.org
14. P. G. Neumann, "The Green Card Flap," *RISKS Forum Digest* 15, April 18, 1994, <http://catless.ncl.ac.uk/Risks/15.76.html#subj2>
15. Choe Sang-Hun, "Computer Networks in South Korea Are Paralyzed in Cyber-attacks" *The New York Times*, March 20, 2013, www.nytimes.com/2013/03/21/world/asia/south-korea-computer-network-crashes.html
16. John Markoff, "Before the Gunfire, Cyberattacks," *The New York Times*, August 12, 2008, www.nytimes.com/2008/08/13/technology/13cyber.html
17. "Estonia Hit by 'Moscow Cyber War,'" *BBC News*, last updated May 17, 2007, <http://news.bbc.co.uk/2/hi/europe/6665145.stm>
18. G. Keizer, "Hacker Contest Weekend," *Web Design & Technology News*, July 2, 2003, www.webdesignsnow.com/news/070203.html
19. Extensible Markup Language (XML), www.w3.org/XML
20. J. Serjeant, "'Botmaster' Admits Infecting 250,000 Computers," *Reuters*, November 9, 2007, www.reuters.com/article/domesticNews/idUSNO823938120071110; Press Release No. 07-143, United States Attorney's Office, Central District of California, www.justice.gov/usao/cac/Pressroom/pr2007/143.html; FBI, "'Bot Roast II' Nets 8 Individuals," www.fbi.gov/news/pressrel/press-releases/bot-roast-ii-nets-8-individuals
21. Camilla Hall and Javier Blas, "Aramco cyber attack targeted production," *Financial Times*, December 10, 2012, www.ft.com/cms/s/0/5f313ab6-42da-11e2-a4e4-00144feabdc0.html
22. Shaun Waterman, "South Korea Cyberattack Not Linked to China After All," *The Washington Times*, March 22, 2013, www.washingtontimes.com/news/2013/mar/22/south-korea-cyberattack-not-linked-china-after-all
23. Brian Grow, Keith Epstein, and Chi-Chu Tschang, "The New E-spyionage Threat" *Bloomberg BusinessWeek*, April 9, 2008, www.businessweek.com/stories/2008-04-09/the-new-e-spyionage-threat

30 · 48 E-COMMERCE AND WEB SERVER SAFEGUARDS

24. Robert Gezelter, "Protecting Web Sites," *Computer Security Handbook*, 4th ed. (2002), Chapter 22, Section 22.3.3, pp. 22–27.
25. TechWise Research, Inc., "Quantifying the Value of Availability," June 2000, <http://techwise-research.com/member-resources/white-papers/id/43/quantifying-the-value-of-availability>
26. J. Vijayan, "TJX Agrees to Pay \$40.9M to Visa Card Issuers in Breach Case," *PCWorld*, November 30, 2007, www.pcworld.com/article/140174/article.html
27. C. Morris, "Sony: PlayStation Breach Involves 70 Million Subscribers," CNBC, April 26, 2011, www.cnbc.com/id/42769019
28. Y. Rekhter, B. Moskowitz, D. Karrenberg, G. J. de Groot, and E. Lear, "Address Allocation for Private Internets," IETF RFC 1918, February 1996, <http://tools.ietf.org/rfc/rfc1918.txt>
29. R. Hinden and B. Haberman, "Unique Local IPv6 Unicast Addresses," RFC 4193, October 2005, <http://tools.ietf.org/rfc/rfc4193.txt>
30. R. Gezelter, "Security Guidelines for New and Existing OpenVMS Applications," November 12, 1996, www.rlgsc.com/decus/usf96/ad020.pdf
31. N. D. Schwartz and E. Dash, "Thieves Found Citigroup Site an Easy Entry," *The New York Times*, June 14, 2011, page A1, www.nytimes.com/2011/06/14/technology/14security.html?pagewanted=all
32. R. Gezelter, "Internet Dial Tones & Firewalls: One Policy Does Not Fit All," *IEEE Computer Society*, Charleston, SC, June 10, 2003, www.rlgsc.com/ieee/charleston/2003-6/internetdial.html; R. Gezelter, "Safe Computing in the Age of Ubiquitous Connectivity," LISAT 2007, May 2007, www.rlgsc.com/ieee/longisland/2007/ubiquitous.html
33. RFC1597, <http://tools.ietf.org/html/rfc1597>; superseded by RFC1918 <http://tools.ietf.org/html/rfc1918>
34. RFC 3927, <http://tools.ietf.org/html/rfc3927>
35. Robert Gezelter, "Internet Dial Tones & Firewalls."
36. Robert Gezelter, "Safe Computing in the Age of Ubiquitous Connectivity."
37. Formerly RFC 1597.
38. Gezelter, "Internet Dial Tones & Firewalls." Gezelter, "Safe Computing in the Age of Ubiquitous Connectivity." Gezelter, "Internet Dial Tones & Firewalls." Gezelter, "Safe Computing in the Age of Ubiquitous Connectivity."
39. RFC 3927.
40. Robert Gezelter, "Stopping Spoofed Addresses Can Cut Down on DDoS Attacks," *Network World Fusion*, August 14, 2000, www.networkworld.com/columnists/2000/0814gezelter.html
41. Frank da Cruz and Christine M. Gianone, "Using C-KERMIT," 2nd ed. Boston: Digital Press, 1997.
42. M. Slatalla and J. Quittner, "Masters of Deception," New York: Harper Collins, 1995.
43. R Gezelter, "Networks Placed At Risk: By Their Providers," Ruminations—An IT Blog, December 2009, www.rlgsc.com/blog/ruminations/networks-placed-at-risk.html
44. W. Townsley, A. Valencia, A. Rubens, G. Pall, G. Zorn, and B. Palter, "RFC 2661—Layer Two Tunneling Protocol 'L2TP,'" August 1999, <http://tools.ietf.org/html/rfc2661>

NOTES 30 · 49

45. K. Hamzeh, G. Pall, W. Verthein, J. Taarud, W. Little, and G. Zorn, “RFC 2637—Point-to-Point Tunneling Protocol,” July 1999, <http://tools.ietf.org/html/rfc2637>
46. K. Timms, telephone interview, summer 2001; K. Bradsher, “With Its Email Infected, Ford Scrambled and Caught Up,” *The New York Times*, May 8, 2000, www.nytimes.com/2000/05/08/business/with-its-e-mail-infected-ford-scrambled-and-caught-up.html
47. T. Kelley, “An Expert in Computer Security Finds His Life Is a Wide Open Book,” *The New York Times*, December 31, 1999, www.nytimes.com/1999/12/13/business/technology-an-expert-in-computer-security-finds-his-life-is-a-wide-open-book.html
48. P. Lewis, “Forget Big Brother,” *The New York Times*, March 19, 1998, www.nytimes.com/1998/03/19/technology/forget-big-brother.html?pagewanted=all&src=pm

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 31

WEB MONITORING AND CONTENT FILTERING

Steven Lovaas

31.1	INTRODUCTION	31·1	31.5	IMPLEMENTATION	31·8
31.2	SOME TERMINOLOGY	31·2		31.5.1 Manual “Bad URL” Lists	31·8
31.3	MOTIVATION	31·2		31.5.2 Third-Party Block Lists	31·9
31.3.1	Prevention of Dissent	31·3	31.6	ENFORCEMENT	31·9
31.3.2	Protection of Children	31·3	31.6.1	Proxies	31·9
31.3.3	Supporting Organizational Human Resources Policy	31·3	31.6.2	Firewalls	31·9
31.3.4	Enforcement of Laws	31·4	31.6.3	Parental Tools	31·10
31.3.5	National Security Surveillance	31·4	31.7	VULNERABILITIES	31·10
			31.7.1	Spoofing	31·10
			31.7.2	Tunneling	31·11
			31.7.3	Encryption	31·11
			31.7.4	Anonymity	31·12
			31.7.5	Translation Sites	31·12
			31.7.6	Caching Services	31·13
31.4	GENERAL TECHNIQUES	31·4	31.8	THE FUTURE	31·13
31.4.1	Matching the Request	31·5	31.9	SUMMARY	31·14
31.4.2	Matching the Host	31·5	31.10	FURTHER READING	31·14
31.4.3	Matching the Domain	31·7	31.11	NOTES	31·14
31.4.4	Matching the Content	31·7			

31.1 INTRODUCTION. The Internet has been called a cesspool, sometimes in reference to the number of virus-infected and hacker-controlled machines, but more often in reference to the amount of objectionable content available at a click of the mouse. This chapter deals with efforts to monitor and control access to some of this content. Applications that perform this kind of activity are controversial: Privacy and free-speech advocates regularly refer to *censorware*, while the writers of such software tend to use the term *content filtering*. This chapter uses *content filtering*, without meaning to take a side in the argument by so doing. For more on the policy and legal issues surrounding Web monitoring and content filtering, see Chapters 48 and 72 in this *Handbook*.

31 · 2 WEB MONITORING AND CONTENT FILTERING

This chapter briefly discusses the possible motivations leading to the decision to filter content, without debating the legitimacy of these motives. Given the variety of good and bad reasons to monitor and filter Web content, this chapter reviews the various techniques used in filtering, as well as some ways in which monitoring and filtering can be defeated.

31.2 SOME TERMINOLOGY

Proxy—a computer that intercedes in a communication on behalf of a client. The proxy receives the client request and then regenerates the request with the proxy’s own address as the source address. Thus, the server only sees identifying information from the proxy, and the client’s identity remains hidden (at least from a network addressing point of view). Proxies are widely used by organizations to control the outward flow of traffic (mostly Web traffic) and to protect users from direct connection with potentially damaging Websites.

Anonymizing proxy—a proxy that allows users to hide their Web activity. Typically, such a proxy is located outside organizational boundaries, and often used to get around filtering rules.

Privacy-enhancing technologies (PET)—a class of technologies that helps users keep their network usage private. These include encryption, anonymizing proxies, mixing networks, and onion routing.

Encryption—reversible garbling of text (plaintext) into random-looking data (ciphertext) that cannot be reversed without using secret information (keys). See Chapter 7 in this *Handbook* for details of cryptography.

Mixing networks—users encrypt their outbound messages with their recipient’s (or recipients’) public key(s) and then also with the public key of a “Mix server.” The encrypted message is sent to the Mix server, which acts as a message broker and decrypts the cryptogram to determine the true recipient’s address and to forward the message.¹

Onion routing—an anonymous and secure routing protocol developed by the Naval Research Laboratory Center for High Assurance Computer Systems in the 1990s. Messages are sent through a network (cloud) of *onion routers* that communicate with each other using public key cryptography. Once the temporary circuit has been established, the sender encrypts the outbound message with each of the public keys of all the onion routers in the circuit.

Specifically, the architecture provides for bi-directional communication even though no one but the initiator’s proxy server knows anything but previous and next hops in the communication chain. This implies that neither the respondent nor his proxy server nor any external observer need know the identity of the initiator or his proxy server.²

A third-generation implementation of onion routing is *Tor* (the onion router).³

31.3 MOTIVATION. As a general concept, an individual or group chooses to monitor network activity and filter content through an authority relationship. The most common relationships leading to monitoring and filtering are:

- Governments controlling their citizenry
- Parents and schools protecting their children
- Organizations enforcing their policies

MOTIVATION 31 · 3

- Governments enforcing their laws
- Governments monitoring terrorists and other nation states

31.3.1 Prevention of Dissent. Probably the most common reason for the opposition to content filtering is that many repressive governments filter content to prevent dissent. If citizens (or subjects) cannot access information that either reflects badly on the government or questions the country's philosophical or religious doctrines, then the citizens probably will not realize the degree to which they are being repressed.

In countries such as the United States, where the Constitution guarantees freedoms of speech and the press, many people feel that only those with something to hide about their activities would fight against censorship. This notion, combined with objectionable content so readily available on the Internet, may be why efforts to question the use of filtering software in libraries and other public places have met with resistance or indifference. Chapter 72 gives many examples of countries using filtering products to limit the rights of their citizens. These examples should be a cautionary tale to readers in countries that are currently more liberal in information policy.

31.3.2 Protection of Children. For the same reasons that convenience stores, at least partially, hide adult magazines from view, the government has required that some classes of information on the Internet be off-limits to children in public schools. This stems from the perception of the schools' role as a surrogate parent and government liability involved in possible failure in that role. A variety of regulations have required the use of content monitoring and filtering in schools and libraries, on the theory that such public or publicly supported terminals should not be using taxpayers' money to provide objectionable content to children. The U.S. Supreme Court has ruled several of these efforts unconstitutional after extreme protest from libraries. Most recently, the Child Internet Protection Act (CIPA)⁴ required school districts that receive certain kinds of government funding to use filtering technologies. Most schools have implemented Web filtering, aggressively limiting the Web content available to students. Estimates of the amount of blocked content vary widely, but one school district claims to filter about 10 percent of its total Web traffic due to questionable content.⁵

The obvious target of filtering technology in schools is to keep students from viewing material considered *harmful to minors*. Another less-publicized reason to filter Web content is to prevent students from doing things on the Web that they might not do if they knew that someone was watching. The hope is to keep students from getting a criminal record before they even graduate from high school, protecting them from their bad judgment while they are learning to develop good judgment and learning the rules of society. Arguably, this is the job of a parent and not the job of a school, but schools providing Internet access nonetheless must provide this kind of protection.

High school students and kindergarteners generally have differing levels of maturity, and a reasonable filtering policy would include a flexible, gradual degree of filtering depending on age. Schools could use less-intrusive methods of controlling access to Web content as well, such as direct supervision of students using computers. Parents also have the opportunity to filter the content that their children view at home.

31.3.3 Supporting Organizational Human Resources Policy. Organizations have a variety of reasons for monitoring and filtering the Web content accessed by their employees. The simplest is a desire to keep employees doing work-related activity while at work. Despite studies indicating that the flexibility to conduct limited personal business, such as banking or email, from work produces happier, more

31 · 4 WEB MONITORING AND CONTENT FILTERING

productive workers, managers sometimes view personal use of business computers as stealing time from the company.

A more pragmatic reason to filter some Web content is to prevent “hostile workplace” liability under Equal Employment Opportunity laws. The problem with this approach is that many kinds of content would potentially be offensive to coworkers, so it is difficult to use filtering technology to guarantee that no one will be offended while still allowing reasonable work-related use of the Internet.

Some organizations choose to monitor Web traffic with automated systems, rather than blocking anything, and to notify users of the monitoring. The notification could be in general policy documents or in the form of a pop-up window announcing that the user is about to view content that might violate policy. Either way, the notion is that the organization might avoid liability by having warned the user, but might also avoid privacy and freedom-of-speech complaints. The monitor-and-notify approach also sends the message that the organization trusts its employees, but wants to maintain a positive and productive work environment.

31.3.4 Enforcement of Laws. Law enforcement agencies rarely engage in content filtering, but traffic monitoring is an often-utilized tool for investigating computer crime. Gathering evidence for these investigations often involves catching the traffic as it actually arrives at its destination, so filtering would be counterproductive. In this case, proving identity is the key to getting usable evidence, so privacy-enhancing technologies are real problems. In some cases, the logs of Web proxies—which would identify the real source address of the client machine—are available with a subpoena. Investigations of this sort often target child pornography, drug production, or the theft of computer hardware protected by asset-recovery software.

31.3.5 National Security Surveillance. While government law enforcement agencies work to control activity that violates the nation’s laws, other government agencies are concerned with the growth of information and communications technologies as tools of terror and international conflict. From the threat of eventual cyber-war to the more immediate concern with terrorism, governments are exploring ways to detect the use of email, cell phones, social media platforms, and even online games to communicate and coordinate plans for hostile activities. Monitoring these media has become a pressing, expensive, and risky priority for national governments.

As this edition was going to press, secret information was leaked about communications-monitoring activities of the U.S. Government’s National Security Agency, rekindling a perennial debate on the balance between privacy and safety in a democratic republic. Despite long-standing precedent and assurances about constitutional protection from unreasonable search, the revelations about the NSA’s PRISM and related programs have made it clear that networked communication is fair game in a nation’s efforts to stay one step ahead of its adversaries (whether they are real, potential, or even imagined ones).⁶ Furthermore, as the U.S. citizenry has discovered, it is very difficult for a government to effectively eavesdrop on hostile foreign electronic communications without gathering *all* communications, at the risk of violating privacy expectations of citizens and allies alike.

31.4 GENERAL TECHNIQUES. Filtering of networked communication can employ two basic tactics: examination of *metadata* and examination of the actual content of a message. The U.S. news media rediscovered the concept of metadata in the wake of the PRISM leak, as NSA representatives insisted that actual message content would

GENERAL TECHNIQUES 31 · 5

only be examined after filtering out irrelevant communications by looking at phone numbers, call durations, IP addresses, and other transactional information that reveals the source, destination, and temporal context of a conversation.⁷ Based on the overwhelming volume of communications traffic traversing the copper wires, fiber-optic cables, and wireless frequencies, it is indeed a daunting task to make filtering (or legal or military) decisions based on examining every word of every message. It is much more efficient, for most purposes, at least to start with these metadata and make assumptions about the relevance or danger of a particular computer, human, or organizational communication partner. Direct examination of the full content of messages is much more resource-intensive, though advances are being made in effective and efficient automated processing of text and images.

31.4.1 Matching the Request. The simplest technique used by filtering technologies is the matching of strings against lists of keywords. Every Web request uses a uniform resource locator (URL) in the general form:

protocol://server.organization.top-level-domain/path-to-file.file-format

Filtering a request for a URL can examine any portion of the string for a match against prohibited strings:

- Filters can match the *protocol* field to enforce policies about the use of encrypted Web traffic (HTTP versus HTTPS). This is more of a general security concern than it is a Web filtering issue, although an organization worried about the need to analyze all traffic could prevent the use of encrypted Web traffic by blocking all HTTPS requests.
- The *server* and *organization* fields describe who is hosting the content. Filtering based in these strings is a broad approach, as it leads to blocking either an entire Web server or an entire organization's Web traffic. See the discussion of server blocking in Section 31.4.2.
- The *top-level-domain* field can be used to filter content; see Section 31.4.3 for attempts to set up an .xxx domain.
- The *path-to-file* field includes the actual title of the requested Web page, so it varies most between individual requests. Whether it is more likely than other fields to contain information useful in filtering depends on the naming convention of the server. This field (and the *file-format* field) is optional, as shown in the request for www.wiley.com, which directs the server to display its default page.
- The *file-format* field tells the Web browser how to handle the text displayed in the page, whether in straight html format or encoded in some other file format (e.g., doc, pdf), or whether to allow dynamic generation of content (e.g., asp). Few filtering products use this field to filter traditional kinds of objectionable content, although enforcing other policies regarding dynamic code in high-security environments can require matching this field.

31.4.2 Matching the Host. Some filtering systems attempt to distinguish between acceptable and unacceptable sources on the Web by inspecting the particular servers or general information portals such as search engines.

31 · 6 WEB MONITORING AND CONTENT FILTERING

31.4.2.1 Block Lists of Servers. Objectionable content tends to be concentrated on individual servers. This naturally leads some organizations to block access to those servers. The two methods of blocking servers are by Internet Protocol (IP) address and by name.

Blocking by IP address simply denies all traffic (or all HTTP traffic) to and from certain addresses. This tactic involves several difficulties based on how addresses are used. First, IP addresses are not permanent. While the numeric addresses of most large commercial servers tend to remain the same over time, many smaller servers have dynamically assigned addresses that may change periodically. Thus, blocking an IP address may prevent access to content hosted on a completely different server than intended. Second, commercial servers often host content for many different customers, and blocking the server as a whole will block all of the contents rather than just the objectionable ones. This is a particular problem for very large service providers like AOL, which grants every user the ability to host a Website. Blocking the AOL servers because of objectionable content would potentially overblock large numbers of personal Web pages. Third, address-based blocking creates the possibility of malicious listing, a practice in which an attacker (or a competitor) spoofs the address of a Web server and provides objectionable content likely to land the server on blocking lists. Some filtering products allow users to submit “bad” sites, providing another opportunity for malicious listing.

Blocking by name involves the Domain Name System (DNS), in which a human-readable name (e.g., www.wiley.com) maps to a computer-readable IP addresses (in this case, 208.215.179.146). DNS allows an organization to change the physical address of its Web server by updating the DNS listing to point to the new IP address, although this does rely on the system as a whole propagating the change in a reasonable amount of time. A device that monitors or filters traffic based on a domain name needs to be able to periodically refresh its list of name-to-address mappings in order to avoid blocking sites whose addresses periodically change. As with address-based blocking, name-based blocking also risks malicious listing as well as both under- and overblocking. Many organizations register multiple names for their servers, for a variety of reasons. For instance, an organization might register its name in the .net, .com, .org, and .biz top-level domains to prevent the kind of misdirection described in Section 31.44 (whitehouse.com). Other reasons for registering a name in multiple domains include preventing competitors from using name registration to steal customers and preventing speculative buying of similar names by individuals hoping to sell them for significant profit. Web hosting companies also provide service to many different organizations, so a variety of URLs would point to the IP address of the same hosting server. Thus, blocking by server name could underblock by not accounting for all the possible registered names pointing to the same server and overblock by matching the server name of a service provider that hosts sites for many different customers who use the provider’s server name in the URL of their Web pages.

31.4.2.2 Block/Modify Intermediaries. The Web has become an enormous repository of information, requiring the development of powerful search tools to find information. Early search engines gave way to more sophisticated information portals like Yahoo!, Google, AOL, and MSN. By allowing advanced searches and customized results, these portals give users easy access to information that would be difficult or impossible to find using manual search techniques. Portals have become wildly popular tools for accessing the Internet as a whole; Google claimed 100 million search queries

GENERAL TECHNIQUES 31 · 7

per day at the end of 2000; in 2013, its site index was estimated to contain over 48 billion pages.⁸ Information access at this scale makes portals natural targets for monitoring and filtering. Few commercial organizations prevent their employees from using popular portals, since they have become so much a part of the way people use the Internet. Some countries, however, have blocked access to certain portals for their citizens, hoping to control access to information that might tend to violate national laws (e.g., access to Nazi memorabilia in France) or inspire citizen resistance to the government (e.g., access to information about the Tiananmen Square massacre in China). For further discussion of these issues, see Chapter 72 in this *Handbook*.

31.4.3 Matching the Domain. From 2004 until its final decision in 2011, the Internet Corporation for Assigned Names and Numbers (ICANN) considered and rejected requests for a new top-level domain, .xxx, which would allow providers of sexually related content to voluntarily reregister their sites. Such a domain, it was argued, would be easy to filter in the URL of the request, which presumably would appeal to supporters of Web filtering. It would also allow content providers to show that they were complying with laws preventing children from accessing inappropriate material, by enabling more effective parental filtering. Nevertheless, the move met resistance on both fronts. Conservative religious groups feared that establishing a .xxx domain would legitimize pornography, while not all sexual content providers agreed that the perceived benefit would outweigh either the increased filtering of their sites or the easier monitoring of their clients' traffic.

ICANN rejected several revisions of the .xxx domain proposal over the years, citing the lack of unanimity in the sex content provider community as well as the fear that ICANN might be placed in the position of regulating content, which is outside the organization's charter. Ultimately, however, in 2011, the .xxx domain was approved.⁹

31.4.4 Matching the Content. String matching is simple to do, but difficult to do without both over- and underblocking. For instance, one of the most common categories for content filtering (particularly in the United States) is sex. Blocking all content exactly matching the word "sex" would fail to match the words "sexy" and "sexual." To avoid this kind of underblocking, word lists need to be very long to account for all permutations. A slightly more effective tactic is to block all works containing the string "sex," but this would overblock the words "Essex," "Sussex," and "asexual." Looking for all strings beginning with the combination "sex" would overblock "sextion," "sextet," and "sextant." Simple string matching also ignores context, so blocking "sex" would match in cases where a survey page asked the respondent to identify gender using the word "sex" or in pages describing inherited sex traits or gender roles or sexual discrimination lawsuits.

Other difficulties in string matching involve the vagaries of language. URLs can be displayed in any language whose character set a computer recognizes, so a filter will underblock requests in a language for which it lacks word lists. More generally, in any language it is possible to obfuscate the contents of a site with a seemingly benign URL to avoid filtering. The classic example of this is the pornography site www.whitehouse.com, presumably set up to catch visitors who mistakenly typed "com" when trying to reach the U.S. White House Website (www.whitehouse.gov). More recently, spam marketing campaigns have been setting up Websites linked in email, with meaningless strings of numbers and characters in the URL (e.g., <http://2sfh.com/7hioh>), making the sites difficult to filter.

31 · 8 WEB MONITORING AND CONTENT FILTERING

In a 2006 study, Veritest compared three of the industry-leading Web filter products (WebSense, SmartFilter, and SurfControl). The winning product underblocked 7 sites, overblocked 8 sites, and miscategorized 10 sites, from a preselected list of 600 URLs. The two competing products fared worse, underblocking 23 and 14 sites, and overblocking 9 and 12 sites out of 600.¹⁰ A meta-study of filtering effectiveness studies performed in 2010 determined that from 2001 through 2008, the average accuracy of filtering products was 78 percent, with some increasing success: During 2007 and 2008, surveys found that the success rate had risen to 83 percent.¹¹ If this is the performance of the industry's leading edge, then clearly the technology is still developing.

Given the difficulties of accurate matching based on text or address related to the Web page request, a natural alternative is to examine the page content itself. Of course, content matching needs to have access to the unencrypted data in transit, so encrypted Web sessions cause a real problem for this tactic. Some organizations allow (or require) that HTTPS sessions terminate on the organization's own proxy server, potentially allowing the proxy to decrypt the data and perform content analysis.

31.4.4.1 Text. It is possible, although resource intensive, to watch the network traffic stream and look for text that matches a list of undesired content. This sort of matching typically does little analysis of context and so is prone to the same kind of false positives (overblocking) and false negatives (underblocking) described in Section 31.4.2. Moreover, as an increasing amount of Web content involves pictures and sounds, text matching becomes less effective.

31.4.4.2 Graphics. A promising new technique, with applications in visual searching as well as visual content blocking, breaks a graphic image into smaller objects by color or pattern. The technique then evaluates each object against a database of reference images for matching with desired criteria. In the case of blocking objectionable sex content, objects can be evaluated for skin tone and either blocked outright or referred to administrators for manual review if the match is inconclusive. Although content-based filtering has not yet developed into a commercial product, the tools exist and the technology seems applicable not only to still images, but also to video and even audio content.¹² In 2006, the NASA Inspector General's office used an image-search program called Web ContExt to snare an employee who had been trafficking in child pornography.¹³

31.5 IMPLEMENTATION. With the exception of content-based matching, which has not yet reached the market in any significant way, most filtering—whether of address, domain, or keyword—involves matching text lists.

31.5.1 Manual “Bad URL” Lists. Many firewalls provide the capability for administrators to block individual URLs in the firewall configuration. Entered manually, these rules are good for one-time blocking when a security alert or investigation identifies sites hosting viruses or other malware. This approach is also useful for demonstrating the general filtering abilities of the firewall and for testing other Web-blocking technologies. For instance, in an organization using a commercial blocking solution on a Web proxy, a simple URL-blocking rule on the organization's border firewall would provide some easy spot testing of the effectiveness of the commercial solution. Given the extreme size and constant growth of the Web, however, the manual approach does not scale well to protect against all the possible sources of objectionable material.

ENFORCEMENT 31 · 9

31.5.2 Third-Party Block Lists. With the enormous size of the Web, the more typical approach is to use a third-party block list. Most of these are commercial products with proprietary databases, developed through a combination of automated “Web crawlers” and human technicians evaluating Websites. Some companies have attempted to prevent researchers from trying to learn about blocking lists and strategies, but the U.S. Copyright Office granted a Digital Millennium Copyright Act (DMCA) exemption in 2003 for fair use by researchers studying these lists.¹⁴ Web-filter companies continue to oppose such exemptions. Two open-source filtering alternatives also exist, with publicly viewable (and customizable) block lists that run on caching proxies: SquidGuard¹⁵ and DansGuardian.¹⁶

31.6 ENFORCEMENT. Filtering of Web traffic typically occurs either at a network choke point, such as a firewall or Web proxy, or on the individual client machine. Economies of scale lead organizations to filter on a network device, while products designed for parental control of children’s Internet use usually reside on individual home computers.

31.6.1 Proxies. A proxy server is a device that accepts a request from a client computer and then redirects that request to its ultimate destination. Proxies serve a variety of purposes for an organization, including reduction of traffic over expensive wide-area network links and Internet connections, increased performance through caching frequently accessed Web pages, and protection of internal users through hiding their actual IP addresses from the destination Web servers. Proxies also represent a natural locus of control for the organization, enabling authentication and tracking of Web requests that go through this single device. Most browsers support manual configuration of a proxy for all Web traffic as well as automatic discovery of proxies running on the organization’s network. Organizations that use Web proxies typically allow outbound Web traffic only from the IP address of the proxy, thus forcing all HTTP traffic to use the proxy. Use of an encrypted Web session (HTTPS) is possible through a proxy, although either at the expense of the ability to monitor content (if the proxy merely passes the traffic through) or at the expense of the end-to-end privacy of the encrypted link (if the proxy decrypts and re-encrypts the session).

Individuals also use proxies to maintain the privacy of their activities on the network, as described in Section 31.7.4. Thus, in addition to serving as a natural vehicle for content-filtering applications, proxies also represent a serious threat to those same applications.

31.6.2 Firewalls. A firewall’s job is to analyze information about the traffic passing through it and apply policy rules based on that information. Maintaining acceptable response time and throughput requires that the firewall do its job quickly and efficiently. In order to do so, most firewalls merely look at network-layer information, such as source and destination addresses and ports. More recently, firewall vendors have been adding more features to increase security and product appeal. Many companies now call their more advanced firewalls “service gateways” or “security gateways” as the notion of Unified Threat Management (UTM) becomes more popular. These UTM devices combine many features that formerly required individual devices, such as anti-virus, intrusion detection, and filtering of both junk email and Web content.

Sophisticated traffic examination increases the demand on firewall hardware. In order to reduce the performance hit caused by increased packet inspection, many

31 · 10 WEB MONITORING AND CONTENT FILTERING

firewalls allow the administrator to define particular rules or protocols for advanced checking. For instance, since viruses are most prevalent in email, Web, and peer-to-peer connections, the firewall administrator might need to configure only anti-virus checking on rules applying to these protocols. Similarly, if the firewall needs only to monitor outbound HTTP requests from a single IP address, the Web proxy, then the extra processing load of the monitoring function can be constrained to that traffic profile.

The decision between filtering Web traffic at the proxy server (letting the firewall just pass the traffic from that address) and filtering Web traffic at the firewall (having the firewall do the URL inspection) depends on the amount of Web traffic and the budget (one device or two). The decision also affects the strength of the assertion that the organization is successfully filtering objectionable content. If the organization's border firewall performs the filtering, then this assertion depends on the firewall being the only way for traffic to leave the organization's network. Other traffic vectors, including wireless networking, protocol tunneling, and anonymizing proxies, may come into play. If the organization relies on client computers to use a Web proxy by policy, then the degree to which users can circumvent this policy should also be a consideration.

31.6.3 Parental Tools. Although client-based Web filtering is not common in large organizations because of the expense and management of such services on a large scale, products enabling parents to block content for their children at home have become a big business. Many large ISPs, such as AOL and MSN, offer parental content-blocking tools as a free feature of their services. Other companies sell stand-alone products that install on a home computer, with password-protected parental administrative access to content-blocking functions. Net Nanny, CYBERSitter, and CyberPatrol are some of the more popular offerings. These products typically reside on individual computers rather than on a network device, although if a home computer is set up as a hub for network connectivity (such as with Microsoft's Internet Connection Sharing), then the controls can filter traffic in the same way as an organizational proxy server. Many of these products also filter other traffic, including email, peer-to-peer file sharing, and instant messaging, as well as offering foreign language filtering, destination-address blocking, and time-of-day access rules.¹⁷

31.7 VULNERABILITIES. No security scheme, whether physical or logical, is completely free of vulnerabilities, and Web filtering is certainly no exception to this rule. Users who want to access blocked content have a variety of tactics available to them, although solutions vary in ease of use. IP spoofing, protocol tunneling, and some forms of encryption are not trivial practices, and therefore are the tools of technically adept users in reasonably small number. Other technologies, however, such as anonymizing proxies, translation sites, and caching services, are easy ways for the average user to defeat filtering. Web-filtering vendors are constantly striving to make their products more effective, while privacy and free-speech advocates support ongoing efforts to defeat what they call censorware.

31.7.1 Spoofing. In an organization that performs Web filtering on a proxy server, the organization's border firewall must allow outbound HTTP requests from the IP address of that proxy. A user who can configure traffic to look as though it is coming from the proxy's IP address might be able to get traffic through the firewall without actually going through the proxy. This tactic, known as address spoofing, takes advantage of lax routing policies on routers that forward all unknown traffic to default gateways without checking to see if that traffic came from a direction consistent with

VULNERABILITIES 31 · 11

its reported source address. The drawback of spoofing, from the attacker's point of view, is that a large organization with significant amounts of Web traffic traversing the network will notice the temporary unavailability of the proxy caused by the spoofing.

An organization can defeat spoofing by configuring internal routers to check their routing tables, to see if the address of a packet coming into an interface is consistent with the networks available via that interface. The router drops packets with source addresses inconsistent with their actual source. This tactic, called *reverse-path filtering*, requires a more capable (and expensive) router, so it is generally not available for the home user trying to set up network-based protection for parental control.

31.7.2 Tunneling. A more problematic tactic, because it relies on the behavior of applications rather than on subverting network-layer protections, is protocol or application tunneling. An application can encapsulate any other application's information as a packet of generic data and send it across the network. Virtual private network (VPN) clients use this approach to send traffic through an encrypted tunnel.

Protocol tunneling (sometimes called dynamic application tunneling)¹⁸ relies on applications that send data on commonly allowed ports. For example, a user might tunnel a Web session through the secure shell (SSH) application, which uses TCP port 22. SSH often is allowed through firewalls, and because it uses both authentication and encryption, it can be hard to monitor the difference between a legitimate SSH session and a covert tunnel. In this example, the Web browser issues a request for an HTTP page on TCP port 80, and another application running on the client system captures and redirects the port 80 request into an SSH-encrypted tunnel. The other end of the SSH tunnel could be the destination server or a proxy server somewhere between the client and server. The client application (in this case, the Web browser) is unaware of the traffic diversion and needs no altered configuration. This is similar to the approach used by traditional VPNs. A final form of protocol tunneling has come into limited use as IPv6 has become more broadly available. Still in its infancy as an Internet-wide protocol, it nevertheless is available in some pockets of the world, and tunnel brokers offer connections that tunnel IPv6 traffic over the IPv4 Internet backbone. This can become problematic if a home or site believes that it is filtering outbound requests but is only looking for content in IPv4 packets.

Application tunneling (also called static application tunneling)¹⁹ requires reconfiguration of the client application to redirect requests through a different port on the client machine. Typically, the user redefines the destination address for the application to a *localhost* address (in the 127.0.0.x range, referring to the local device), and an application then sets up a connection from that local port to the destination on an allowed network port. This approach requires alteration of the client application's configuration. Some Secure Sockets Layer (SSL) VPN products use this approach to tunnel one or more protocols, or all traffic, across an encrypted HTTP tunnel.

Tunneling via protocol or application generally requires access to either configure existing application settings or install extra software. Thus, tunneling is unavailable to users in organizations that give end users limited control over their computers. Tunneling is more of a problem for parents, whose technically adroit children can install tunneling applications to get around parental controls.

31.7.3 Encryption. Security professionals generally encourage the use of encryption, since it protects sensitive information in transit across a network. However, when users encrypt data to hide transactions that violate organizational policy, encryption can become a liability instead of an asset from the organization's perspective.

31 · 12 WEB MONITORING AND CONTENT FILTERING

In an encrypted Web session using HTTPS (which is HTTP over SSL), the contents of the session are encrypted and thus unavailable for monitoring or filtering based on URL or data content. However, the source and destination IP addresses of the session, which are visible to TCP as the session is set up, remain visible on the network. This is necessary for routers to be able to get the encrypted data packets from one end of the transaction to the other. Thus, while an HTTPS session is immune to filtering by URL text matching or content-based filtering, blocking the destination server is still effective.

As mentioned in the previous section, VPN technologies represent another use of encryption to protect the content of a transaction. Again, although the VPN encrypts the contents of the transaction, the IP addresses of the endpoints must remain visible in order to transport the data to their destination.

Another more complicated version of encryption, called *steganography*, actually embeds data inside other information to avoid detection. For instance, a user could embed a text message within the information used to encode a picture. For more details on steganography and other types of encryption, see Chapter 7 in this *Handbook*.

31.7.4 Anonymity. Most Web monitoring and content filtering relies on the identity of the user. Content filtering can happen without identity information, but unless an organization or country chooses to impose draconian broad-based filtering of all Web requests (and some do choose this approach), the organization might like to enforce filtering requirements for only some users. For instance, a school district might wish to impose stricter filtering on elementary school students than on high school students and use a more permissive policy for teachers and administrative employees.

The bane of this approach is anonymity. When privacy concerns by individuals lead to the use of anonymizing technologies on a scale that makes identity of users difficult to determine, then the only way to comply with policies and laws requiring filtering of content for some groups is to filter for all groups at the level of the strictest requirement. In the case of a school district, if the network cannot distinguish between student traffic and teacher traffic, then the district must impose the requirements of student content-filtering on teachers as well. In fact, many school districts have made this choice, due to both the extreme numbers of hard-to-find anonymizing proxies and the technical difficulty of separating student and teacher use of a common school network infrastructure. Even so, external anonymizing proxies bedevil network administrators' attempts to force compliance with filtering regulations.

Other uses of privacy-enhancing technologies (PET) include network-based anonymity schemes, such as mixing networks and onion routing. The project, known as Tor (originally an acronym for the onion router), has been gaining popularity in recent years. For more information about anonymity and identity, see Chapter 70 in this *Handbook*. For more about PET, see Chapter 42 in this *Handbook*.

31.7.5 Translation Sites. Language-translation sites, such as BabelFish (formerly at <http://babelfish.yahoo.com>),²⁰ also offer the possibility of avoiding content filters. The user enters a URL and clicks a button to request a translation of the text of the site. The user's session is between the client computer and Yahoo!, with the requested URL merely passed as keystroke data within the HTTP session, so the potentially blocked site is made available so long as the user can get to Yahoo! This is a special case of a proxy, in that the request typed into BabelFish generates a request to the server at the other end, but the client does not receive the exact results of that reply; the display, instead, includes all the graphics of the original, but the text is translated

THE FUTURE 31 · 13

into the requested language. So, from the point of view of a monitoring tool, the client is always connected with Yahoo!, and even a content filter looking at text within HTTP packets could be stymied by the foreign language. Recognizing the potential for abuse, Yahoo! published a Terms of Use document for the BabelFish site prohibiting using the service for “items subject to US embargo, hate materials (e.g. Nazi memorabilia), ...pornography, prostitution, ...[or] gambling items,”²¹ among many other classes of activities or products.

Translation sites are also annoying to system administrators because translation of a site in a particular language *into* that language usually results in unchanged content. For example, a translation of an English-language Website from French into English simply passes the content through without alteration.

31.7.6 Caching Services. One of the primary uses for proxy servers has been to reduce network traffic by saving local copies of frequently requested pages. This behavior can circumvent Web filtering, so long as the user can get to the caching server. Google, for example, caches many pages in an effort to provide very fast response to search requests. Large graphic and video files take up the most bandwidth, so Google’s Image Search feature often caches them. Often, graphic files are available even for sites that no longer exist. Users wanting to bypass blocks on sexually explicit content often use Google Image Search. Google does provide a feature called Safe Search, with three levels of voluntary filtering. The default (middle) setting filters “explicit images only,” while the strict setting filters “both explicit text and explicit images.”²² Google notes, “no filter is 100% accurate, but Safe Search should eliminate most inappropriate material.”²³ The Safe Search feature is configurable per user, and offers no password protection, so it is not a Web-filtering technology so much as it is a voluntary sex-based search filter.

31.8 THE FUTURE. As more ways of communicating and distributing content emerge, the content-filtering industry will doubtless evolve to cover the new technologies. At present, vendors sell filtering products for email, Web chat, newsgroups, instant messaging, peer-to-peer file sharing, and FTP, as well as filtering of Web requests. New features will appear in “traditional” Web filtering as well, including filtering of IPv6 tunneled traffic. The latest versions of the home filtering products Net Nanny and McAfee Parental Controls now offer the ability to force safe search options in the major search engines (like Google Safe Search, described in Section 31.7.6), and provide “object recognition,” which recognizes certain versions of Web objects (like visit counters) that are commonly used in pornography sites.²⁴

Supporters of content filtering, and those who are required to use it and need a reliable product, will be encouraged by the growth of the industry but perhaps disappointed that the problem never seems to be completely solved. Advances in image recognition could provide much better filtering, but may well spawn new ways to alter content to circumvent these tools. Those who decry these products as censorware will point out that, historically, most attempts to censor speech have failed in the end. Eventually, the two sides will probably work out an uneasy compromise, as has happened regarding sales of “adult” content in print and video. As long as some people insist on their right to distribute information that other people find offensive, this conflict is likely to continue. The debate about the balance between freedom and safety will probably become increasingly colored by concerns with illegal activity, terrorism, and cyberwar, which could lead to advances in available tools as well as attempts to define and regulate appropriate use of those tools.

31 · 14 WEB MONITORING AND CONTENT FILTERING

31.9 SUMMARY. For a variety of reasons, some better than others, groups of people with power over, or responsibility for, other groups of people want to control the kind of information to which the other groups has access and to gain knowledge of who is trading information. Monitoring and content-filtering products (or censorware) provide this kind of control using computer technology to examine information flowing across the Internet. Free speech and privacy advocates argue that content filtering prevents legitimate, legal access to information. Even should one grant the legitimacy of filtering in some cases, current technologies are prone to error, both failing to block some objectionable content and blocking some sites that contain no such content.

Most filtering techniques involve examining metadata, matching a string of text or numbers to determine the source or destination of a request or message, or some characteristic of the communication context. This metadata can be used to block access: servers can be blocked by address or by server name. The recent addition of a .xxx domain has opened the possibility of filtering entire top-level domains. Other methods focus on blocking content, examining text or nontextual parts of a Web page, including graphics. Research into image recognition and flesh-tone matching is progressing, and government agencies have used some image-recognition tools in prosecuting cases, but image recognition has not yet entered the commercial market in a large way.

With every protective, or overprotective, strategy comes a group of people dedicated to its defeat. Content filtering has a number of vulnerabilities, chief among which is the use of anonymity via privacy-enhancing technologies such as anonymizing proxies and onion routing. Other ways to defeat Web filtering include the use of protocol and application tunneling, encryption, Web translation sites, and caching services. Filtering technologies have been improving over the years, as has the inventiveness of those dedicated to thwarting them. As information outlets continue to proliferate and new communications media appear, this kind of conflict between protective technologies and privacy-enhancing circumvention is likely to continue.

31.10 FURTHER READING

- Brennan Center for Justice, “Internet Filters: a public policy report,” 2006. www.fepproject.org/policyreports/filters2.pdf
- The Censorware Project, <http://sethf.com/freespeech/censorware/essays/censorware.org.php>
- Electronic Privacy Information Center. “Censorware: A Post-CDA Solution?” http://epic.org/free_speech/censorware
- Secure Computing. “Best Practices for Monitoring and Filtering Internet Access in the Workplace,” www.securecomputing.com/pdf/scc_bestpractices.pdf (note: registration required).
- “Seth Finkelstein’s Anticensorware Investigations—Censorware Exposed,” <http://sethf.com/anticensorware>

31.11 NOTES

1. D. Chaum, “Untraceable Electronic Mail, Return Addresses, and Digital Pseudonyms,” *Communications of the ACM* 24, No. 2 (February 1981), www.cs.utexas.edu/~shmat/courses/cs395t_fall04/chaum81.pdf
2. D. M. Goldschlag, M. G. Reed, and P. F. Syverson, “Hiding Routing Information,” Workshop on Information Hiding—*Proceedings*. May 1996, Cambridge, UK.
3. www.onion-router.net
4. Codified at 47 U.S.C. §254(h) and 20 U.S.C. §9134.

NOTES 31 · 15

5. Poudre School District (Fort Collins, Colorado), Information Technology Services, accessed June 1, 2013, www.psdschools.org/department/information-technology
6. “US Intelligence Outlines Checks It Says Validate Surveillance,” *The Guardian*, June 16, 2013, www.guardian.co.uk/world/2013/jun/16/nsa-the-nsa-files
7. “NSA Collecting Phone Records of Millions of Verizon Customers Daily,” *The Guardian*, June 6, 2013, www.guardian.co.uk/world/2013/jun/06/nsa-phone-records-verizon-court-order
8. “WorldWideWebSize,” www.worldwidewebsize.com
9. Internet Corporation for Assigned Names and Numbers (ICANN), “18 March 2011 Draft Rationale for Approving Registry Agreement with ICM’s for .xxx sTLD,” March 18, 2011, www.icann.org/en/minutes/draft-icm-rationale-18mar11-en.pdf
10. “Websense: Web Filtering Effectiveness Study,” January 2006, accessed April 7, 2007, www.lionbridge.com/NR/rdonlyres/websensecontentfilte7fmfspvtsryjhojtsecqomzmiriqoefctif.pdf (URL inactive).
11. “Filtering Studies and Their Findings, 2001–2008.” Librarian in Black blog, May 7, 2010, <http://librarianinblack.net/librarianinblack/2010/05/filtering.html>
12. “Using eVe for Content Filtering,” eVision Visual Search Technology, accessed June 1, 2013, www.evisionglobal.com/business/cf.html
13. “NASA HQ Raided in Kiddie Porn Probe,” *The Smoking Gun*, March 31, 2006, www.themokinggun.com/documents/crime/nasa-hq-raided-kiddie-porn-probe
14. S. Finkelstein, “DMCA 1201 Exemption Transcript,” April 11, 2003, http://sethf.com/anticensorware/hearing_dc.php
15. SquidGuard Website, accessed June 1, 2013, www.squidguard.org
16. “DansGuardian: True Web Content Filtering for All,” accessed June 1, 2013, <http://dansguardian.org>
17. Top Ten Reviews, “Internet Filter Review 2007,” accessed April 4, 2007, <http://internet-filter-review.toptenreviews.com>
18. SSH Communications Security, “Secure Application Connectivity,” accessed June 1, 2013, www.ssh.com/manuals/client-user/53/tunnel-dynamic.html
19. SSH Communications Security, “Secure Application Connectivity.”
20. Named after the tiny Babel fish in Douglas Adams’s science-fiction classic *Hitchhiker’s Guide to the Galaxy* (New York: Del Rey, 1995). The Babel fish, when inserted in a person’s ear, would instantly enable the person to understand any spoken language. The fish name is used for a translation technology developed by AltaVista, then owned by Yahoo!, and formerly offered as a free service by Yahoo! It is also the name of a commercial translation firm that has registered the trademark BabelFish.com and which runs the translation service at www.babelfish.com (accessed June 1, 2013).
21. “Yahoo! Search Builder Terms of Use,” accessed April 7, 2007, <http://help.yahoo.com/help/us/ysearch-01.html?fr=bf-home>
22. “Google Search Settings,” accessed April 7, 2007, <http://images.google.com/preferences?q=we+live+together&um=1&hl=en>
23. “Google Help Center,” accessed April 7, 2007, <http://images.google.com/intl/en/help/customize.html#safe>
24. Internet Filter Review, “Internet Filter Terms,” accessed April 7, 2007, <http://internet-filter-review.toptenreviews.com/short-definitions.html>

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 32

VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

Justin Opatrny and Carl Ness

32.1 INTRODUCTION	32·1	32.3.2 Remote-Site VPNs	32·8
32.1.1 Borders Dissolving	32·2	32.3.3 Information-Assurance Considerations	32·9
32.1.2 Secure Remote Access	32·3		
32.1.3 VPNs	32·3		
32.2 REMOTE ACCESS VPNs	32·4	32.4 EXTRANETS	32·15
32.2.1 IPsec	32·4	32.4.1 Information-Assurance Goals	32·15
32.2.2 Transport Layer Security	32·5	32.4.2 Extranet Concepts	32·16
32.2.3 User-Authentication Methods	32·6	32.4.3 Types of Extranet Access	32·16
32.2.4 Infrastructure Requirements	32·6	32.4.4 Information Assurance Considerations	32·17
32.2.5 Network Access Requirements	32·7		
		32.5 CONCLUDING REMARKS	32·20
32.3 SITE-TO-SITE VPNs	32·7	32.6 FURTHER READING	32·20
32.3.1 Multiprotocol Label Switching	32·7	32.7 NOTES	32·20

32.1 INTRODUCTION. The rise of the Internet created a new chapter in human civilization. People are no longer tied to slowly updated information sources such as libraries. The exponential growth in the number of people looking for wide varieties of content also spurred the desire for mobility. If a person can search for information residing half way around the world from home, why not be able to do the same from the local coffee shop or while traveling during a business trip? This information revolution offered an opportunity to provide information and services to consumers, businesses, and employees at virtually any point on the globe and on a multitude of mobile platforms.

The days of focusing on protection at the network perimeter are over and have been for longer than people may realize. Laptops lead the charge of the initial mobile force. Although access to the full processing power of the laptop is useful in some instances, even the lightweight versions are bulky to travel with and have relatively limited battery

32 · 2 VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

life. Smartphones and tablets, fitting in the palm of a hand, a pocket, or a purse, now provide the power to access more commonly used applications.

32.1.1 Borders Dissolving. Widespread Internet access redefined the dynamics of network and perimeter protection. Previously, companies focused on protecting the internal network and systems exposed to the Internet. A perimeter firewall was sufficient to keep the digital predators at bay. The greater challenge is now maintaining the security of the internal network when employees use mobile technologies from home or while traveling. Further complicating the issue is how to allow other business partners—but not intruders—secure access to the systems and information you are responsible for protecting.

Organizations large and small cannot expect to maintain a competitive advantage by restricting employee network and information access to the confines of the workplace. Traveling employees can now maintain direct communications with those at the office. Employees stranded at home due to inclement weather or illness can continue to function, with direct access to necessary information. This mobility provides organizations a viable method to maintain a geographically diverse workforce that allows people to work in closer proximity to customers and partners.

Organizations are also deploying flexible mobile solutions capable of working on smart phones and tablets. Instead of bringing a corporate laptop on a customer visit, a sales associate can use a tablet to present product information, input/change order information, and access other internal resources. There is also a considerable movement to allow employees to use their personal laptop or smart device to access internal resources—and evidence that this policy may actually increase information technology costs instead of decreasing them.¹ This trend increases the potential risk profile, as the organization must provide security at different layers with little to no control over the security protections of the personal device(s) that may connect to internal resources. *Bring Your Own Device* (BYOD) is of increasing importance in discussions of productivity and of security.²

As an organization's network of vendors, suppliers, and partners grows, so does the necessity to share information. The organization must look for methods that allow these outside entities to access relevant information without exposing nonrelevant information. Information sharing is not the limit of partner involvement. There may be situations where in-house staff does not have the necessary experience to develop or support certain information systems. Although having a consultant on-staff or on-call may be cost-prohibitive, the organization may opt to allow a vendor to access specific parts of the internal network remotely, to provide the necessary levels of service. These parts of the network are called *extranets*.

Customers continue to demand increasing levels of convenience. The banking industry is a useful example of how consumer demands increase, and how business responds to meet these needs.³ Instead of having to go to a local bank, the ability to conduct certain inquiries and transactions were possible over the telephone. The limitations of this technology demanded an even greater opportunity, provided by Internet access to accounts and other banking services. Now, the consumer has a visual, point-and-click/touch interface instead of a computerized voice and numeric menus. These same convenience requirements are present in the sale of goods and services. Consumers may not want to go to a brick-and-mortar store to make purchases, but by having an online presence, an organization can meet this demand by providing an e-commerce Website that gives the consumer the ability to view and purchase items at the consumer's time and place of choosing.

INTRODUCTION 32 · 3

32.1.2 Secure Remote Access. Although the implications of not meeting partner and/or consumer demands are obvious, organizations conduct a great deal of planning and review to ensure that meeting these demands does not jeopardize the underlying safety of the information systems that drive the business. The two primary technologies leading to secure remote access are the virtual private network (VPN) and the extranet.

A VPN is a virtual network overlaying an existing set of physical and logical networks. A common VPN implementation is a secured connection allowing a remote client to access an organization's internal network through an *encrypted tunnel*. VPNs also have the ability to extend all or only selected portions of the network to other locations in the organization.

An extranet meets some of the same goals as the VPN with the emphasis being information sharing and e-commerce. Instead of connecting to the internal network, the client establishes an encrypted connection to a Web application server residing in an external service zone. Extranets can introduce additional complexity because the outward-facing systems normally require access to one or more information assets within the internal network.

32.1.3 VPNs. The essence of a VPN is simple: Create an encrypted tunnel into the internal network to protect the transmitted data. However, the concepts underlying VPN technologies are complex and require a great deal of planning to ensure the best possible implementation. VPNs normally fall into two categories: remote access VPNs and site-to-site VPNs. Each VPN type has multiple avenues of implementation, with each having unique requirements.

Allowing people to access internal network resources remotely has significant information-assurance implications. The main information assurance goals of VPNs are securely extending the internal network, protecting data during transmission, and minimizing the security impact of the process.

Smart phones and tablets often also have the ability to establish a VPN tunnel. No matter the type of mobile device, the organization must meet the requirements for remote access while accomplishing the least amount of additional risk.

When properly architected, managed, monitored, and secured, the internal network can provide a relatively safe environment for sharing data. However, the Internet is a hostile environment. The mobile device can no longer count on the internal network protections as a digital comfort zone. Internet hotspots for public access, such as those in restaurants, airports, and hotels, create even greater challenges to information assurance. These public locations (many configured for open access) create sites for malicious actors to intercept network traffic to gather information such as credentials, financial information, and intellectual property. Properly deployed and managed VPNs and extranets have the ability to thwart this type of monitoring by sending data through encrypted channels.

As the need to increase stricter security measures continues, so does the potential to interfere with the end-user experience. It does not take long for a user to figure out that disabling local security measures can have perceived advantages. Without the client firewall running, the connection speed may increase. Worse yet, by disabling security protections, the user may be able to install and run software otherwise blocked from use. The goal of the security managers must be to maintain high security that is transparent to the user. The process of establishing a VPN tunnel or accessing internal resources should seem easy to the authorized user, with the background security infrastructure

32 · 4 VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

providing the necessary levels of protection without interfering with legitimate users' productivity.

This chapter provides a foundation for understanding of the different types, terminologies, and uses for VPNs and extranets.

32.2 REMOTE ACCESS VPNs. The most common use of a VPN is for secure, client remote access. This allows an authorized user to connect to a VPN and gain access to internal resources from anywhere with an Internet connection, while maintaining the security of the transmission. Remote Access VPNs are waging an on-going battle between traditional Internet Protocol Security (IPsec) and the increasingly popular and powerful Transport Level Security/Secure Sockets Layer (TLS/SSL) VPNs. Each VPN type has distinct advantages and each warrants thorough investigation prior to selecting the technology that best fits organizational needs.

32.2.1 IPsec. IPsec is a suite of network layer protocols designed to set up and protect VPN transmissions. Traditional IPsec VPN implementations use a client-resident application or embedded operating-system service to establish a VPN tunnel into to the internal network.

32.2.1.1 Key Exchange and Management. One of the most complex aspects of IPsec is key exchange and management. IPsec uses Internet Key Exchange (IKE) to facilitate the establishment and management of a Security Association (SA).⁴ Key points of IKE Phases 1 and 2 include the following⁵:

- Phase 1 can use main mode or aggressive mode to create the initial IKE Security Association. Main mode—consisting of three pairs of packets—is the most common implementation.
- The first pair negotiates the four-parameter protection suite:
 - Encryption algorithm (e.g., Advanced Encryption Standard, AES),
 - Integrity protection algorithm (e.g., Hashed Message Authentication Mode [HMAC] with Secure Hash Algorithm, SHA-256),⁶
 - Authentication method (e.g., pre-shared key or PKI certificate), and
 - Diffie-Hellman group⁷
- The second pair exchanges encryption keys using Diffie-Hellman.
- The third pair authenticates each side of the connection to the other.
- Aggressive mode accomplishes the same task by only using three packets:
 - The first two messages negotiate the IKE SA parameters and perform a key exchange.
 - The second and third messages authenticate the endpoints to each other.
- Internet Key Exchange Version 2 (IKEv2), defined in RFC 5282,⁸ is an updated version, though not interoperable with version 1. This new version attempts to correct some of the shortcomings of the original implementation of IKE.
- Some of the improvements of IKE version 2 include:
 - Substantially fewer RFCs,
 - A reduced number of message exchanges,
 - Fewer cryptographic methods,

REMOTE ACCESS VPNS 32 · 5

- Built-in denial-of-service (DoS) protections,
- State detection and management, and probably most importantly,
- Network Address Translation (NAT)⁹ traversal via User Datagram Protocol (UDP).¹⁰

If an organization largely depends on IPsec, one or more of these improvements may make IKEv2 extremely beneficial and create a much more resilient, higher functioning of VPN service.

- Phase 2 uses quick mode to establish the IPsec source addresses (SAs).
- Each side of the connection will maintain an IPsec SA in its Security Association Database (SAD).
- The initiating device creates and sends its SA proposal to the VPN device.
- The VPN device replies with its SA selection and another hash to authenticate the connection.
- The initiating device then replies with the hash it generates from the previous request.
- If the hash matches the challenge from VPN device, the SA goes into the SAD and the connection proceeds.

32.2.1.2 Authentication Header versus Encapsulating Security Payload.

IPsec provides two security protocols for protecting encapsulated data. Authentication Header (AH) protects the integrity of the packet header and payload by using cryptographic hashing to ensure data does not change. Encapsulating Security Payload (ESP) is the more common implementation because it not only provides the integrity protection of AH but also protects the confidentiality by encrypting the entire original packet and creating a new IP header.

AH and ESP can transmit data in either *transport* or *tunnel* modes:

- Transport mode preserves the original IP header information while providing payload integrity and confidentiality payload protection.
- Tunnel mode provides integrity and confidentiality protection for the IP header and payload.

Since transport mode uses the original IP header information, this mode creates incompatibilities with Network Address Translation (NAT) because of layer-four integrity checks. NAT causes the IP address of the packet to change during transit, and in turn, will cause the integrity check to fail. Thus, tunnel mode is the primary method for host-to-gateway and gateway-to-gateway VPN connections.

32.2.2 Transport Layer Security.

The Transport Layer Security (TLS) protocol provides a method for protecting client/server communications. One of the most identifiable implementations of TLS is *Secure Sockets Layer* (SSL), which provides the basis for the HTTPS protocol. Although this chapter uses TLS and SSL interchangeably as many products support both, there are some differences. TLS is the next evolution of the Netscape-developed SSL because it is an open-standards protocol supported by the Internet Engineering Task Force (IETF) and provides enhanced security capabilities.¹¹

32 · 6 VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

The implementation of TLS/SSL is much less complicated (though not necessarily more secure) than IPsec and provides up to 256-bit encryption capabilities that already exist on virtually all Internet-connected systems. The client provides SSL-related parameters to negotiate an HTTPS connection with the server. The server then replies with the negotiated SSL-related parameters as well as its digital certificate. The client then uses the certificate to authenticate the server. If authentication is successful, the client and server establish the encryption keys used to protect the session and begin encrypted communications.

Although very specific conditions must be present, it is possible to conduct a man-in-the-middle (MITM) attack during the described the negotiation process, rendering the encryption ineffective. See Section 32.3.3.1.7 for more information.

32.2.3 User-Authentication Methods. Before being allowed unfettered access to internal resources, the client must identify and authenticate to the VPN termination device. The simplest authentication method is the password associated with the username as identifier. Some typical methods for validating user credentials are through Remote Authentication Dial-In User Service (RADIUS),¹² Lightweight Directory Access Protocol (LDAP),¹³ native Microsoft Active Directory (AD),¹⁴ or Kerberos.¹⁵ To combat the ease of stealing or guessing user names and passwords, most vendors provide alternative or complementary authentication mechanisms. An existing Public Key Infrastructure (PKI)—see Chapters 7 and 37 of this *Handbook*—opens the possibility for using dual-factor authentication by requiring a user or machine PKI certificate. Another multifactor option is to use a cryptographic token such as an RSA SecurID¹⁶ or smartcard.

Access to internal resources should not be wide-open to remote VPN users, nor should a one-size-fits-all mechanism be deployed for VPN users post-authentication. Often applied to remote access VPNs, user authorization based on the extended directory attributes of the authenticated user can further control the level of access to internal assets. For example, based on the user's group membership or role in an LDAP directory, VPN connectivity should be granted only to specific resources on the internal network. Further, a certain class of users or department should be assigned a defined, specific subnet of VPN address space.

Further protections such as downstream firewalls or host-based firewalls may control the remote users' access to resources that should or should not be accessed from off-site by policy. Access control lists (ACLs) applied to the VPN endpoint itself can limit what resources on the internal network may be accessed by remote users as well. VPN users should never be considered of equal security levels; all possible points of control, especially in a role-based access control environment, should be applied.

32.2.4 Infrastructure Requirements. VPN networks are an opportunity to use the External Service Zone. Since the VPN connection device must be Internet facing, zoning requires two different networks connected to the firewall. The external, Internet-facing (VPN untrusted) interface would contain inbound and outbound encrypted traffic only. The second network (VPN trusted) is for unencrypted traffic moving to and from the internal network. By using this principle, the firewall policy would restrict the external, inbound connections to the VPN connection device using only the few required protocols. The firewall policy would restrict all external connections to the unencrypted network while allowing authenticated VPN clients the appropriate level of access into the internal network. This unencrypted network also provides a security inspection point that is unimpeded by encryption.

SITE-TO-SITE VPNs 32 · 7

32.2.5 Network Access Requirements. IPsec and TLS/SSL remote access VPNs provide avenues for secure remote access, but each has unique challenges and opportunities.

32.2.5.1 IPsec. IPsec VPNs typically provide full access to the internal network upon successful connection and authentication. IPsec client VPN implementations using IP protocol 50 (ESP) or IP protocol 51 (AH) do not behave well behind an NAT device, which often limits usefulness for users in hotels, coffee shops, and the like. However, configuring the IPsec to use NAT transversal via UDP or TCP connection may overcome the original limitations.

It is also possible to restrict access to internal hosts and networks through a RADIUS policy, VPN connection device, and/or a firewall policy. One consideration that affects client traffic is *Split Tunneling*.¹⁷ With Split Tunneling enabled, only traffic destined for the internal network flows through the encrypted tunnel. Network traffic bound for the Internet takes the most direct route instead of traveling through the tunnel. This can help to reduce bandwidth by not having nonessential traffic flowing through the tunnel. However, the main disadvantage is losing the ability to inspect the Internet-bound network traffic.

32.2.5.2 TLS/SSL. The original implementation of TLS/SSL VPN was far from its IPsec predecessor. These early VPN connections focused on network pass-throughs to specific hosts and protocols, as well as crude, difficult-to-configure portal of simple links to internal files shares and Websites. The portal concept for remote access is unique to the TLS/SSL VPN. The current portals are more robust and flexible, allowing the administrator to feed specific content and network access depending on a user's classification. The true breakthrough for TLS/SSL VPNs was its match of IPsec's ability to allow full network access, without the necessity for a full client-resident VPN application.

The popularity of these VPNs increased rapidly because they reduce overhead, typically behave well with NATed networks, and are easier to configure than IPsec. Most modern SSL VPN implementations now also support a lightweight remote-access client, some of which are built into current smartphones and tablet devices that replicate the fully managed client/server environment traditionally used in IPsec remote access VPN deployments.

32.3 SITE-TO-SITE VPNs. VPNs are not purely for client remote access. Site-to-Site (S2S) VPNs provide the ability to facilitate internal communications needs. VPN Technologies exist that can create meshed, virtual wide-area networks (WANs) in addition to providing alternatives to traditional WAN implementations.

32.3.1 Multiprotocol Label Switching. Multiprotocol Label Switching (MPLS)¹⁸ continues to extend its dominance in the WAN world. Though it is not a traditional encrypted VPN, it does provide a similar type of service. This section will not cover most MPLS intricacies because deployment must meet organizational requirements, and because each service provider has different offerings. However, MPLS does have several distinct, advantageous characteristics.

Purpose. MPLS creates a meshed, routed virtual WAN network at the service provider level. The MPLS network is then free to route packets directly from one WAN endpoint to another within the confines of its virtual network. WAN sites are able to

32 · 8 VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

communicate directly, increasing the WAN's resilience by eliminating one location from creating as a single point of failure. MPLS networks may also provide multiple levels of quality of service (QoS). QoS provides the ability to prioritize network traffic, allowing certain protocols more bandwidth during high utilization periods. It is also possible to establish IPsec tunnels through the MPLS network to increase the confidentiality of traversing data.

Requirements. As with all WAN routing technologies, the routing hardware and software must support the necessary routing and/or MPLS protocols. In addition, the service provider must already have the technology available. Moving to MPLS may require a significant re-architecture of the existing routing infrastructure.

32.3.2 Remote-Site VPNs. Another common use of an IPsec VPN is for intersite (site-to-site, S2S) communications through the Internet.

Purpose. Although geographic diversity provides the organization closer proximity to customers and materials, traditional WAN implementations may not be available or practical. Remote-site VPNs provide an avenue for bridging this gap.

Alternative WAN. It can be cost-prohibitive to run a traditional WAN connection (e.g., T1/T3, Metro Ethernet, etc.) into a leased or owned location. In addition, even a lower-bandwidth leased line may be overpriced and underpowered to provide direct WAN connectivity to a small branch office. This type of S2S VPN may provide the necessary level of connectivity without the high cost of a leased line. Since this S2S VPN will travel over an Internet connection, it may be possible to provide a higher bandwidth for the price. However, the increased bandwidth may come at the detriment of less predictable network performance such as no end-to-end QoS.

Backup. A S2S VPN over the Internet is also an alternative means for WAN backup. Although a cellular backup card can provide increased bandwidth at a lower cost point than Integrated Services Digital Network (ISDN), the cellular networks typically have higher latency and may be infeasible depending on the building location in relation to a cell tower or provider. For isolated locations, a satellite connection may be feasible but typically has even higher latency than a cellular card. These solutions may be insufficient if the site requires heavy access to systems and information across the WAN. An S2S VPN over the Internet provides a means for higher bandwidth using existing hardware. The S2S connection provides additional redundancy by providing a standby, routable link. If architected properly, the site should notice little to no outage if the primary WAN link fails.

Vendor/Partner Connections. Business needs may necessitate the interconnection of the internal network with a vendor/partner network. It is essential to conduct proper planning, configuration, maintenance, and monitoring of these connections to ensure neither side introduces a negative situation leading to compromise and/or information loss. This setup will typically sit behind a GSD to control access to only necessary protocols, networks, and systems. However, depending on the protocols necessary (e.g., Remote Desktop Protocol, or RDP), it can be difficult to track what occurs from here. The RDP system can act as a jump-off point to other portions of the network not necessarily protected to the same level as they may exist on the internal network.

Requirements. Remote site VPNs require an Internet connection as well as VPN and encryption-enabled endpoints. A typical S2S deployment could be between two routers. Other potential deployments include a router to a VPN-enabled Gateway Security Device (GSD), or from GSD to GSD.

SITE-TO-SITE VPNs 32 · 9

32.3.3 Information-Assurance Considerations. Although the ability for people and remote sites to connect to internal resources provides mobility and increases productivity, the organization must carefully consider the information-assurance caveats. Each implementation will have its own unique requirements, but addressing the main concerns presented in the following sections can provide for more informed decisions.

32.3.3.1 Remote-Access VPN Considerations. Remote-access VPNs provide the most incremental risk. All uncontrolled networks deserve treatment as insecure, thus hostile, computing environments. Since there is uncertainty as to the security of the network the end-user system will be using to connect, it is imperative that the host have sufficient self-protection mechanisms to make up for the loss of the additional defense-in-depth layers provided by the internal network. This understanding will require the security administrator(s) to focus on the necessary precautions to protect all of the organization's systems.

32.3.3.2 Fidelity of the Mobile Device. Since there is no absolute control when mobile systems leave the confines of the internal network, the fidelity of the mobile device is a concern. The risk of compromise at every level—from physical to operating system to applications and beyond—can expose the internal network to increased risk during a VPN session.

There are several ways to reduce the likelihood of a loss of fidelity. Protections may include any combination of host-based security software (e.g., firewall, antimalware, intrusion prevention), current patch levels, and secure configurations. Proper deployment and configuration of these protection measures will help to reduce the possibility of compromise.

Even if the fidelity of the device is acceptable during the initial connection, that does not ensure this status will not change at a later time during the connection. Another method for ensuring the health of the device is through Network Access Control (NAC).¹⁹ NAC provides the ability to interrogate a connecting device before and during the connection to determine such information as patch levels, status of security protections, and memory-resident malware. Unfortunately, NAC is not just a plug-and-play proposition. There are many other implications and requirements—such as determining what to evaluate and/or enforce, inline versus passive deployments and client-based versus client-less solutions—associated with effectively deploying NAC at any point in the network.

32.3.3.3 VPN Client Management. Client administration plays a crucial role in the success of a VPN solution.

IPsec requires the use of a client-side application or embedded operating system mechanism. The primary implications of this include configuring and maintaining the client. The intricacies surrounding a hard-client installation include the necessity of local administrative permissions, potential user interaction, and dealing with a corrupt VPN client application. As with all other applications, updates are necessary due to new enhancements, unsupported older versions, and vulnerabilities. In keeping with the goal of minimizing negative user experience, administrators need a way to update client software with little to no user interaction or side effects. One primary question is to determine whether the VPN vendor has a mechanism to push configuration and other client software updates on connection. Without this functionality, it becomes

32 · 10 VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

more difficult to change client-side parameters. For example, a VPN using a preshared key for VPN authentication gets hard-coded into the VPN software configuration. This makes it difficult to change it without having the end-user enter a new key, in addition to being an avenue for compromise if the key is recoverable.

TLS/SSL VPN clients do not have the same pitfalls as IPsec client management. Instead of a hard client, the TLS/SSL VPN normally uses a small Java or ActiveX-based dynamic client or a lightweight OS-specific installation. In certain VPN systems, the administrator may have the choice to leave the dynamic client resident or remove it on disconnection. If the client becomes corrupted, the user can delete the control, and it will automatically download on the next successful connection.

Having the dynamic client remain resident has two advantages. First, this may provide the ability to have a local, encrypted workspace to store working documents and other files, or to provide access to other resources. Second, this keeps the remote system from having to download the client with each connection. The TLS/SSL client also has an advantage by always checking to see if a resident client is current. If not, the VPN system will force a download of the newest client, helping to minimize client-side management.

Proper monitoring of remote access VPNs is also essential. A compromised VPN connection or credential is extremely dangerous and an increasingly more popular target for attackers. Organizations should monitor logs and connections for suspicious activity and have procedures in place for such problems. For example, if a user ID logs in from the main office in California and less than two hours later the same user ID logs in from halfway around the world, it is likely the second connection is not intended or legitimate. This is especially true of single sign-on (SSO) environments.

32.3.3.4 Protection of the VPN Device. The heart of being able to provide VPN services revolves around the existing network security infrastructure and has direct implications regarding perimeter protection, since the VPN connection devices must be Internet *and* internal facing. These devices must also follow the external service zone discussed in Chapter 26 of this *Handbook*. The device itself is another part of the overall network security profile.

It is important to remove all ancillary exposures. The primary protection method is to configure the firewall to allow access only to the necessary ports on the VPN device. All unnecessary protocols, services, and configuration options should be shut down. If a network management tool uses ICMP and/or SNMP to monitor the VPN device, firewall rules or onboard ACLs should only allow access to these protocols from a specific set of monitoring hosts. Although network reconnaissance and attacks are still possible on these ports, this method drastically reduces the attack profile.

Administrative device access should never use insecure protocols such as Telnet, FTP, or HTTP. Instead, administration should only occur using properly configured encrypted protocols such as Secure Shell (SSH) and Hypertext Transfer Protocol Secure (HTTPS). Administrative access should at a minimum be protected with a user name and strong password, or better yet, add dual factor authentication using a certificate or token combined with a PIN or password. To minimize administrative access, only specific network(s) or host(s) should be able to access the administrative console directly, and if possible, never allowed on any Internet-facing interfaces. To reduce administrative exposures further, it may be possible to use a jump host for all administrative access. However, this creates the potential for a single point of administrative failure if the jump host is not available.

32.3.3.5 Cryptographic Options. Without carefully controlling the VPN's cryptographic protocols, it is possible to choose a combination that creates additional risk of compromise or information disclosure. The strength and proper implementation of the VPN tunnel's cryptographic protection is essential to protecting the data-in-motion and must be sufficient to meet the organization's security requirements. The selection of the available cryptographic protocols for IPsec occurs during the IKE negotiation. The cipher suite determines which cryptographic protocols are available to protect VPN sessions. For example, if 3DES and SHA1 are considered insufficient protection, these and all other weaker alternatives must be removed from the negotiable IKE or cipher suite options. This ensures that incompatible clients are unable to connect and only sufficient protections are in-use during the VPN session.

32.3.3.6 Traffic Inspection. Inspection of network traffic plays a crucial role in keeping the internal network secure. Since the goal of the VPN is to provide secure data transport, this is a hindrance when attempting to evaluate inbound, encrypted traffic. Network traffic inspection can occur on the VPN connection device or on the unencrypted side leading to the internal network. Although the network traffic is one step closer to the internal network, this provides a valid inspection point in which to detect and protect against malicious traffic. Certain VPN devices and GSDs may provide the ability to view data inside the tunnel, evaluate network traffic (providing a troubleshooting point), and inspect data before passing onto the unencrypted network. However, doing all of the inspection on the VPN device can have significant performance implications.

32.3.3.7 Processing Power. Encryption and decryption of VPN tunnels requires sufficient processing power. This issue escalates as the number of clients increases. For situations where there are only a few remote-access clients, it may be possible to use a router to terminate VPN connections and provide security services. The higher the number of remote-access clients, the greater the need for dedicated VPN devices (including GSDs). It is possible to install hardware encryption accelerators to offload cryptographic calculations and management from the shared CPU to a dedicated processing device within the router or GSD, thus easing the processing demands on these units. Larger organizations may have thousands of concurrent users, requiring multiple VPN devices.

32.3.3.8 Interception. Since VPN sessions traverse uncontrolled networks, it is possible for malicious actors to intercept the encrypted communications; however, the datagram protocol underlying TCP/IP makes capturing sequential packets difficult anywhere but the source and destination of the session. If the attackers can capture the data stream, it is increasingly feasible for them to crack the encrypted communications using distributed computing environments with massively parallel processing. However, a more practicable method is an MITM attack while the VPN session is being established. Both IPsec and TLS/SSL VPNs are potentially susceptible to MITM attacks. These MITM attacks typically start with the malicious actor on the same network as the victim and conducting an Address Resolution Protocol (ARP)-poisoning attack²⁰ to have all traffic flow through the actor's system before being passed on to the rest of the network.

The main IPsec MITM vector is a preshared key (typically a group name and password) to authenticate the IPsec connection itself. This key is typically readily

32 · 12 VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

recoverable from the VPN client configuration file(s), although this approach requires the malicious actor to exploit access to the client side to determine this information. With the group name and password determined, a tool such as *Fiked* will emulate the VPN termination device, authenticate the IPsec connection, and collect the user's credentials as it attempts to authenticate to gain access to the full VPN connection.²¹ Although this step is not a compromise of the confidentiality of data-in-motion, it provides the actor the ability to connect to the internal network at will with the harvested credentials. It is possible to defend against this attack using a client or computer certificate for the IPsec authentication.

The ability to compromise the confidentiality of TLS/SSL with relative ease became reality in 2009 with the introduction of *SSLStrip*.²² SSLStrip uses the ARP-poisoning attack to intercept the clients request to use TLS/SSL. Once intercepted, the tool will establish a valid TLS/SSL sessions with the server and create a rewritten, unencrypted HTTP session with the client. This method does not cause a certificate error message that may alert the user to their being an issue with the connection. The tool also uses a padlock favicon to deceive the user into believing the connection is secure even though the URL is actually using HTTP, not HTTPS. It is possible to mitigate this attack by only allowing direct TLS/SSL connections instead of redirecting from HTTP to HTTPS. In addition, previously established TLS/SSL sessions are not susceptible to this attack.

32.3.3.9 Site-to-Site VPN Considerations

32.3.3.9.1 *Infrastructure Design.* MPLS can be a wholesale change in how the WAN functions as well as introducing troubleshooting and security implications.

Since an MPLS network can provide any-to-any, mesh functionality, it becomes more difficult to troubleshoot network issues. In a hub-and-spoke WAN topology, it is possible to deploy centrally located probes to monitor and troubleshoot WAN issues. It is unlikely that organizations using MPLS with many locations would invest the time and money to deploy these types of monitors in all locations. However, the service provider may be able to provide additional services that provide this type of troubleshooting capability for an additional fee.

MPLS networks also add to the overall network security burden by requiring an increase in the number of network security devices to provide protection at each individual location. This is largely because the sites can establish direct communications instead of going to a hub first.

MPLS places the security of the WAN infrastructure in the service provider's hands. Ensure the MPLS provider has a robust network and strict processes for keeping MPLS networks separate. One added protection against provider-side errors is to deploy Border Gateway Protocol (BGP)²³ routing security, where both sides of the connection would authenticate before becoming part of the routing table.

Remote site VPNs create additional risk by placing the Internet directly against the internal network at remote sites. If the Internet connects directly to the site's router, there are several protection options. An ACL can restrict communication to this interface only by the other side of the remote site VPN connection. Another option would be to enable additional security features on the router. However, this may require a code upgrade and increase the processing demands on the router. A more secure but expensive option would be to deploy a GSD that provides security and VPN services—creating an additional another layer of protection between the Internet and the internal network.

SITE-TO-SITE VPNs 32 · 13

Multipoint VPNs are also an important design consideration. Where it has become common, if not necessary, for highly available, high-bandwidth VPN connections, the typical central-site or hub-and-spoke VPN design may not be desirable. Many designs rely on a central point to establish a one-to-many mesh of VPN connections. However, if the central point is unavailable, the entire VPN mesh is down.

Multipoint VPNs remove the dependency on the central VPN endpoint. If the central endpoint becomes temporarily unavailable, the other endpoints are able to communicate and negotiate connections among themselves without the presence of the central (or master) node. This allows the rest of the VPN mesh to stay operational and reroute traffic to other endpoints that would otherwise be unavailable without the central endpoint.

32.3.3.9.2 Cost. MPLS adds cost on many fronts. These include the cost of new or converted circuits and additional services such as QoS or monitoring capabilities. In addition, the equipment running the network must be able to support the new MPLS and routing architecture. There also may be a significant time investment required to redesign the network routing infrastructure.

S2S VPNs using routing devices require sufficient processing power (e.g., random-access memory—RAM—and encryption modules). In certain cases, it may cost more to get to a level of code that supports encryption. The cost of S2S VPN deployments using GSDs increases as the number of security services increases, as well as the number of clients that can connect. If deploying GSDs to remote sites, there is not only a monetary cost, but also higher administrative costs for managing this new piece of the routing and security infrastructure.

32.3.3.9.3 Availability. VPNs are less of a convenience and more of a necessity. A VPN outage can cause severe interruptions for a highly mobile workforce. Even with layers, it can be difficult to eliminate all points of failure due to the associated costs. If the organization considers VPN as a necessity, high-availability solutions are a requirement.

One possible solution is to load-balance inbound connections to distribute them across multiple devices. If one device fails, the other devices will continue to service previously connected clients and will provide a place for new clients to connect. Another option is active/standby failover: If one device fails, the other will be available for new connections. A preferable method would be for the active member to replicate connection information to the standby member. If the active member fails, the standby member can continue servicing the existing VPN tunnels and can accept new connections. This whole strategy also hinges on the main VPN connection points having redundant Internet links, uninterruptable electric power, environmental controls, and back-end network infrastructure.

Cloud-based VPN is also becoming a popular alternative to provide primary or backup availability of VPN remote access.²⁴ These services eliminated the need for a dedicated VPN infrastructure to be built and maintained by organization. Instead, a network may be extended logically to a cloud provider that integrates these connections into a distributed, highly scalable VPN infrastructure hosted in the cloud. This class of service is popular when there are many geographically diverse VPN endpoints or remote users, where low latency is necessary or uptime is critical.

32.3.3.10 Implications of Elusive VPNs. VPNs and VPN-like services occur in more places than most people realize. These unseen VPNs range from legitimate

32 · 14 VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

VPN services such as GotoMyPC²⁵ to malicious VPNs such as botnet command-and-control channels and reverse tunnels. Although VPN sites such as GoToMyPC can be useful, it is up to the organization to decide if this is an acceptable means of *ad hoc* remote access. Some of the key issues to consider with this type of VPN service are that it bypasses internal security controls such as packet and content inspection, and the third party is creating and managing the VPN connection.

Malicious VPNs are a source of continuing concern. Botnet command and control channels can use encryption to evade network defenses that protect information. A reverse tunnel works by redirecting network traffic destined for a listening port to the same port on another host. For example, if a compromised FTP server had an active reverse tunnel, any user attempting to connect to that FTP server would redirect to another FTP server. Beyond the compromised host theory, malicious people already on the internal network may establish an outgoing VPN connection in an effort to hide network traffic.

The most elusive of all are applications that provide embedded VPN-like services. Peer-to-Peer (P2P) networks tunnel file sharing through proprietary and encrypted protocols. Skype²⁶ not only has the ability to create encrypted voice calls, but also has VPN-like capabilities that allow encrypted file transfers. Network defenses should be able to detect the setup of these sessions and stop them before the connection is completed.

It is up to the organization to develop internal network security defenses and policies necessary to protect against all elusive VPNs.

32.3.3.11 Impact of IPv6. Mass migration to IPv6 is still occurring at a slower pace than one may expect, even though IPv6 was developed in the late 1990s and the allocation of the last remaining available IPv4 address blocks to regional Internet registries occurred in 2011. As of the writing of this *Handbook* in early 2013, only about 12 percent of networks in global BGP tables are advertising IPv6 prefixes. Although most organizations continue to ignore the risks and ramifications of delaying enterprise IPv6 transition plans, it is likely that the protocol is already in their environment. Almost all modern operating systems have dual IPv4 and IPv6 network stacks, usually enabled by default, and typically prefer IPv6 by default. IPv6 and its associated transition technologies have specific implications for VPNs.

Without ubiquitous end-to-end IPv6 connectivity, there are several IPv6 transition technologies (such as *6to4*²⁷ and *Teredo*²⁸) that allow IPv6 capable systems at both ends to tunnel communication over legacy IPv4 networks. The 6to4 method encapsulates the IPv6 packet in the payload of IPv4 protocol type 41 packets and requires the use of relay routers along with at least one globally assigned IPv4 address.²⁹ Teredo encapsulates IPv6 packets in the payload of UDP datagrams and also uses relay routers. UDP allows for better traversal of environments using NAT.

These technologies lessen the barrier to adoption of IPv6 by enhancing reachability of IPv6-enabled services; however, these technologies are also often not well understood or managed. Most network and security administrators learn the hard way that when working as designed, a connection (VPN or local) to an IPv6-enabled resource over a transition technology may traverse the organization's border more than once—adding severe latency and passing data through several relay routers not controlled by the organization.

These transition technologies have serious security impacts on both the local network and VPN connection—misconfigurations and lack of experience being the most

EXTRANETS 32 · 15

dangerous. Without properly configuring (and testing) both endpoints, it is possible for the IPv6 traffic to take an unsecure path before (or even after) entering the VPN tunnel, putting data in the clear and therefore at risk. It is critical for network administrators and VPN users to have a valid, tested IPv6 configuration to ensure IPv6 traffic (both native and tunneled) is traversing *only* the secured VPN tunnel and not traversing any encapsulated tunnel not controlled by the organization. Unless this threat is managed, data may be vulnerable even when an organization thinks its VPN is providing added security to a remote site or worker.

These transition technologies, many of which enabled by default, also tend to set up tunnels and connections automatically. Endpoint configuration, especially on client workstations, must be centrally configured or managed to disable these automatic connections or to specify specific configurations so that unintended tunnels are not configured.

IPv6 also has native support for IPsec (both AH and ESP) through the use of extension headers³⁰ (versus proprietary IPsec clients), though its use is optional. Configuration and use of IPsec in an IPv6 environment requires the same discipline in choosing and configuring cryptographic options as IPv4. It is also possible to use IPv4 IPsec to protect IPv6 transition technology tunnels, as unencrypted tunnel sessions would otherwise be vulnerable to interception and/or manipulation without this protection. Both of these options are generally misunderstood and often unconfigured on networking and VPN devices. Organizations should only consider native IPsec for IPv6 after a thorough education and planning process.

32.4 EXTRANETS. As the necessity for convenient access to data increases, so do the challenges associated with providing the protections required to secure this type of access.

32.4.1 Information-Assurance Goals. Allowing external entities to access internal data remotely has major information-assurance implications. The main information-assurance goals of extranets are protecting shared information assets, preventing information exposure, and minimizing ancillary risks.

Protecting Shared Information Assets. Much of the current information security market focuses on protecting hosts and networks from digital threats. Although the loss of a server or portion of the network undoubtedly causes issues, there may be areas of greater concern. Organizations develop and sustain competitive advantage by manipulating data collected or created into valuable and actionable information. By mitigating the risks associated with external access, there is less chance for a breach of information security to occur.

Preventing Information Exposure. The concept of least privilege is essential when designing system access. One would not expect a financial analyst to have access to proprietary product specifications. The same concepts hold true when assigning access to external entities. Identity management provides an accountability mechanism for providing rights and monitoring actions of issued accounts. Access management provides the enforcement mechanism that keeps those accounts from viewing or changing data not explicitly necessary. The fusion of these two principles reduces the risk of inappropriate or unexpected access.

Minimize Ancillary Risks. It does not take long for Internet-facing systems to come under attack. Minimizing Internet attack vectors occurs from the Network Layer to the Application Layer. The extranet servers are not all-purpose systems. The network protection policies should restrict access to these services so that the Internet

32 · 16 VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

community only has access to the specific service(s) advertised. Extranet services require greater consideration, as network security devices still lag in protecting the protocols and customizations routinely occurring at the Application Layer.

32.4.2 Extranet Concepts. *External Service Zone.* As described in Chapter 26 of this *Handbook*, extranet networks are another opportunity to use zoning, largely because extranets require robust protections and make extensive use of resources in other zones. Zoning allows administrators to be more granular on the services allowed into the network, as well as providing another location for additional security mechanisms.

N-Tier Architecture. If an extranet system provides all aspects of a service, there exists a single point of failure and compromise. N-Tier architecture provides the ability to distribute different aspects of the service offering to additional systems. The extranet system will provide the front-end user interface; the application server facilitates the logic and processing necessary to access data; and the back-end database server stores and provides the data. Since the extranet server resides on an external service zone, the firewall policy must allow that server to connect to its next tier servers on the internal network. When possible, restrict access from the extranet server to the internal network to only necessary systems and services. This architecture helps to create a more robust infrastructure that reduces unnecessary exposures.

TLS/SSL Encryption. The most visible marker of TLS/SSL encryption is the little golden padlock favicon that appears in an Internet browser when connecting to a server using HTTPS. When a host attempts to connect to an extranet server, several steps occur before the encrypted session can begin. The connecting host provides SSL-related parameters to negotiate a HTTPS connection with the server. The server then replies with the negotiated SSL-related parameters and a digital certificate. The client then uses the certificate to authenticate the server. If authentication is successful, the client and server establish the encryption keys used to protect the session and begin encrypted communications.

32.4.3 Types of Extranet Access. The type of systems and services able to go onto the extranet continues to grow. Providing extranet services is always subject to the needs of the organization, and the following sections provide some common examples of extranet systems.

Vendor/Partner Information Sharing. Businesses survive and thrive on competitive advantage, and increasingly need to share information with strategic partners. By providing access to internal information, these partners can better understand and meet the organizational needs. Enterprise Resource Planning (ERP) systems contain the data that serves as the digital lifeblood of an organization. On the extranet, vendors and business partners could have access to work with information stored in the ERP system.

E-Commerce. Conducting business transactions digitally can create efficiencies and reduce costs. E-commerce occurs at the business and consumer level. Businesses can use Electronic Data Interchange (EDI), which allows two entities to exchange standardized, digital versions of documents such as a purchase order and payments without human interaction. The extranet provides a means to allow business partners to exchange EDI data without granting access to the internal servers that actually process and store data. Consumers access e-commerce Websites to view and purchase goods and services.

EXTRANETS 32 · 17

Employee Productivity. Providing extranet services is not restricted to business partners and consumers. It is possible to enhance employees' access to specific internal services and systems by allowing access without requiring a full VPN connection. Email is the universal business tool, and providing Webmail services can allow employees to access email from almost anywhere in the world on any device with an Internet browser. Employees may find it helpful to be able to access and change benefits information from home. This may also be an opportunity to allow access to intranet information through a content management system.

Outsourced. Enterprises not possessing the requisite internal human or technological expertise to achieve the organization's IT goals have long looked to outsourced solutions to meet their needs. Software as a service (SAAS) provides offerings such as productivity applications, collaboration, and email (e.g., Microsoft Office365) and CRM (e.g., salesforce.com). Infrastructure and platform as a service (IAAS and PAAS, respectively) provide on-demand storage and computing (Amazon Web Services—S3/EC2 or Rackspace Open Cloud). Each of these outsourced offerings continues to mature and redefine how enterprises develop, manage, and present their information.

32.4.4 Information Assurance Considerations

32.4.4.1 Technical Security. Securing an extranet is not possible at a single layer or point in the network. Each layer is interdependent and requires specific protections to ensure the security of the extranet network.

32.4.4.2 Infrastructure. Since the extranet sits behind a GSD, network layer security measures are normally the protections encountered first.

Traffic Inspection. As reiterated in multiple parts of this and other chapters, some security devices have a distinct inability to evaluate encrypted network traffic. Extranets can create additional issues as the encrypted session may travel from the client directly to the extranet server. This potentially requires traffic inspection to occur on the extranet server. Although possible, this method will increase the processing requirements of the server.

Another option is to terminate the SSL connections on an upstream device. This device would then use unencrypted protocols on the downstream side when communicating with the actual extranet server. This method adds an additional layer of complexity when troubleshooting, but it can relieve the extranet server of the encryption and security inspection processing burdens.

Internal Network Exposure. Although the compromise of an extranet server is significant, it should not provide unfettered access to internal resources. Although the firewall policy allows access to specific extranet services from the Internet, the firewall policy must also granularly restrict access from the extranet server to internal systems.

32.4.4.3 Server. Server-level protections are the true first line of defense. It is difficult for network-layer security mechanisms to protect a server that is advertising a vulnerable service. System hardening is an approach to reducing the overall risk profile of the server. Nonessential protocols and services are shutdown. The operating system (OS) may have additional parameters such as file and account permissions that further reduce the risk profile. A script can help to ensure these changes occur in a consistent manner. One well-known system-hardening script is Bastille Linux.

32 · 18 VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

Due to the complexities of operating systems and applications, vulnerabilities will occur. The simplest method of remediating vulnerabilities is to update the server with a patch. It is important to test patches for adverse effects before deployment to the production environment. This practice will help to reduce the potential for unexpected results and downtime.

Vendors can take weeks and months to develop and release a patch. In certain cases, third-party patches will become available for unpatched, actively exploited vulnerabilities. Although the ability to use these patches is an option, it is inadvisable due to the inability to verify their integrity and safety.

Even though the goal is to minimize the risk profile of the server through all of these means, it may also be necessary to provide additional server-level protections, including host-based firewall and intrusion-prevention software. These protections have the advantage of being contextually aware of operating systems and application functions, unlike those residing at the network layer.

Virtualization is already a deeply entrenched aspect of internal server infrastructures and on the extranet. Although virtualization provides economies of scale by allowing multiple extranet systems to run on a single server or cluster, there are many security implications. Without clustering, a single point of failure exists if the physical server goes down. Since the virtual switch directs most of the layer 2 traffic generated by the VMs, there is minimal native visibility and security enforcement of the network traffic generated between the VMs. VM-embedded security enforcement is continuing to become more readily available to protect at the hypervisor level. There are also ongoing concerns of escaping virtual sessions and allowing unauthorized access to another virtual session.

32.4.4.4 Application. Some common application layer vulnerabilities include buffer overflows, Structured Query Language (SQL) injection, and cross-site scripting (XSS). Buffer overflows existed decades ago but continue to be a favorite amongst malicious code developers. This condition exists because the application does not do proper bounds checking on input fields, or uses functions that fail to do the same. SQL injection is a malicious attempt to insert an SQL query into server-side requests. By failing to check for and protect against this type of attack, the back-end server may return the information the malicious query requests. XSS is a type of code injection commonly used for phishing. A malicious actor may send an unsuspecting person a link that redirects some or all content elsewhere.

Vulnerabilities at the application layer can easily defeat even the most secure network layer configurations. Applications developed by multibillion-dollar organizations can be just as vulnerable as something developed internally in a small facility. Developers must understand and follow secure coding practices to help minimize the number and severity of application vulnerabilities. See Chapters 38, 39, and 52 in this *Handbook* for discussions of programming and application security, and Chapters 30 and 31 for Web-based security.

32.4.4.5 Policies. Since extranets provide external entities access to internal information and systems, it is advisable to have a policy that governs the use of these systems. The policy may include information about requirements such as needing a confidentiality agreement before allowing a business partner access to the extranet or mandating the use of TLS/SSL encryption for all extranet communications. Policies do not provide an active protection mechanism but do establish expectations when using extranet systems and consequences of not following those requirements.

EXTRANETS 32 · 19

32.4.4.6 Access and Identity Management. Access management is crucial to ensuring that only necessary and relevant information gets to the requesting user. The permissions given to this user will dictate the data the user is able to access and should follow the concept of least privilege. Identity management ensures that only properly credentialed entities can gain access to the extranet and underlying systems.

Typically, an external user will authenticate using a username and password. Unfortunately, this is a poor way to validate an entity's identity because of the ease of stealing or guessing passwords. A more reliable method for identifying an external entity would be to require multifactor authentication using a username/password combination and issuing the user a digital certificate. The certificate would then become a part of the connection and authentication process. Although not impossible, it is vastly more difficult to spoof and/or steal a digital certificate.

Federated Identity Management (FIdM)³¹ is a common method to provide single sign-on authentication capabilities between two or more partner organizations. Protocols such as Security Assertion Markup Language (SAML)³² create a standards-based method for each side to communicate and validate authentication information. Due to the intricate details of FIdM, this chapter only mentions its potential use as an authentication mechanism for extranet services. Since each federated organization must place a certain level of trust in the other, it is important to go through a thorough planning process to identify and mitigate the risks associated with this technology.

Chapters 28 and 29 of this *Handbook* contain detailed information about identification and authentication.

32.4.4.7 Availability. Extranet systems are an important part of business operations. When extranet access is unavailable, certain transactions do not occur that can cause issues at points throughout the organization. Extranet availability is a function of the network and server infrastructures.

A single Internet connection has an implied single point of failure if that connection goes down. If the network infrastructure leading to the extranet only has a single path to follow, this becomes another single point of failure. By adding redundancy to the network infrastructure, such as a second Internet connection and secondary routing paths, there is less likelihood of downtime.

Server infrastructures also have the same potential to be a single point of failure. Instead, using a single server to provide an extranet service, it is possible to deploy a cluster of servers and load balancing. The cluster will protect the availability of the service if one or more servers go down. In addition, this provides the ability to perform server maintenance with little to no effect on availability. Load balancing helps to optimize the number of connections going to each cluster member.

32.4.4.8 Impact of IPv6. The largest impact of IPv6 on the extranet is infrastructure and application support. As noted previously, operating systems and applications must be able to accept, process, and reply to the IPv6 packets; the routing infrastructure must be able to communicate using IPv6 natively or through an IPv6 transition tunneling mechanism (e.g., Teredo or 6to4); and lastly, the security infrastructure must be able to evaluate IPv6 packets. A breakdown in any one of these will render the extranet vulnerable or inaccessible.

32 · 20 VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

32.5 CONCLUDING REMARKS. VPNs and extranets have many powerful features that can enhance an organization's ability to conduct and improve employee mobility and business functionality. Support and input must come from all levels of the organization to ensure that the secure remote access systems meet the organization's specific needs and expectations. Most important, these systems must adhere to and enhance, not diminish, the organization's overall security profile.

32.6 FURTHER READING

- Stewart, J. Michael. *Network Security, Firewalls, and VPNs*. Jones & Bartlett Learning, 2010.
- Whitman, Michael E., Herbert J. Mattord, and Andrew Green. *Guide to Firewalls and VPNs*, 3rd ed. Delmar Cengage Learning, 2011.

32.7 NOTES

1. Tom Kaneshige, "BYOD Planning and Costs: Everything You Need to Know," *CIO*, December 13, 2012, www.cio.com/article/723864/BYOD_Planning_and_Costs_Everything_You_Need_to_Know
2. Caroline Baldwin, "BYOD Increases Productivity, but IT Departments Need To Be Prepared." *ComputerWeekly*, August 2, 2012, www.computerweekly.com/news/2240160757/BYOD-increases-productivity-but-IT-departments-need-to-be-prepared
3. Nasim Z. Hosein, "Internet Banking: Understanding Consumer Adoption Rates among Community Banks," Academic and Business Research Institute Website, January 25, 2010, www.aabri.com/LV2010Manuscripts/LV10038.pdf
4. C. Kaufman, P. Hoffman, Y. Nir, and P. Eronen, "Internet Key Exchange Protocol Version 2 (IKEv2): Request for Comments 5996," Internet Engineering Task Force, September 2010, <http://tools.ietf.org/html/rfc5996>
5. S. Frankel, K. Kent, R. Lewkowski, A. Orebaugh, R. W. Ritchey, and S. R. Sharma, "Guide to IPSEC VPNs: SP 800-77," NIST Special Publications, December 2005, <http://csrc.nist.gov/publications/nistpubs/800-77/sp800-77.pdf>
6. S. Kelly and S. Frankel, "Using HMAC-SHA-256, HMAC-SHA-384, and HMAC-SHA-512 with IPsec: Request for Comments 4868," Internet Engineering Task Force Datatracker, May 2007, <http://datatracker.ietf.org/doc/rfc4868>
7. M. Friedl, N. Provos, and W. Simpson, "Diffie-Hellman Group Exchange for the Secure Shell (SSH) Transport Layer Protocol: Request for Comments 4419," Internet Engineering Task Force, March 2006, <http://tools.ietf.org/html/rfc4419>
8. D. Black and D. McGrew, "Using Authenticated Encryption Algorithms with the Encrypted Payload of the Internet Key Exchange Version 2 (IKEv2) Protocol: Request for Comments 5282," Internet Engineering Task Force, August 2008, <http://tools.ietf.org/html/rfc5282>
9. Internet Engineering Task Force, "Network Address Translators (NAT)—Description of Working Group," Internet Engineering Task Force Datatracker, November 2001, <http://datatracker.ietf.org/wg/nat/charter>

NOTES 32 · 21

10. J. Postel, "User Datagram Protocol: Request for Comments 768," Internet Engineering Task Force, August 28, 1980, <https://tools.ietf.org/html/rfc768>
11. T. Dierks and E. Rescorla, "The Transport Layer Security (TLS) Protocol v1.2—Request for Comment 5246," Internet Engineering Task Force, August 2008, <http://tools.ietf.org/html/rfc5246>
12. A. DeKok and G. Weber, "RADIUS Design Guidelines: Request for Comments 6158," Internet Engineering Task Force, March 2011, <http://tools.ietf.org/html/rfc6158>
13. K. Zeilenga, "Lightweight Directory Access Protocol (LDAP): Technical Specification Road Map (Request for Comments 4510)," Internet Engineering Task Force, June 2006, <https://tools.ietf.org/html/rfc4510>
14. Microsoft, "So What Is Active Directory? (Windows)," Windows Dev Center—Desktop, 2013, [http://msdn.microsoft.com/en-us/library/windows/desktop/aa746492\(v=vs.85\).aspx](http://msdn.microsoft.com/en-us/library/windows/desktop/aa746492(v=vs.85).aspx)
15. Jason Garman, *Kerberos: The Definitive Guide*, O'Reilly Media, 2003.
16. EMC², "RSA SecurID," EMC² Security, 2012, www.emc.com/security/rsa-securid.htm
17. Tom Shinder, "Split Tunneling versus Force Tunneling for DirectAccess Clients," Microsoft TechNet, March 30, 2010, <http://social.technet.microsoft.com/wiki/contents/articles/135324.aspx>
18. Internet Engineering Task Force, "Multiprotocol Label Switching (mpls): Charter for Working Group," Internet Engineering Task Force Datatracker, November 2012, updated July 2013, <http://datatracker.ietf.org/wg/mpls/charter>
19. J. Sabir, "Network Access Control (NAC) Vendor List," JafSec.Com Website, 2013, <http://jafsec.com/Network-Access-Control/Network-Access-Control-A-B.html>
20. Corey Nachreiner, "Anatomy of an ARP Poisoning Attack," WatchGuard Infocenter, 2011, www.watchguard.com/infocenter/editorial/135324.asp
21. Irongeek.com, "FIKED," Irongeek.com Website, August 4, 2009, www.irongeek.com/i.php?page=backtrack-r1-man-pages/fiked
22. Irongeek.com, "Using SSLStrip To Proxy an SSL Connection and Sniff It," Irongeek.com Website, 2013, www.irongeek.com/i.php?page=videos/sslstrip
23. CISCO, "Border Gateway Protocol," CISCO DocWiki, October 16, 2012, http://docwiki.cisco.com/wiki/Border_Gateway_Protocol
24. VPN Reviews, "Best VPN Reviews—Compare and Find VPN Account Providers," Best VPN Reviews Website, 2013, www.vpnreviews.com
25. GoToMyPC, "Access Your Mac or PC from Anywhere," Citrix GoToMyPC Website, 2013, www.gotomypc.com/remote_access/remote_access
26. Skype, "Wherever You Are, Wherever They Are—Skype Keeps You Together," Skype, 2013, <http://beta.skype.com/en>
27. B. Carpenter and K. Moore, "Connection of IPv6 Domains via IPv4 Clouds: Request for Comments 3056," Internet Engineering Task Force, February 2001, www.ietf.org/rfc/rfc3056.txt
28. Microsoft, "Teredo Overview," Microsoft Technet Windows, January 15, 2007, <http://technet.microsoft.com/en-us/library/bb457011.aspx>
29. Carpenter and Moore, "Connection of IPv6 Domains."

32 · 22 VIRTUAL PRIVATE NETWORKS AND SECURE REMOTE ACCESS

30. Charles M. Kozierok, “IPv6 Datagram Extension Headers,” The TCP/IP Guide Website, September 20, 2005, www.tcpipguide.com/free/t_IPv6DatagramExtensionHeaders.htm
31. EDUCAUSE, “Seven Things You Should Know about Federated Identity Management,” EDUCAUSE, September 2009, <http://net.educause.edu/ir/library/pdf/EST0903.pdf>
32. Thomas Hardjono, Nathan Klingenstein, Anil Saldhana, Hal Lockhart, and Scott Cantor, “OASIS Security Services (SAML) TC,” *OASIS: Advancing Open Standards for the Information Society*, 2013, https://www.oasis-open.org/committees/tc_home.php?wg_abbrev=security

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 33

802.11 WIRELESS LAN SECURITY

**Gary L. Tagg, CISSP and
Jason Sinchak, CISSP**

33.1	INTRODUCTION	33·2	33.3.7 Temporal Key Integrity Protocol (TKIP) 33·17	
33.1.1	Scope	33·2	33.3.8 Counter-Mode/CBC-MAC Protocol (CCMP) 33·19	
33.1.2	Corporate use of Wireless LANs	33·2	33.3.9 Fast Secure Roaming 33·20	
33.1.3	Functional Benefits of Wireless	33·3		
33.1.4	Security Benefits of Wireless	33·3		
33.1.5	Centralized Management	33·4		
33.1.6	Overview and History of the IEEE 802.11 Standards	33·4		
33.2	802.11 SECURITY FUNDAMENTALS	33·6	33.4	FUNDAMENTAL WIRELESS THREATS 33·21
33.2.1	Terminology	33·6	33.4.1 Unauthorized Network Extensions 33·21	
33.2.2	Authentication and Access Control	33·7	33.4.2 Layer 1 Attacks 33·23	
33.2.3	Data Confidentiality	33·8	33.4.3 Endpoint Attacks 33·24	
33.2.4	Key Management	33·8		
33.3	IEEE 802.11 ROBUST SECURITY NETWORK	33·9	33.5	SPECIFIC WIRELESS SECURITY ATTACKS 33·24
33.3.1	Features	33·9	33.5.1 MAC Address and IP Spoofing 33·24	
33.3.2	802.1X Overview	33·10	33.5.2 EAP (WEP/WPA/WPA-Enterprise) RADIUS Impersonation Attack 33·25	
33.3.3	EAP, EAPoL, and PEAP	33·11	33.5.3 Management Frames 33·26	
33.3.4	Insecure Legacy EAP protocols	33·12	33.5.4 Hotspot and Cafe Latte Attacks 33·27	
33.3.5	Detailed EAP/EAPoL Procedure	33·12	33.5.5 WEP 33·29	
33.3.6	RSN Key Hierarchy and Management	33·16	33.5.6 WPA/WPA2 Pre-Shared Key 33·31	
			33.5.7 WPA/WPA2: Wi-Fi Protected Setup (WPS) 33·32	

33.2 802.11 WIRELESS LAN SECURITY

33.5.8	MS-CHAP-v2— Divide and Conquer	33.33	33.7.1	Benefits of Wireless Controller Architecture	33.41
33.6	MITIGATING CONTROLS	33.34	33.7.2	Network Segmentation	33.42
33.6.1	Improvements	33.34	33.7.3	User Segmentation	33.44
33.6.2	Authentication Server and Client Certificates	33.35	33.8	SECURITY AUDITING TOOLS	33.44
33.6.3	Endpoint Supplicant Configurations	33.36	33.9	CONCLUDING REMARKS	33.46
33.6.4	Wireless Intrusion- Detection Systems	33.37	33.10	ABBREVIATIONS AND DEFINITIONS	33.46
33.7	SECURE ENTERPRISE DESIGN	33.40	33.11	FURTHER READING	33.50
			33.12	NOTES	33.50

33.1 INTRODUCTION. IEEE 802.11 wireless local area networks (LANs) are now ubiquitous and have major benefits such as mobility, flexibility, rapid deployment, and cost reduction over traditional wired networks. However, as with any networking technology, wireless LANs create opportunities for unauthorized individuals to potentially access the enterprise network and the information carried over it.

This chapter provides an overview of wireless LAN technologies, security threats and attacks, and how to address them. It is structured as follows:

- 802.11 history and technological overview
- 802.11 security fundamentals
- Detailed coverage of 802.11 security covering both the original legacy functionality and the upgraded system first defined in 802.11i-2004
- Fundamental wireless-medium security threats
- Specific technical wireless-LAN security attacks
- Technical mitigating controls for specific security attacks
- Overarching secure enterprise design principals

33.1.1 Scope. The scope of this chapter is the security of ANSI/IEEE standard 802.11 wireless LANs. This chapter does not consider any other wireless systems, such as mobile telephone networks, or other wireless standards such as HomeRF, Bluetooth, WiMax, or HiperLAN. This chapter provides a high-level overview of fundamentals with technical deep dives in specific areas necessary to comprehend threats. Many of the basic terms and concepts are documented in Chapter 32 in this *Handbook*.

33.1.2 Corporate use of Wireless LANs. Corporations have been using wireless LANs since the 1990s. However, to begin with, the market was fairly small and the technologies proprietary. In the late 1990s and early 2000s, the groundwork was laid for the mass adoption of wireless LANs. The starting point was the publication of ANSI/IEEE standard 802.11 that provided a baseline design enabling manufacturers to develop interoperable products at lower costs.¹

INTRODUCTION 33 · 3

33.1.3 Functional Benefits of Wireless. The main advantages of implementing wireless networks are mobility, flexibility, and cost reduction.

- **Mobility:** Wireless technologies enable staff to access network information via mobile terminals as they move around the office campus; examples include warehouses, shop floors, and hospitals. Within an office environment, wireless technologies provide a flexible alternative or addition to the wired network. Often, desks and meeting rooms have a limited number of Ethernet connections; wireless technologies can cost effectively provide additional network connections as required.
- **Flexibility:** Public wireless networks (hotspots) allow staff to utilize idle time between meetings, in airports, coffee shops, and even on airplanes in flight.² Typical uses include access to the corporate LAN, along with information on the Internet. Public hotspots can also be trunked over the enterprise wireless LAN to provide internet access for consultants and visitors.
- **Cost Reductions:** Costs can be lowered by not having to install physical network links between buildings separated by a road, river, railway tracks, or even a city block. A wireless link can be set up between two buildings provided there is an uninterrupted line of sight between them. In addition, economies of scale can be realized with a wireless medium. A single access point (AP) (thick or thin) can service one or many end users and scale appropriately by bumping excessive users to neighboring APs. This capability is exclusive to wireless. A wired network has a fixed capacity for per-port access and virtual LAN (VLAN) assignment. The use of wireless with Service Set Identifier (SSID) VLAN support can reduce the networking hardware volume. In the case of dynamic VLAN assignment within a wired LAN, endpoint or server moves require proper VLAN pruning at the switch layer to retain security yet provide the required VLANs. Cost savings can be realized within a wireless LAN by reducing network-management tasks and overall complexity to maximize resources and hardware usage.

33.1.4 Security Benefits of Wireless

- **Physical Security:** An AP and supporting hardware can be hidden from end users to protect it from physical attack. A wireless AP can be hidden above the ceiling in contrast to physical endpoint network jacks, which must be accessible by all users who need access to the internal network. This leaves the door open for unauthorized access to a network port in the absence of enterprise wide 802.1X port security.
- **Segmentation Visibility:** Wired networks commonly assign VLANs on a per-port basis, or in advanced configurations, using one per Media Access Control (MAC) address. In port-based VLAN assignment, management relies heavily on proper access-layer switch configuration and the pruning of sensitive VLANs not necessary for a specific business unit or location. MAC address-based VLAN assignment requires the use of MAC authentication-aware switch hardware and a backend authentication and RADIUS server for MAC → VLAN mappings. MAC-based VLAN assignments can improve management capability and oversight, but pose a significant security risk as they are susceptible to MAC-address spoofing.

In a wireless environment, VLANs can be assigned on a per-SSID basis and administrative efforts can be greatly reduced. A user is no longer confined to a particular

33 · 4 802.11 WIRELESS LAN SECURITY

physical location for access to a necessary VLAN, as long as the SSID is available on accessible APs. In the event of expansion, deploying a new AP in close proximate will deliver the necessary VLANs with the corresponding SSID that a client supplicant profile is configured to search for. Since VLAN assignment can be controlled through minor client supplicant configurations and backend RADIUS privileges, segmentation can be determined during the asset-provisioning process and require relatively minor management thereafter.

In an environment which relies heavily on access-layer switches for endpoint connectivity, it is not uncommon to find switching or upstream router protocols making their way to access-layer switch ports. These include protocols such as Dynamic Trunking Protocol (DTP), VLAN Trunking Protocol (VTP), Hot Standby Router Protocol (HSRP), Cisco Discovery Protocol (CDP), Open Shortest Path First (OSPF) routing, and so on. A user with port access can exploit multiple known vulnerabilities in these protocols.

33.1.5 Centralized Management. Thin-client AP environments leverage wireless controllers to provide a central configuration for all associated APs. Thin-client APs are configured on the controller, and intelligence is provided by a backend user directory (Extensible Authentication Protocol–Remote Authentication Dial In User Service, or EAP-RADIUS) and controller itself. The widespread yet central connection of a mesh of APs can additionally be leveraged for security monitoring of the wireless airwaves in a Wireless IDS fashion as described in Section 33.4.

33.1.6 Overview and History of the IEEE 802.11 Standards. The first IEEE 802.11 standard was published in 1999. Following publication, work continued on developing 802.11 with amendments published on a regular basis. The publication of the 802.11b standard increased WLAN bandwidth from 2Mb/s to 11Mb/s, making it a possible technical replacement for a wired LAN. Following that were the 802.11a and 802.11g standards that increased throughput further to 54Mb/s.

These and other amendments were then brought together in the revised standard 802.11-2007. Development continued via further amendments, which were once again brought together with the issue of 802.11-2012.³

The first version of the standard contained authentication and confidentiality services in an attempt to provide similar levels of security as wired LANs. The authentication service consisted of two systems called *Open Authentication* and *Shared Key Authentication*, and the data confidentiality service was called *Wired Equivalent Privacy* (WEP). These services are often referred to as the *legacy security services*.

In 2004 the 802.11i standard was released and defined the *Robust Security Network* (RSN) system, which provided enhanced authentication and confidentiality services. The authentication service has two options: the first is the use of pre-shared keys targeted for homes and Small Office Home Office (SOHO) users and the second uses the 802.1X/EAP framework for enterprise use.

Following the issue of the 802.11-2007, amendments were introduced to provide security mechanisms for mesh wireless networks and the use of 802.11 for transportation applications. There were also enhancements to the core 802.11i algorithms to secure some management frames, and enable fast transition to other APs for time critical services such as Voice over IP (VoIP). Throughput was further improved to a maximum of 600Mb/s by the release of 802.11n-2009.

Exhibit 33.1 provides an overview of the development of the 802.11 standard, the amendments, and identifies those with security functionality using shading.

INTRODUCTION 33 · 5**EXHIBIT 33.1 802.11 Standards**

Standard	Description	Security Functionality Description
802.11-1999	The original ANSI/IEEE standard 802.11, 1999 edition. Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications.	Defined open and shared key authentication and WEP
802.11a-1999	An extension to the 802.11 standard that provides an improved throughput of up to 54Mb/S in the 5GHz frequency band.	None
802.11b-1999	The ANSI/IEEE 802.11b specification is an extension to the 802.11 standard that specifies a higher speed physical layer (up to 11Mb/s).	None
802.11d-2001	802.11d extends the original 802.11 standard to define physical layers, for adoption in additional countries.	None
802.11e-2005	802.11e extends the original 802.11 MAC to provide quality of service support for time-critical applications.	None
802.11g-2003	This standard provides an improved throughput of up to 54Mb/S in the 2.4GHz band.	None
802.11h-2003	802.11h provided spectrum and transmit-power enhancements to the 802.11a 5GHz standard to allow its use in Europe.	None
802.11i-2004	802.11i enhanced the 802.11 MAC to provide improvements in security. This standard defined the RSN security system. The two main encryption algorithms are also known as WPA and WPA2.	Defined the RSN security system
802.11j-2004	802.11j enhanced the 802.11 standard, to add channel selection for 4.9 GHz and 5 GHz in Japan.	None
802.11-2007	The ANSI/IEEE standard 802.11-2007 rolled up all the above amendments into a revised standard.	Security functionality defined in 802.11i-2004
802.11k-2008	802.11k enhanced 802.11-2007 to improve the management and performance of wireless networks. Access Points collect information from clients to control channel selection, roaming, and transmit power control.	None
802.11n-2009	802.11n provides an improved throughput of up to 600Mb/S in both the 2.4 and 5GHz bands.	None
802.11p-2010	802.11p enhanced the 802.11 standard for intelligent transportation system applications.	Dedicated security functionality defined in the standard

(continued)

33 · 6 802.11 WIRELESS LAN SECURITY

EXHIBIT 33.1 *(Continued)*

Standard	Description	Security Functionality Description
802.11r-2008	802.11r provided fast transitions between APs and to establish the required Quality of Service before the transition.	Redefined the security key negotiation protocol to avoid the lengthy 802.1X process
802.11s-2011	802.11s defines a self-configuring multihop wireless distribution system to provide Mesh networking.	Defines a password-based authentication and key establishment protocol known as SAE
802.11u-2011	802.11u provided improvements to the user experience with wireless hotspots, addressing issues such as enrollment, network selection, and authorization for roaming.	None
802.11v-2011	802.11v extends the management of attached wireless clients established in 802.11k to include the ability to configure the clients.	None
802.11w-2009	802.11w introduced protection for some wireless management and control frames to address potential denial of service deauthentication spoofing attacks through the use of a Message Integrity Code (MIC).	Extended 802.11i to also protect some management frames
802.11y-2008	802.11y produced an enhancement to the existing 802.11 standard to enable high-powered wireless LAN operations in the 3650-3700 MHz band for the USA.	None
802.11z-2010	802.11z defines Tunneled Direct Link setup that enables client devices to establish a direct link with each other whilst remaining associated to an AP. This can more than double throughput between the clients.	None
802.11-2012	The ANSI/IEEE standard 802.11-2012 rolled up all the above amendments to 802.11-2007 into a revised standard.	RSNA defined in 802.11i and mechanisms designed in the above amendments

33.2 802.11 SECURITY FUNDAMENTALS. This section describes the underlying 802.11 security fundamentals in preparation for the details covered in later sections.

33.2.1 Terminology. To better understanding this chapter, readers will find clarification of a few commonly misunderstood 802.11 terms. Section 0 has a more extensive glossary.⁴

Authentication: Authentication is the first step in a two-step process of client connection to an AP. This step validates the client's authority to start associating

802.11 SECURITY FUNDAMENTALS 33 · 7

with the AP. Network-based authentication such as username and password takes place after layer 2 authentication and association.

Association: Association is the process of the AP accepting the client connection and allocating resources for it. This includes things such as adding client specific information such as supported data rate, data protocol 802.11 b/g/n, and MAC address information.

802.1X: 802.1X provides the encapsulation of Extensible Authentication Protocol (EAP) implementations within 802 communication mediums. 802.1X does not define a specific authentication method, but provides a vehicle for EAP implementations and their underlying methods. 802.1X and ultimately 802.11i are major components of the modern 802.11 RSN system. 802.1X allows or denies client access to requested resources until the client is successfully authenticated.

RADIUS: *Remote Authentication Dial In User Service* (RADIUS) is a network protocol used for authentication, authorization, and accounting (AAA). Historically, RADIUS servers leveraged a flat-file directory for user-based access decisions, but modern implementations leverage a dedicated directory such as Windows Active Directory, Lightweight Directory Access Protocol (LDAP), or a relational database such as Microsoft SQL Server or Oracle. In a Robust Security Network (RSN), a RADIUS server is only responsible for brokering the authentication/authorization of a user requesting access to an AP. Authentication data is transparently sent from the client to the RADIUS server by the AP, and the RADIUS server leverages an external directory to determine a response.

SSID versus BSSID: An SSID is a unique identifier used by a client to establish connectivity to a particular wireless network. An AP can provide multiple SSIDs on the same channel through the use of the same or multiple interfaces. A BSSID is the unique identifier for a Basic Service Set (BSS). A BSS consists of an AP and associated clients or clients.

An **Extended Service Set** (ESS) is a series of BSSIDs (AP interfaces) sharing the same SSID. This helps enable a wireless client to seamlessly move between APs using the same SSID. A BSSID is a separate interface with its own MAC address; multiple SSIDs can share the same interface and MAC address. On a commercial AP the first SSID/VLAN pair will use the BSSID interface/MAC address and each SSID after that will use a virtual MAC address, which increments the BSSID by a small value for each SSID. A BSSID can be used to reference a unique interface or AP (assuming the AP only has one interface). Depending on the vendor of a particular AP, a single BSSID will send broadcast beacons with all SSIDs in one sweep. If an SSID is placed on its own BSSID, it will have a dedicated beacon, which may improve compatibility.

33.2.2 Authentication and Access Control. In the absence of 802.1X port security, wired-LAN access control is primarily reliant on physical security. To access the LAN, an attacker first needs to have physical access to a connection point. In contrast, the nature of wireless LANs means any wireless client within radio range can potentially connect to the internal LAN network, bypassing physical controls.

33 · 8 802.11 WIRELESS LAN SECURITY

To address this issue 802.11 provides an optional native-authentication service. The legacy 802.11 standard includes two protocols, open authentication and shared key authentication. Neither of these protocols was adequate for secure access.

Open authentication is a null-authentication service allowing any and all clients to connect and associate. Shared key authentication (SKA) requires the client to use a cryptographic key to successfully authenticate. This method did put a lock of sorts on the network, but it was soon exploited by attackers who acquired keying material to help compromise the system.⁵ SKA lacked unique authenticated user tracking, and forced an out-of-band key management process. SKA has since been deprecated and should not be used, except for necessary backwards compatibility with older devices. The legacy 802.11 security system did not provide any access-control functionality. Although most equipment included access filtering based upon MAC address, this control was not part of the standard and the legacy service is therefore easily defeated.

The RSN defined by 802.11i is a much stronger system, and provides multiple mechanisms to authenticate the device and the user.⁶ The system defines a personal profile for home/SOHO use based on PSKs, as well as an enterprise profile based upon the 802.1X/EAP framework, which allows the use of a backend authentication server. The 802.1X authentication system also allows for access control decisions to be taken by the network to restrict network resources available to an authenticated user through technologies such as VLAN segmentation.

33.2.3 Data Confidentiality. Data confidentiality on a wired network is provided by physical security and layer-2 boundaries, which limit the accessibility of the data. Unless an attacker is physically connected to the wired LAN and logically resides between a sender and recipient on the network, through ARP poisoning or physical position, the data traveling over the network cannot be captured. A wireless LAN uses a physically public medium for data transfer; therefore, every packet traveling between a wireless client and an AP is transmitted via radio signals and can be captured by any client within radio range. Although captured data may be encrypted and not readily viewable, it should be noted that any client can obtain the ciphertext in some manner despite layer-2 or logical boundaries.

To address this issue the 802.11 standard provides a data confidentiality service. The legacy system provided the Wired Equivalent Privacy (WEP) protocol, which encrypted each message with a symmetric key before transmission. However, this protocol was successfully attacked and automated tools built to enable WEP keys to be cracked by anyone with basic IT skills and access to the freely available toolset.⁷ Overtime, improvements have been made to the WEP protocol, moving it toward a more respectable enterprise solution, but these fell short and were soon superseded by newer, ground-up protocols. Additional controls are needed to protect the network if WEP must be used, such as with older existing equipment.

The RSN system provides two new data confidential protocols called *Temporal Key Integrity Protocol* (TKIP)⁸ and *Counter Mode with Cipher-Block Chaining Message Authentication Code Protocol* (CCMP).⁹ As well as confidentiality, both protocols provide message integrity as well.

33.2.4 Key Management. Secret-key encryption itself is a relatively simple process designed to protect data long enough for the encryption keys to be changed at a determined interval. Designers start by choosing a suitable proven algorithm, protocol, and key length that they are confident will protect user data for a defined

IEEE 802.11 ROBUST SECURITY NETWORK 33 · 9

period. The user (or user program) then provides that algorithm with the plaintext data and an encryption key, and the data are then encrypted and ready for transmission. Due to the rapidly advancing speed of processing (of individual CPUs and CPUs running in parallel), no level of any encryption can provide definitive protection for an unlimited period of time—the time required for brute-force testing of all possible keys in the keyspace—necessitating the need for a key scheduling and management system.

Securely and repeatedly establishing mutual encryption keys is the single most complex issue plaguing cryptographic communications between two parties in physically separate locations. The repeatable process of creating keys, mutually (and securely) exchanging different keys or independently deriving the same key, and finally destroying them, is key management.

The legacy 802.11 algorithms did not provide any specific key-management functionality and vendors were left to design their own. The most common system in early standalone APs was to manually enter a defined-length static WEP key into each client and AP. For home users, the RSN system defines the manual entry of a common variable length PSK or passphrase in the client and AP. From this PSK, a key-management process derives working cryptographic keys, which are changed for each message.

For enterprises, RSN uses the 802.1/EAP framework to establish a secure channel during the user- and device-authentication phase, allowing a pairwise master key (PMK) to be set up between the client and the AP. From this PMK, a key-management system establishes working cryptographic keys which are changed for each message.

33.3 IEEE 802.11 ROBUST SECURITY NETWORK. In June 2004, the IEEE released the 802.11i standard to improve the security of 802.11 networks. This new system is called the *Robust Security Network* (RSN) and is designed for both personal and enterprise users. Enterprise use is based on the 802.1X protocol to provide authentication and establish a security context. A “personal” profile uses a *pre-shared key* (PSK) based on a password provided for consumers and SOHO users who do not require the necessary 802.1X backend authentication infrastructure.

33.3.1 Features. The core protocol of RSN is IEEE 802.1X, which forms *RSN Associations* (RSNA) with the wireless network. RSN provides the following features:

- Mutual authentication mechanisms. These mechanisms can authenticate users as well as the network client or machine. The AP and backend authentication server can also be authenticated to the client defeating rogue AP and man-in-the-middle attacks.
- Key-management algorithms
- Cryptographic key establishment (through PMK and Pairwise Transient Key [PTK] establishment)
- Cryptographic message integrity codes to defeat the bit-flipping attacks possible in the original standard (WEP)
- Two data privacy protocols which also implement message integrity:
 1. *Temporal Key Integrity Protocol* (TKIP), which is an optional protocol specifically designed so that existing WEP based hardware can be upgraded to use it.

33 · 10 802.11 WIRELESS LAN SECURITY

2. *Counter Mode with CBC-MAC Protocol* (CCMP), which is mandatory for RSNA compliance. It uses the Advanced Encryption Standard (AES) in Counter mode for confidentiality and CBC-MAC for authentication and integrity. CCMP is a strong protocol that has been designed for the next generation of wireless equipment.

TKIP was designed as a temporary stopgap that would work on existing WEP-based hardware until new hardware containing the CCMP protocol became commonplace. Whereas TKIP continued to use the existing RC4 encryption algorithm that was at the heart of WEP, CCMP uses the AES algorithm and requires more powerful hardware—no longer a problem at present.

33.3.2 802.1X Overview. 802.1X was originally designed for port-based network-access controls for IEEE 802 LAN infrastructures. These infrastructures include Ethernet, token ring, and wireless networks. 802.1X authenticates and authorizes devices attached to an LAN port, and will not allow a device to access the network if authentication fails.

802.1X defines three roles:

1. **Authenticator.** The device that authenticates a network device before allowing it to access network resources. In an 802.11 BSS network the AP is the authenticator.
2. **Supplicant.** The device that wants to access network resources and needs to be authenticated.
3. **Authentication Server (AS).** The AS performs the actual authentication of the supplicant on behalf of the authenticator. The AS can be located with the authenticator, but is commonly an external system such as a RADIUS server.

The 802.1X standard defines the object *Port Access Entity* (PAE) which operates the authentication algorithms and protocols in the supplicant and authenticator. An overview of the 802.1X architecture is shown in Exhibit 33.2 below.¹⁰

The authenticator has two logical ports; the first is an uncontrolled port that allows access to required functionality such as the authenticator PAE. The second port is the controlled port that allows access to the rest of the network. The status of the controlled port is set by the authenticator PAE and is dependent upon the outcome of the authentication between the supplicant and the authentication server.

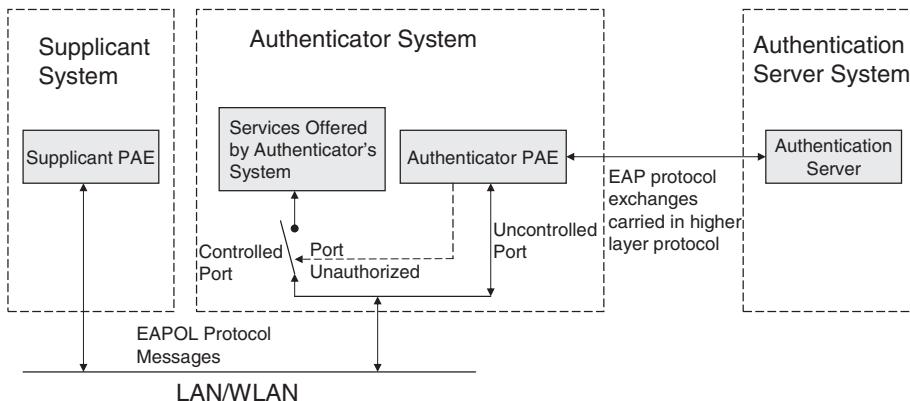


EXHIBIT 33.2 802.1X Architecture

IEEE 802.11 ROBUST SECURITY NETWORK 33 · 11

The messages between the supplicant and authenticator use the Extensible Authentication Protocol (EAP) over LAN (EAPoL) framework defined in 802.1X. Communications between the authenticator and the AS leverage the EAP framework carried in a higher-layer protocol such as RADIUS.

33.3.3 EAP, EAPoL, and PEAP. 802.1X relies on EAP¹¹ to perform authentication of the supplicant and authenticator. The EAP protocol is a series of method interfaces that make up the framework known as EAP. EAP was originally designed for use on modem dial-up networks; therefore, the 802.1X specification details the expansion of EAP across Ethernet/Token Ring networks through the EAPoL (EAP Over LAN)¹² extension. In addition to placing EAP methods within an Ethernet payload, EAPoL specifies a number of additional functions to assist in the authentication process during discovery and key exchange. There are currently over 40+ different implementations of the EAP framework. The primary difference in EAP implementations centers on how the supplicant and authenticator are authenticated.

Standard EAP is not a mutual authentication protocol; only the supplicant is authenticated. This makes supplicants vulnerable to rogue AP attacks. Additionally, due to its original design for physical dial-up connections, EAP does not protect its authentication messages from eavesdropping. Therefore, the current EAP/EAPoL implementations establish secure tunnels to provide security prior to exchanging authenticating material. The following are the most common EAP/EAPoL enterprise 802.1X implementations ranked by the level of afforded security:

1. **EAP-TLS**¹³: This EAP implementation only allows mutual certificated-based authentication through *Transport Layer Security* (TLS) X509 certificates. Authentication of both the backend authenticating server and wireless client provides a strong level of wireless security, but requires a full enterprise *Public Key Infrastructure* (PKI) implementation to securely distribute and update client keys on a regular basis. An issue with EAP-TLS is most organizations do not have the necessary PKI to issue the supplicant client TLS certificates.
 - **EAP-TTLS (*Tunneled TLS*)**¹⁴: EAP-TTLS is similar to EAP-TLS and supports mutual certificate authentication but does not require client-side certificates. EAP-TTLS creates a TLS tunnel prior to starting any network authentication process and can therefore tunnel any password authentication mechanism, even insecure legacy mechanisms such as PAP.
2. **EAP-PEAP (*Protected EAP*)**¹⁵: In its native form, EAP-PEAP does not support mutual certificate-based authentication. Native EAP-PEAP uses TLS to authenticate the backend directory server only and leverages a password-based challenge-response process to authenticate the client. Later PEAP extensions help to mitigate this weakness. EAP-PEAP is native to most Microsoft Windows versions and is therefore very prevalent across the wireless industry. There are currently three primary types of EAP-PEAP, which provide varying levels of protection:
 - **EAP-PEAP-MS-CHAP-v2 (*Microsoft Challenge Handshake Authentication Protocol*)**¹⁶: The most common EAP-PEAP inner authentication method uses Microsoft's CHAP-v2 protocol to provide user-identifier (userID) and password challenge-response-based authentication. The MS-CHAP authentication process has historically allowed an attacker with physical access to the AP, to capture the provided challenge and challenge-response hash in plaintext. Due to recent cracking advancements, the time to crack an

33 · 12 802.11 WIRELESS LAN SECURITY

MS-CHAP-v2 challenge-response hash has been reduced to days or less and should be avoided at all costs unless proper client supplicant configurations can be established.¹⁷ This process is described in section 0 of this chapter.

- **EAP-PEAP-TLS¹⁸:** PEAP-TLS is the second PEAP inner protocol defined by Microsoft. PEAP-TLS tunnels the EAP-TLS protocol within PEAP to provide mutual X509 certificate-based authentication.
 - **EAP-PEAPv1 (EAP-GTC):** The third implementation is defined by Cisco, which allows authentication using generic token cards such as RSA's SecurID token as well as user ID and password.
- 3. EAP-LEAP¹⁹:** A proprietary protocol developed by Cisco, which performed challenge-response MS-CHAP-v2, based username/password authentication in cleartext. An attacker targeting a network employing EAP-LEAP only needs to monitor traffic to capture challenge-response hashes for offline dictionary attack. This is in contrast to EAP-PEAP-MS-CHAP-v2, which requires a rogue masquerading AP for an attacker to intercept MS-CHAP-v2 challenge-response hashes.
- 4. EAP-FAST²⁰:** (Flexible Authentication via Secure Tunneling) was developed by Cisco as a replacement to the vulnerable LEAP protocol. EAP-FAST introduced secure pre-authentication tunnels without using certificates. EAP-FAST has secure mutual authentication capabilities, but also has an automatic PAC provisioning option—which is vulnerable to man-in-the-middle attacks.

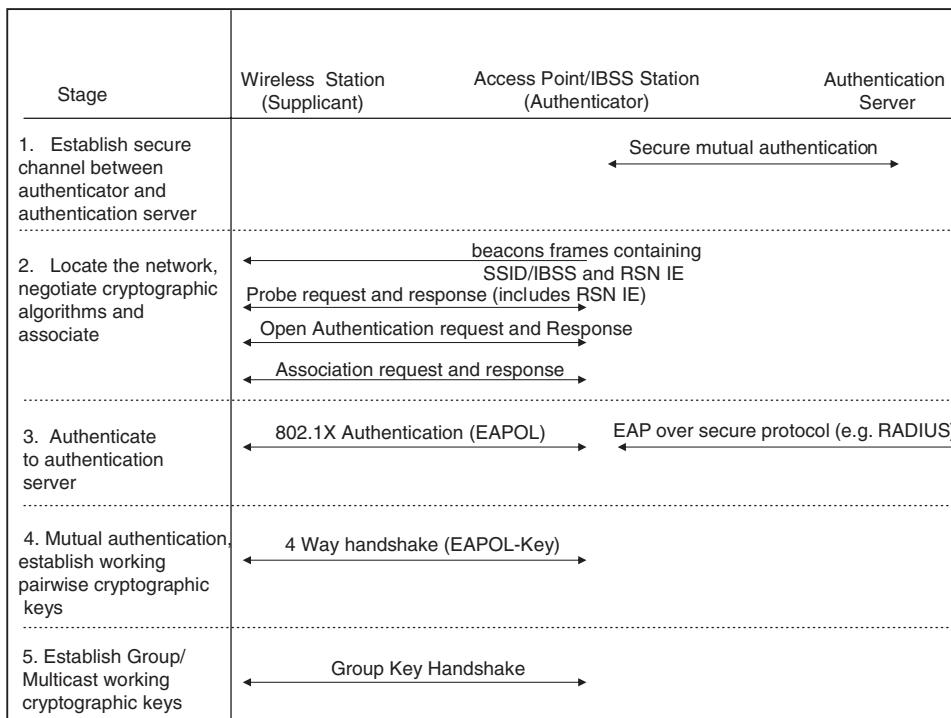
33.3.4 Insecure Legacy EAP protocols. When EAP was first integrated into 802.11, the EAP-MD5 and Cisco's EAP-LEAP protocols were commonly used. Unfortunately, these turned out to be vulnerable to attack, which spurred the development of the secure protocols covered above. EAP-MD5 and EAP-LEAP should not be used for enterprise wireless deployment.

33.3.5 Detailed EAP/EAPoL Procedure. Exhibit 33.3 shows the high-level RSN security association management flow. It consists of five stages that:

1. Establish a secure channel between the authenticator and authentication server
2. Locate the network, negotiate cryptographic algorithms, and associate to it
3. Provide 802.1X authentication to the authentication server
4. Provide Mutual Authentication and establish pairwise cryptographic keys
5. Establish Group/Multicast cryptographic keys

The following sections describe these stages. There are two aspects that vary the security association flow. The first is whether the wireless network contains an Access Point (a *Basic Service Set* or BSS) or whether it is an *independent BSS* (IBSS), also known as an *ad-hoc network*, which is a peer-to-peer network topology.²¹ The second is whether the master cryptographic key is a global PSK or if it is established during the 802.1X authentication protocol. Any variations caused by BSS/IBSS and 802.1X/PSK are described within the sections.

Stage 1—Establish a Secure Channel between Authenticator and Authentication Server

IEEE 802.11 ROBUST SECURITY NETWORK 33 · 13**EXHIBIT 33.3** RSN Security Association Management

In this stage the authenticator and AS mutually authenticate one another and establish a secure channel between them using a protocol such as RADIUS, IP Security (IPSec), or TLS. This channel is used to securely carry the authentication exchanges between the supplicant and the AS. This stage is not required if the network uses a PSK.

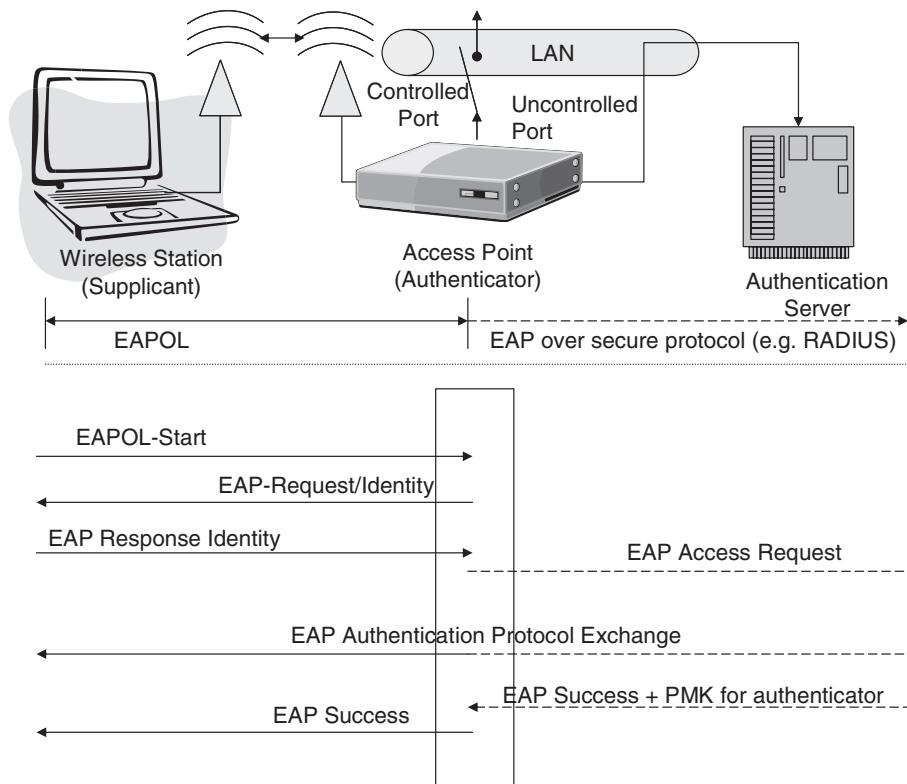
Stage 2—Locating the Network and Associating To It

This stage is mostly the original 802.11 functionality for locating, authenticating, and associating to a wireless network. The key difference in RSN is the beacon frames, probe responses, and association requests contain Information Elements that indicate supported and available authentication and privacy protocols. In addition, the fast BSS transition protocol defined in 802.11r²² enhanced these frames again to speed up AP roaming.

Stage 3—802.1X Authentication to the Authentication Server

The purpose of this stage is to mutually authenticate the supplicant and AS to one another and independently generate the PMK for use in stage 4. The EAP described above is used to achieve this. The messages exchanged between the supplicant and the AS are defined by the EAP method. An overview of this stage is given in Exhibit 33.4.

For a BSS network, the wireless client is the 802.1X supplicant and the AP is the authenticator. The AP relays the authentication messages between the supplicant and the authentication server, which can either be a separate service or built in.

33 · 14 802.11 WIRELESS LAN SECURITY**EXHIBIT 33.4** 802.1X Authentication to Authentication Server

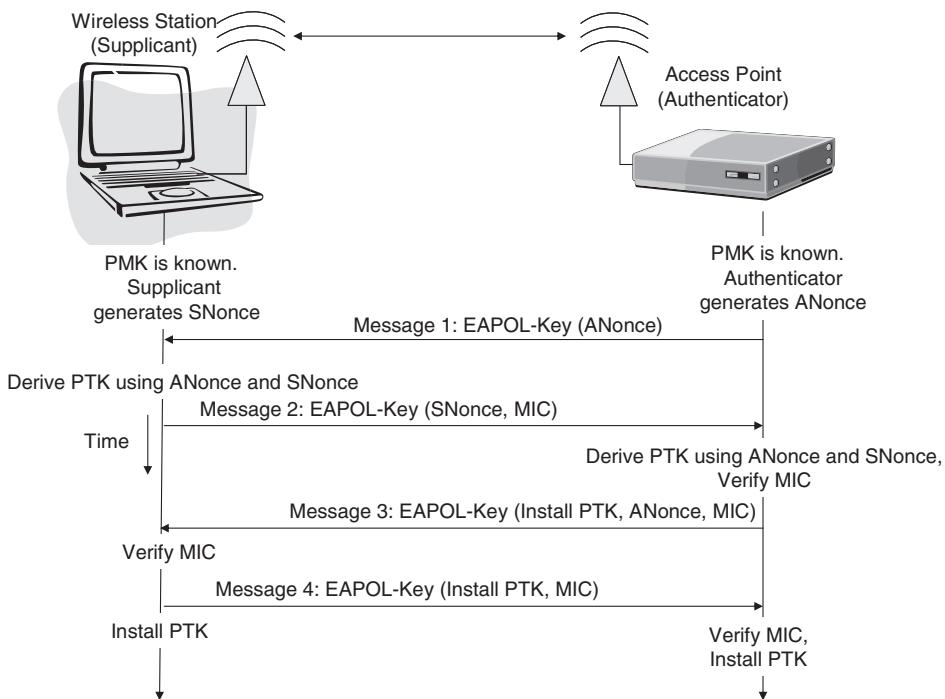
For an IBSS, the client wanting to associate to another client is the supplicant, and the target client is the authenticator. This means that clients in an IBSS can simultaneously be supplicants and authenticators depending on who initiated the association. Additionally, each IBSS client will need an authentication server unless the network uses PSKs.

As covered in stage 1, a secure channel has to be established between the authenticator and authentication server prior to this exchange; this serves two purposes. The first is to protect the integrity and authenticity of the authentication exchange, the second is to allow the AS to securely send the PMK to the authenticator once authentication is complete.

Stage 4—Mutual Authentication and Establish Pairwise Working Keys

Stage 3 established a PMK in both the supplicant and authenticator. If the network uses a *pre-shared key* (PSK) then the PSK is the PMK. This stage is called the 4-way handshake and has the following purposes:

1. To mutually authenticate the supplicant and authenticator to one another by confirming they both have the same valid PMK.
2. To generate a *Pairwise Transient Key* (PTK) from the PMK and fresh temporal keys bound to their MAC addresses.
3. To synchronize the installation of the keys in both devices.

IEEE 802.11 ROBUST SECURITY NETWORK 33 · 15**EXHIBIT 33.5** The 4-Way Handshake To Authenticate and Establish PTK

The 4-way handshake is implemented by EAPoL-Key messages exchanged between the supplicant and the authenticator and consists of the following as shown in Exhibit 33.5.

1. The authenticator and the supplicant both generate nonces for use in the authentication protocol. The authenticator's nonce is called ANonce and the supplicant's nonce is called SNonce.
2. The authenticator sends an EAPoL-Key message containing the ANonce (Message 1).
3. The supplicant derives the PTK using the ANonce and SNonce and calculates EAPoL-Key Encrypting Key (KEK) and EAPoL-Key Message Integrity Code (MIC).
4. The supplicant sends an EAPoL-Key message containing the SNonce and an MIC calculated using the EAPoL MIC Key (Message 2).
5. The authenticator can now derive the PTK because it has both the ANonce and SNonce. It then calculates the EAPoL-Key KEK and EAPoL-Key MIC Key and verifies the MIC in Message 2.
6. The authenticator sends a message containing the ANonce, and a flag instructing the supplicant to install the key. The message is also authenticated by an MIC (Message 3).
7. The supplicant verifies the MIC and sends a message to the authenticator confirming the installation of the key. This message is authenticated by an MIC and encrypted using the EAPoL-Key KEK (Message 4).

33 · 16 802.11 WIRELESS LAN SECURITY

8. The authenticator installs the new keys and starts the last stage to establish the group keys.

Stage 5—Establish Group/Multicast Cryptographic Keys

In Stage 4, the PTK and temporal keys were established. These keys are used to secure an additional message pair in which the authenticator sends the group temporal key (GTK) in an encrypted form to the client. The group keys are used to secure broadcast messages such as ARP requests and multicast traffic.

In a BSS network, all clients have the same group/multicast key, which is sent to each client by the AP. However, for IBSS networks, there is no AP to set a common key, so each client has its own group transmit key that it sends to all clients in the IBSS. This is achieved by executing the 4-way handshake and group key handshake in both directions.

33.3.6 RSN Key Hierarchy and Management. A key-management system derives working (temporal) keys from a root master key. The exact key hierarchy varies slightly between TKIP and CCMP, but broadly follows the same system. This section describes the key-management system and notes the differences between TKIP and CCMP.

There are two key hierarchies: the first contains the pairwise keys, which are shared between two wireless devices (e.g., between two clients in an IBSS or between a client and an AP in a BSS). The second hierarchy is the group/multicast keys that are used for network broadcasts such as ARP requests or multicast traffic.

33.3.6.1 Pairwise Key Hierarchy. There is a maximum of four levels of cryptographic keys in the pairwise key hierarchy.

1. **The 802.1X Authentication Keys.** These keys only exist if the supplicant and authentication server mutually authenticate one another using preinstalled keys. An example is EAP-TLS, which requires both the supplicant and authentication server to have PKI credentials. The 802.1X keys are used to establish the pairwise master key (PMK).
2. **Pairwise Master Key.** The PMK is 256 bits and is either the key established in level 1 or a PSK installed in the devices. The PMK, device MAC addresses, and nonces are fed into a variable length *pseudorandom function* (PRF) to generate a pairwise transient key (PTK). The nonces are derived from a 256-bit key counter that is initialized at system start-up from a random number, the time, and the MAC address. The nonce is then incremented every time a key is changed.
3. **Pairwise Transient Key.** The PTK varies in length from 128 bits to 512 bits, and is split up into the required temporal keys. The key lengths and number of keys are dependent on the algorithm. For TKIP, the PTK is 512 bits and this is split into 5 temporal keys. For CCMP, the PTK is 384 bits and is split into 3 temporal keys.
4. **Temporal and per-packet keys.** The temporal keys are mixed with variable data such as packet counters, which has the effect of creating a fresh key for every packet. See Exhibits 33.6 and 33.7 for a description of the TKIP and CCMP temporal keys.

IEEE 802.11 ROBUST SECURITY NETWORK 33 · 17**EXHIBIT 33.6 TKIP Temporal Keys****TKIP Temporal Keys**

Key Name	Length	Purpose
EAPol-Key MIC Key	128 bits	This key is the first 128 bits of the PTK and authenticates the EAPol-Key messages that establish the temporal keys. The key is used to calculate a message integrity code (MIC).
EAPol-Key Encryption key	128 bits	This key is bits 129–256 of the PTK and is used to encrypt the contents of EAPol-Key management messages.
Temporal Encryption Key	128 bits	This key is bits 257–384 of the PTK and is used to encrypt a packet using WEP.
Temporal MIC Key 1	64 bits	This key is bits 385–448 of the PTK and is used to authenticate messages in one direction.
Temporal MIC Key 2	64 bits	This key is bits 449–512 of the PTK and is used to authenticate messages in the other direction.

Exhibit 33.8 shows the entire key hierarchy for TKIP.

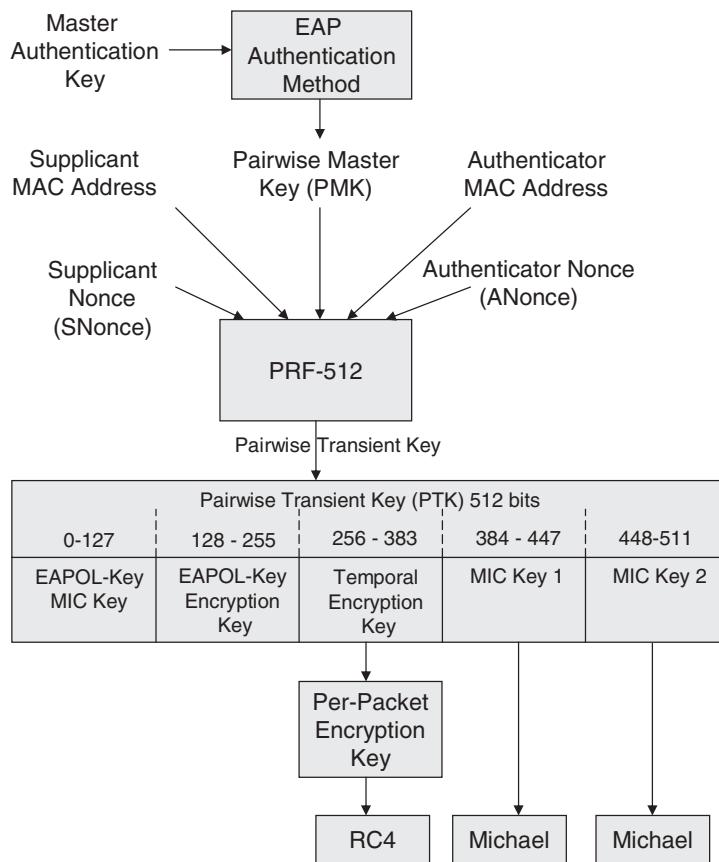
33.3.6.2 Group Key Hierarchy. The group key hierarchy is similarly structured to the pairwise key hierarchy. The authenticator creates a group master key. This master key, the authenticator's MAC address, and a *group nonce* (GNonce) are processed by the PRF to produce a group transient key (GTK), which is then split up into temporal keys.

33.3.7 Temporal Key Integrity Protocol (TKIP). The Temporal Key Integrity Protocol (TKIP) is a suite of algorithms to resolve the known key management issues in WEP for the existing 802.11 equipment already in the field. A key management mechanism was never specifically identified in the WEP protocol, which TKIP was designed to address.

To ensure compatibility with the existing base of wireless products, TKIP needed to take into account the hardware architecture of the existing wireless products.²³ Wireless products have two CPUs; the first is in the MAC chip that implements the wireless protocol, the second is the host CPU.

EXHIBIT 33.7 CCMP Temporal Keys**CCMP Temporal keys**

Key Name	Length	Purpose
EAPol-Key MIC key	128 bits	This key is the first 128 bits of the PTK and is used to authenticate the EAPol-Key management messages by calculating a MIC
EAPol-Key Encryption key	128 bits	This key is bits 129–256 of the PTK and is used to encrypt the contents of EAPol-Key management messages
Temporal Key 1	128 bits	This key is bits 257–384 of the PTK and is used for message encryption, authentication, and integrity. The nature of the CCM algorithm does not require a separate MIC key.

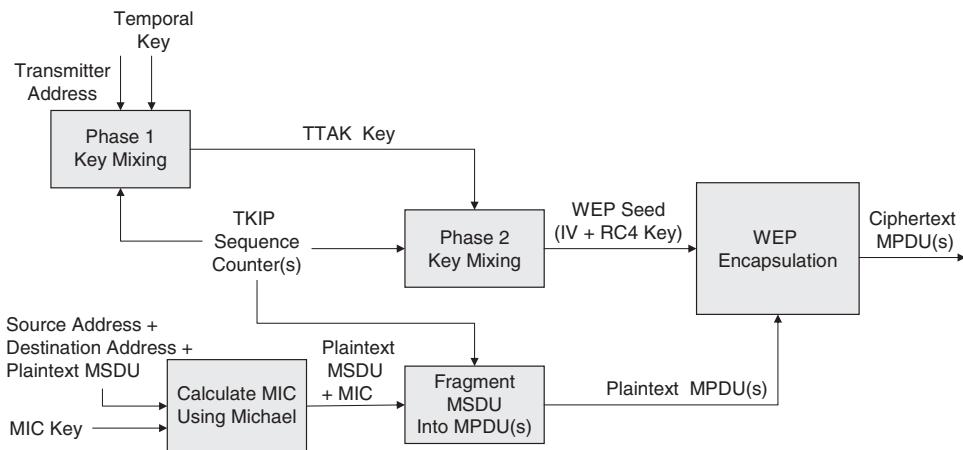
33 · 18 802.11 WIRELESS LAN SECURITY**EXHIBIT 33.8** TKIP Key Hierarchy

For APs, the host CPU is a dedicated CPU that runs the AP. For wireless LAN PC cards the host CPU is that of the host PC. For performance reasons most MAC chips have a hardware RC4 encryption engine to perform the encryption. To avoid an unacceptable performance impact by moving encryption to the host CPU, TKIP needed to use this existing MAC encryption engine. This required TKIP to continue to use RC4 and WEP but in a way to overcome the identified problems.

The designers came up with a solution to the issues, and the result was TKIP with the following functionality:

- Each MAC service data unit (MSDU) is authenticated and integrity protected using a keyed cryptographic MIC. An MSDU is a message to be sent to another client. An MSDU may be fragmented into more than one MAC protocol data unit (MPDU), which are packets/frames sent by the physical layer. When all the MPDUs have been received, the recipient reconstructs the MSDU and passes it up the protocol stack. The source and destination addresses as well as the MSDU plaintext are included in the MIC calculation. This prevents forgery and masquerading attacks.
- However, the strength of the MIC is limited and can be compromised by trial and error. To overcome this shortcoming, TKIP provides optional countermeasures that stop communications for a period of time (60 seconds) when an invalid MIC

IEEE 802.11 ROBUST SECURITY NETWORK 33 · 19

**EXHIBIT 33.9** TKIP Algorithm in the Transmitting Client

is received, and then immediately forcing a re-key of all clients. The designers estimated that the MIC with countermeasures would resist attack for about one year before an attacker would correctly guess the MIC.

- Each TKIP MPDU (packet) has a sequence number encoded in the WEP Initialization Vector. Any MPDUs that arrive out of order are dropped.
- TKIP mixes a temporal key, the transmitter's address, and a sequence counter in a two-phase protocol to form the RC4 WEP seed. The mixing is done in a way to defeat weak-key attacks.
- RSNA uses 802.1X EAPoL-key message to regularly change the temporal keys so that RC4 key streams are not reused. The key change is triggered automatically when the sequence counter is close to exhaustion.

These controls largely address the issues identified in the legacy system by introducing cryptographic message authentication, preventing key stream reuse, and never using a weak IV. However, TKIP is still not considered a strong solution because of the strength of the MIC. TKIP should only be used as an interim measure until existing equipment can be replaced with CCMP-capable equipment. Exhibits 33.9 and 33.10 show the TKIP process in the transmitting and receiving clients.

33.3.8 Counter-Mode/CBC-MAC Protocol (CCMP). CCMP is the mandatory protocol defined in RSN to provide confidentiality, authentication, integrity, and replay protection for the next generation of wireless equipment. It is not possible to upgrade existing equipment to use this protocol due to resource requirements on hardware, which implements AES.

CCMP uses AES encryption in counter mode to provide confidentiality and CBC-MAC (AES-CCM) for message authentication and integrity. The inputs to the algorithm include:

- A 128-bit block cipher encryption key (the temporal key).
- A nonce based on an incrementing packet number that is used only once with the encryption key to encrypt a message. Reusing a nonce to encrypt more than one message destroys its security properties.

33 · 20 802.11 WIRELESS LAN SECURITY

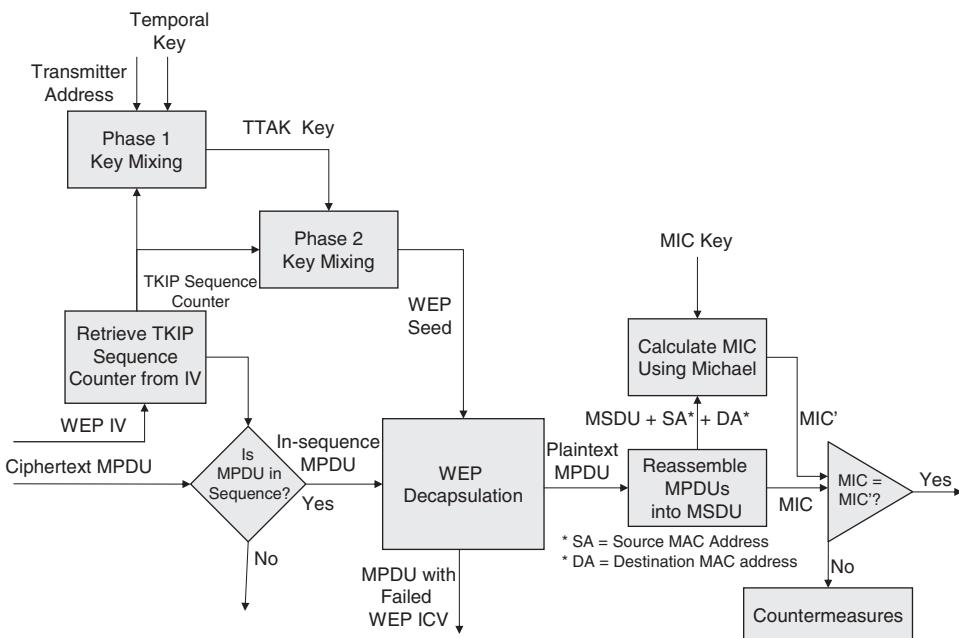


EXHIBIT 33.10 TKIP Algorithm in the Receiving Client

- The message to be encrypted.
- Additional data to be authenticated but not encrypted, such as the packet header containing the source and destination MAC addresses.

33.3.9 Fast Secure Roaming. After studying the RSN authentication system, the reader is likely to be struck as to the complexity of the system and the number of messages required to establish an association to the network. In an enterprise scenario using the 802.1X/EAP framework, roaming between APs could take up to one second to complete. This creates usability issues for time-sensitive applications such as mobile voice.

RSN contains the PMK caching and preauthentication protocols to try and address this problem; however, both protocols have limitations. PMK caching involves the client and APs remembering the keys they used with each other when the client roams away. But this is only useful when roaming back to APs you have roamed from. Preauthentication addresses this issue by completing the 802.1X authentication and establishing keys with APs as they come across them, even if the client never roams to that AP. This creates significant unnecessary authentication traffic that could overwhelm the authentication server.

The 802.11 committee decided to look at the problem again and developed the 802.11r standard that enhanced RSN PMK caching to establish a cached PMK with a central controller, which then distributes a variant of the PMK to APs as needed. This protocol is called Opportunistic Key Caching (OKC) and means that the lengthy 802.1X authentication is no longer required, reducing roaming times to approximately 100 milliseconds. 802.11r also enhanced the authentication and association messages to include the OKC parameters saving the need for additional message pairs.

FUNDAMENTAL WIRELESS THREATS 33 · 21

33.4 FUNDAMENTAL WIRELESS THREATS. It's important to ensure an enterprise is prepared to manage the many fundamental threats which are enabled by the existence of a wireless medium. As an enterprise examines their infrastructure and needs as a whole, there is an inherent lack of consideration for the increased threat exposure due to the elimination of physical layer controls. These types of controls are typically taken for granted and not considered despite the change to a completely new physical medium. It is important to understand the fundamental differences between a wired and wireless medium and be prepared to take action through implementing necessary controls to secure a wireless network to a level equivalent to a wired LAN.

The implementation of a wireless network expands network availability beyond that of facility walls or even parking lots. It's not uncommon for adversaries to sit within a corporate lot and use special antennas (called *cantennas*) to eavesdrop on wireless communications from a safe distance. Wireless networks rely primarily on logical controls which occur at data link layer and above in the form of authentication and encryption to compensate for the fundamentally public availability of the medium.

Driven by the ease of accessibility of wireless communications, the threat model an organization faces due to the existence and use of a wireless network may change as previous threats are eliminated (such as unauthorized network jack use) and others are created. The increasingly complex threat landscape and security-monitoring responsibilities can affect an enterprise's decision to implement a wireless solution for a particular use or limit its scope to designated functional uses.

Within this section, we will cover the fundamental threats enabled by the use of a wireless network and we will take a deeper dive into a subset of specific technical threats, which should be a priority consideration for an enterprise. It's not necessary for an enterprise to understand every threat they may face, but it is necessary to acknowledge their presence and take actionable steps toward countering these threats.

33.4.1 Unauthorized Network Extensions. When transitioning to a wireless environment from a wired network or a combination of the two, threats which were seemingly easy to identify in the past are now increasingly difficult to manage. Due to the inherent lack of transparency in regulated and unregulated airwaves, it is possible for a nonauthorized access point (AP) to hide amongst authorized APs through SSID spoofing. In a wired environment, it is more difficult to hide an unauthorized device due to physical appearance. An enterprise's ability to identify an unauthorized AP will require additional controls to perform this operation after authorized wireless access has been deployed.

33.4.1.1 Rogue Access Points. A rogue AP is attached to the network without authorization or which impersonates an attached device. Rogue APs come in multiple forms:

- **Nonmalicious Internal Rogue AP:** This type of rogue AP is commonly implemented by an employee or staff without malicious intent, but may be abused by an outside attacker. By connecting a Small Office Home Office (SOHO) grade AP to the internal network, enterprise grade protection mechanisms designed to keep attackers out, such as perimeter firewalls, can be bypassed with direct access to the internal LAN. Endpoint security mechanisms are also invalidated, such as host-based antivirus and firewalls, due to the Network Address Translation (NAT) capability of an SOHO router. Consumer-grade security controls present

33 · 22 802.11 WIRELESS LAN SECURITY

on an SOHO AP, combined with the lack of proper security configuration by the installer, can lend themselves to successful compromise, which would result in direct access to the wired internal network. Nonmalicious rogue APs can typically be identified through a nonenterprise standard SSID such as “Linksys” or “netgear.”

- **Attacker Malicious Internal Rogue AP:** An attacker can masquerade an AP as a legitimate enterprise AP by using the same SSID as authorized enterprise APs. In this situation an unsuspecting endpoint client will automatically connect to any AP, which broadcasts the requested SSID (such as “employeeWireless”) and will select the AP with the strongest signal. An attacker can then perform denial of service attacks to force clients onto their malicious AP.

A malicious rogue AP with internal access can be leveraged by an attacker with physical access to the facility to capture credentials or other data from unsuspecting endpoint clients, which leverage the AP to access internal resources such as administrative panels or employee portals. An attacker with an internally connected AP has multiple benefits as they are able to seamlessly allow a user to access all regularly accessible resources. This type of attack would require physical access to the facility and corresponding network jacks.

- **Attacker Malicious External Rogue AP:** An external rogue AP uses similar tactics to an internal rogue AP to force unsuspecting endpoint clients to connect to it. The primary difference is the physical proximity of the AP and the network it provides access to. It is much easier for an attacker to deploy an external rogue AP in a nearby corporate parking lot than gain access to the facility and plant an internally connected AP. The downside of an external rogue AP is the inability to provide requested internal resources to a connected client. If the desired resource is the Internet, this can be provided through a forwarding process allowing the attacker to view all traffic.

Both types of malicious APs allow an attacker to capture sensitive network traffic and attack the endpoint itself by attempting to exploit vulnerabilities in the endpoint’s operating system or network services. Client configurations typically mandate that a client wireless radio connect to the strongest signal strength that matches the desired SSID and Authentication requirements; e.g., WPA/WPA2-PSK, WPA/WPA2-Enterprise (EAP-PEAP, EAP-TLS).

Threats enabled by rogue APs include:

- An attacker’s gaining access to the wired network by compromising weak or absent security controls on a nonmalicious employee-installed rogue AP
- An attacker’s obtaining sensitive data through network captures of clients connected to an attacker-provided rogue AP (such as domain credentials). This threat is usually more effective if used in conjunction with social engineering tactics, such as injecting IFRAMES within HTTP traffic that directs them to a malicious site, which contains browser exploits for various plugins such as ActiveX or Adobe.
- An attacker’s gaining access and penetration testing enterprise endpoint machines connected to an attacker-provided rogue AP at a different physical location.
- Internal users’ bypassing wired 802.1X port security configurations through the use of NAT, which allows multiple users to communicate through the WAN port of an installed SOHO AP.

FUNDAMENTAL WIRELESS THREATS 33 · 23

33.4.2 Layer 1 Attacks. Wireless mediums can be heavily relied upon by staff or business applications to provide connectivity to wired network resources in locations where physical connectivity is difficult to implement. The absence of an enclosed private physical medium allows for tampering of the layer 1 wireless carrier signal regardless of upper layer protective controls. This is a fundamental threat to any form of wireless communication and not necessarily specific to 802.11, although each wireless protocol such as CDMA, 3G, or 802.11 may include a mechanism to detect tampering, they cannot prevent it from occurring.

- The 802.11 protocol functions within 2.4 GHz band. The 2.4 Ghz consumer band can be interfered with by using something as simple as a microwave oven
- The 2.4 GHz spectrum only provides three nonoverlapping channels. The 802.11 spectrum runs between the range of 2.412 GHz and 2.462 GHz. Overlapping channels are side-by-side channels, which through propagation and other wireless tendencies, can bleed into the neighboring 20 MHz bandwidth allocation. A nearby AP can cause a minor denial of service condition by using an overlapping channel.
- Management frame spoofing
 - 802.11 management frames are part of the MAC layer of the 802.11 protocol. Management frames are used to bootstrap the authentication and association of clients to an AP. Management frames “manage” the connection between an AP and all connected clients. There are various types of management frames, some of which are received by any client within range of the AP, and others which are targeted at connected/connecting clients.

Most significant management frame types:

- Probe Request/Response
- Authentication
- Association Request
- Association Response
- Dissociation

As with a wired LAN, the unauthenticated nature of initial data link layer protocols, such as 802.11 management frames, can allow an attacker to spoof and interfere with an existing session or the session initiation process. This fundamental issue is comparable to the unauthenticated Ethernet → IP ARP protocol, still vulnerable over 802.11, but easily abused across a wired connection.

The security afforded to wired bootstrap protocols relies on their physical controls, none of which are afforded to wireless. This void prompted IEEE 802.11 family of standards to adopt the 802.11w specification. This specification details the use of MICs to authenticate management frames to their source. This would effectively eliminate the ability for an unauthenticated attacker to force clients off of a legitimate AP to an attacker provided rogue AP. 802.11w sounds great on paper, but in reality, implementation will take years. The process of authenticating material for each management frame increases the load on an AP's already busy processors, which means a software upgrade may not be feasible for all devices. In addition, all client devices must support 802.11w to create outgoing and process incoming MIC data. Client support would require new

33 · 24 802.11 WIRELESS LAN SECURITY

driver firmware for wireless radios and potentially more hardware acceleration for potentially all connected clients depending on the AP's implementation of 802.11w.

33.4.3 Endpoint Attacks. A new class of endpoint attacks is enabled by wireless endpoints. Endpoints can include anything from a laptop, wireless printer, or handheld scanners. In a pure wired environment the threat of wireless endpoint compromise is slightly reduced due to the absence of corporate wireless authenticating materials stored within the endpoint's wireless supplicant. Authenticating materials include items such as domain username/password pairs, certificates, or even pre-shared keys. Storing authenticating materials on an endpoint means the "keys to the kingdom" can be obtained and used to access the network without even entering the facility.

While the ability for an attacker to compromise an endpoint with a hotspot style attack, is still a viable threat to the integrity of an endpoint, endpoint security measures such as host-based firewalls, IDS, and regular patching can mitigate this threat to an acceptable level. It is necessary to understand that an endpoint is still susceptible to the same attacks in an enterprise wireless environment, but the outcome of these attacks could result in the capture of authenticating materials used to access the corporate wireless. The capture of authenticating materials avoids the need to compromise the endpoint to gain internal access to a corporate network. If employees do not use wireless enterprise access, they are less likely to succumb to social engineering attacks combined with a hotspot attack. An example of this situation involves an attacker using a hotspot attack to gain connectivity to the endpoint, intercepting their Internet browsing, and prompting the user with a spoofed corporate intranet site, which requests remote VPN credentials through a simple Web form.

Endpoint attacks enabled by enterprise wireless access:

- Employer identification through descriptive SSID probe requests.
- Increased effectiveness of Hotspot + Social Engineering attacks based on the fact that a user regularly accesses the enterprise wireless, in contrast to a user whom has never had wireless access.
- RADIUS impersonation for capturing WPA/WPA2-enterprise EAP credentials.

33.5 SPECIFIC WIRELESS SECURITY ATTACKS

33.5.1 MAC Address and IP Spoofing. MAC address-based access control mechanisms are the lowest hanging fruit of wireless security threats. This form of access control is not enterprise grade and primarily originated in SOHO grade APs and hotspots for tracking authorized connections without issuing unique keying material. Pure MAC address authentication is based solely on the connecting devices MAC address; if the identified MAC address is included in a predefined list of allowed MAC addresses, the device is allowed to associate. This is obviously not a scalable means of authentication, due to constant updates in the event of static mappings, and is commonly used to track valid client connections after initial Web-based authentication.

This type of situation is commonly referred to as a Captive Portal or hotspot-based authentication. The AP operates in open authentication mode, allowing anyone to connect to it, but users who have yet to authenticate cannot access the Internet and are dropped into a restricted penalty box. Users are prompted with a Web interface to enter credentials and, upon successful network-based authentication, gain access to the Internet or other protected resources. Captive Portal authentication takes place at

SPECIFIC WIRELESS SECURITY ATTACKS 33 · 25

layer 4 (Application) and relies upon layer 3 (IP address) and layer 2 (MAC address) as reliable mechanisms to track users which have successfully authenticated.

This attack is enabled by the fact that all traffic in a Captive Portal system is transferred unencrypted. No encryption necessary to allow clients to access the Web-based authentication page with no previously established encryption keys or previous knowledge of the AP. An attacker now has the ability to capture plaintext traffic occurring on the AP to and from each client. This is performed through setting a compatible wireless card into monitor mode. Monitor mode allows an attacker to see all traffic on a specific 802.11 channel without needing to associate with an AP. On wireless networks which deploy an encryption mechanism such as WPA/WPA2-Enterprise, an attacker will see no plaintext payloads for 802.11 frames.

An attacker can leverage this to monitor traffic and identify MAC address and IP addresses of connected/authenticated clients. This information can then be spoofed by an attacker to gain access to the portal's restricted content (Internet or other resources) without being directed by the portal to the authentication page. Authenticated clients are maintained by the AP via MAC address, IP address, or a combination of both.

33.5.2 EAP (WEP/WPA/WPA-Enterprise) RADIUS Impersonation Attack.

Wireless implementations, which leverage RSN with clients which have not been configured to trust a specific AP type, and do not employ client mutual certificate-based authentication, are inherently vulnerable to a deadly combination of rogue AP and authentication server (e.g., RADIUS) spoofing. This attack relies on the ability for an attacker to force clients to connect to a malicious AP and begin the expected EAP authentication method. In a typical EAP RSN deployment, in response to an authentication request a RADIUS server will send back a RADIUS access challenge, which is forwarded by the AP to the client. In this attack, the malicious AP is configured to use an attacker provided RADIUS server. This RADIUS server carries out the EAP authentication process and forces the client into providing their password hash in response to a fixed, known, challenge. Depending on the EAP implementation in use, the username associated with the client may be provided in plaintext within the EAP-Identity-Response sent prior to authenticating material exchange.

Since the attacking rogue AP uses a fixed challenge we now possess the ability to brute force the password using standard dictionary attacks. This attack only applies to EAP implementations, which do not use client certificates for authentication, but leverage username and password combinations, such as EAP-PEAP with MS-CHAP-v2.

If a client connects to a rogue AP, the spoofed authentication server certificate will likely not validate correctly and the user will either not see anything or be prompted with a small message. See Section 6.2 for additional information pertaining to proper client supplicant configuration for server-side certificate verification.

EAP challenge/response hashes can be captured and cracked using the following tools:

- **FreeRADIUS WPE (Wireless Pwn Edition):** Used to readily deploy a RADIUS server to accept connections from the rogue AP using a spoofed version of the expected EAP type.
- **AsLEAP:** Used to crack EAP (PEAP, LEAP, etc.) challenge/response hashes using the password hash, known challenge, and word dictionary as input.

33 · 26 802.11 WIRELESS LAN SECURITY

33.5.3 Management Frames. 802.11 frames carry protocol and data for upper layer protocols. Primarily, two types of 802.11 frames are commonly used, data and management frames. Data frames are used as carriers for upper layer protocols, while management frames carry specific information regarding link operations. Historically, management frames are a fundamental weakness to 802.11 security. Across the majority of wireless implementations management frames are never secured and therefore always accessible despite authentication or encryption mechanisms, which may secure the payload of a data frame's contents. As described previously in section 4.1, 802.11w attempts to address this issue through the use of MICs, but like most wireless improvements, the industry lags behind fully adopting these controls due to costly client upgrades and new AP devices and/or firmware updates. The lack of adoption is likely due to the fact that management frame spoofing is not readily seen as a direct threat, despite its deadly impact when used in combination with other attacks.

Management frame spoofing and manipulation enables the following types of threats:

- **Beacon frame spoofing:** Allows an attacker to masquerade as any AP. Crafting a beacon frame will make you show up in a network client “scan” of available networks. This vulnerability enables hotspot and RADIUS impersonation attacks to take place and is largely unstoppable due to the necessity for unauthenticated beacon frames to be accepted by any unassociated client within range.
- **Authentication/Deauthentication:** These frames are heavily used by 802.11 during the secure authentication process. Used to relay the intention of a client to connect to an AP. Even in situations where an AP is using “open” authentication and requires no keying material an authentication frame procedure is performed. In the event of WEP pre-shared key authentication these frames are used to carry traffic associated with this network-based authentication.
- **Association/Reassociation/Disassociation:** Association response frames are used to synchronize radios between an AP and a client radio. This includes details such as supported data rates. This process allows the AP to allocate memory for the client and begin processing communications. Reassociation frames are used when a client roams to another AP within the same BSSID, this indicates to the original AP that it should forward any remaining buffered frames. A dissociation frame is sent by the client to gracefully remove itself from an AP and allow the AP to garbage collect any necessary memory.
- **Request to Send (RST) and Clear To Send (CTS):** RST frames are optional and provide a way to survive in a condition where one client cannot see another client's traffic due to physical proximity to the AP. This frame is sent to the AP and the AP will respond with a CTS frame. The CTS frame includes a narrow time frame for the client to send data where all other connected clients must wait. This procedure reduces collisions from hidden AP clients.²⁴

The manipulation of management frames occurs by injecting them into the airwaves and spoofing the source as if they came from the AP BSSID. These frames are not viewable from a standard packet capture on the host system, as the wireless card driver will not pass this traffic to the underlying operating system. Spoofing management frames does not affect an AP, as it is unaware of their presence.

SPECIFIC WIRELESS SECURITY ATTACKS 33 · 27

The following Aircrack²⁵ tools can be used with a compatible wireless card set into monitor mode to manipulate management frames²⁶:

- **Aireplay-ng -0:** Spoof deauthentication to remove a client from the AP by sending dissociate packets to one or more clients, which are currently associated by forging a disassociate packet with the source MAC address of the target.
- **Aireplay-ng -1:** Spoof authentication with an AP using a specified source and destination MAC address by sending association management frames.

33.5.4 Hotspot and Cafe Latte Attacks. Depending on the wireless chipset and driver, wireless clients typically actively probe for all networks they have associated with in the past and when an AP responds they will automatically associate with it. The root of this vulnerability stems from the behavior of wireless clients and their lack of AP authentication during the early stages of authentication. Only the SSID and authentication/encryption type must match the “known” state for the client to begin the association process. Due to thin AP wireless deployments, a wireless client cannot be specific or remember the MAC address of the AP it requests to connect to, even though this would greatly curb this attack.

A wireless client probing for an SSID last used from a public Wi-Fi hotspot can be identified by an attacker and masqueraded as if they were the public Wi-Fi hotspot, even though the victim is nowhere within the location of the previously remembered public hotspot. Since most public hotspots don’t require authentication, it’s safe for an attacker to assume they only need to match the requested SSID and no other traits to facilitate a successful connection by the targeted victim. If the AP a client is probing for requires WEP or WPA, it will not automatically connect to any AP with the same name. The data transfer rates and encryption levels must align correctly, as described within the process of association.

An attacker sniffing the wireless for SSID probe requests will not necessarily know which networks require encryption and which don’t, but trial and error will eventually yield the correct answer.

The following is the process associated with hotspot spoofing where an attacker recognizes a particular wireless client would like to connect to an AP, that may or may not actually be present (this is the nature of Wi-Fi autoconnect), and actively impersonates the requested AP to satisfy the client and gain momentary connectivity (in the event the client disconnects due to an invalid key management process).

33.5.4.1 Hotspot Attack Procedure

1. Victim machine sends SSID probe requests for APs it has communicated with in the past (supplicants keep this history to enable fast connections).
2. Attacking machine monitors airwaves for SSID probes and takes note of requested AP names (e.g., “coffeeshop,” “airportWifi,” etc.).
3. Attacking machine creates an AP with the SSID name requested by the client (e.g., “coffeeshop,” “airportWifi,” etc.).
4. Victim machine automatically associates with new AP and encryption/authentication levels are accurate.

At this point it is up to the attacker to correctly guess the encryption the client expects (WPA/WPA2/WEP) or open (none). It is very common for

33 · 28 802.11 WIRELESS LAN SECURITY

Microsoft Windows and third-party supplicants to autoconnect to a known AP in the background and show no warning of the connection with the exception of a “connected” symbol in the taskbar.

5. The victim now has a functional connection to the attacker’s AP and the attacker can establish *man-in-the-middle* (MITM) traffic by relaying it to the legitimate AP on a neighboring interface or the Internet, or attack the victim by exploiting vulnerabilities on the connected system.

An older yet valuable extension to this attack revolves around the SSID, which is broadcasted by a “parked” wireless card and is prevalent in older Windows wireless drivers. During Windows Wireless Zero Configuration, a wireless card will attempt to connect to all previously known networks before placing itself in a “parked” state. The “parked” condition consists of the adapter ceasing to attempt connections and using a dynamically generated SSID as a default value. Once this default value is assigned, the wireless card will no longer actively seek networks but will continue its normal beacon behavior, this time with the parked SSID. This behavior introduced a new vulnerability, which allowed an attacker to always have an SSID to rely on, even if the client has no previous networks to seek for. This vulnerability has since been addressed through setting an encryption level (WPA/WPA2) and random key for use by the wireless card during park mode.²⁷

33.5.4.2 The WEP Cafe-Latte Factor. An added twist to hotspot spoofing takes it a step further by actively attacking a client that broadcasts a particular SSID from a network which deploys WEP encryption. This attack leverages a multitude of WEP vulnerabilities to crack a WEP key passed by a client. All clients will store the WEP key within the supplicant or registry of some sort for later use. An attacker can force the client to encrypt certain packets using this key while performing the aforementioned hotspot spoofing attacks to draw the client onto an AP that is masquerading as a requested one employing WEP encryption (albeit with an incorrect matching key). Specific packets are always sent by a client upon associating and during initial authentication as part of “quick reconnect,” which are encrypted using the client’s stored WEP key. These include ARP, DHCP, and specific 802.11 management frames.

Of particular interest are gratuitous ARP packets sent by a client upon associating to an AP it believes is legitimate. Through bit flipping techniques an attacker can forge encrypted ARP requests to a connecting client by reusing previously captured gratuitous ARP packets and turning them into a request. This technique leverages WEP’s absence of integrity checking and reliance on a CRC error-checking mechanism as the sole integrity control. CRC algorithms are not intended to provide anything outside of error-based integrity and prone to abuse across WEP attacks. The compilation of hotspotting and WEP vulnerabilities can yield the WEP key to a corporate network within six minutes, without even needing access to the actual network. An attacker can now take his latte, your corporate WEP key, locate the corporate network, and connect from a parking lot.

33.5.4.3 Caffe-Latte Attack Procedure²⁸

1. Victim machine sends SSID Probe requests for APs it has communicated with in the past (almost all supplicants keep this history to enable fast reconnections).

SPECIFIC WIRELESS SECURITY ATTACKS 33 · 29

2. Attacking machine monitors airwaves for SSID probes and takes note of requested AP names (e.g., “coffeeshop,” “airportWifi,” etc.).
3. Attacking machine creates an AP with the SSID name requested by the client (e.g., “coffeeshop,” “airportWifi,” etc.).
4. Victim’s machine attempts to connect but fails, likely due to an encryption mismatch.
5. Attacking machine modifies AP to possess WEP encryption with a false key.
6. Victim machine attempts to connect to the AP again and successfully authenticates due to a properly matching SSID name and encryption grade.
7. Victim machine pauses during further network-based authentication due to a key mismatch and never fully associates.
8. Attacker captures gratuitous ARP, DHCP, and other packets sent by the victim after association, which are encrypted with the victim’s stored WEP key.
9. Attacker provides victim-encrypted gratuitous ARP packet to ./cafe-latte resulting in bit-flipping to produce a CRC-correct encrypted ARP request.
10. Attacker replays encrypted ARP request to solicit responses from victim with new IVs until enough IVs are captured to exhaust 40-bit key space.
11. Attacker executes Aireplay-ng (or comparable WEP cracking utility) providing IVs captured during ARP request replays.
12. Attack receives plaintext hex of WEP key and uses to gain access to target network.

The following tools can be used to perform a hotspot and cafe latte–style attack:

- **Airbase-ng:** Create an AP with specified SSID
- **Cafe-latte:** WEP cracking script, which takes the captured gratuitous ARP packet and bit-flips it into a request that can be repeated to the client thousands of time.

33.5.4.4 What about WPA/WPA2? WPA and WPA2 implement cryptographic integrity checks, which remove any ability to deduce a plaintext key, but still suffer from a hotspot style of attack. A WPA/WPA2 tailored attack involves masquerading as a known AP where a client expects WPA or WPA2 encryption. Upon providing the correct encryption level the client will try to connect and fail due to the AP not possessing the correct key. Failure is expected, as the pre-shared key hashes with a salt of the SSID has already been sent from the client and captured by the attacker. This hash can then be brute forced.

33.5.5 WEP. Despite the public and lengthy exposure of WEP vulnerabilities and little security it actually affords, its widespread use at the enterprise level is shocking yet usually unknown due to its use in many small crevices of the business that it may have wedged itself in years ago.

WEP was conceived during a time when computing power, and in particular embedded computing power, was hard to come by. The computing power required to carry out the WEP RC4 encryption and key management is very minimal, which coincided perfectly with the computing power of many handheld and other embedded devices. This very principle is still the reason why WEP remains in many environments and will for the foreseeable future. Sunk costs in infrastructure devices incapable of performing

33 · 30 802.11 WIRELESS LAN SECURITY

resource heavy functions such as AES and CCMP in WPA2. WPA's implementation of TKIP was precisely done to enable backward compatibility with WEP hardware. In today's enterprise environment it is commonplace to see WEP fully expunged from employee WLANs but prevalent in retail or distribution centers on handhelds and other portable devices.

Further compounding WEP's poor implementation of the RC4 stream cipher is the ambiguous definition within the standard for specific encryption function processes, largely leaving the decision to the vendor. This is particularly obvious in the way Initialization Vectors are generated; some vendors use a Pseudo-Random-Number-Generator (PRNG) and others merely start at 000000 and work their way to FFFFFF. We will briefly dive into the following tried-and-true core vulnerabilities affecting WEP, and its failure to provide proper key management, integrity, or encryption, and finish with a brief description of WPA/WPA2's advancements in this area.

33.5.5.1 (WEP) Repeating Initialization Vectors (IVs). WEP is home to a number of vulnerabilities including FMS, KoreK, PTW, and ChopChop. We are only going to cover the most famous and original attack, called FMS. WEP has been sunset ages ago in favor of alternative technologies; as such we will dive into the details of this vulnerability only to shed light on the advancements that have been made since its discovery.

WEP's underlying encryption revolves around the RC4 stream cipher. An initialization vector (IV) is used to ensure two identical inputs do not result in the same cipher text output. The IV is leveraged to lengthen the key life by uniquely generating an IV for each frame. WEP keys come in two forms depending on the version, an effective 40 bits (64 including IVs) with some implementations going as high as an effective 104 bits (128 including IVs) WEP uses a 24bit IV which, on a busy network, there is a 50 percent probability that the IV will repeat after 5,000 packets.

A fundamental solution to the WEP IV key scheduling vulnerability is the use of session keys established at the onset of authentication and encryption, which are used for a finite period of time during the encryption session. The process used to derive the temporary or session key is crucial to the success of a key scheduling process. WEP suffered from concatenation of RC4 output with an IV, and was improved upon with TKIP's "mixing" of these values prior to RC4. If a solid key scheduling and management function can be established, the encryption routine itself will deliver as expected. Attacks which target vulnerabilities in WEP's key scheduling process include the famous FMS attack, named after its inventors Fluhrer, Mantin, and Shamir.²⁹

33.5.5.2 (WEP) Key Management. WEP does not provide any form of key management, such as the ability to rotate keys with clients. The result of no key management structure is the reuse of stale keys as input into the key scheduling and underlying encryption process for an extended period of time, further compounding the repeating IV vulnerability mentioned earlier. WPA and WPA2 provide session-based key management through the implementation of TKIP (WPA) and later CCMP/AES (WPA2). It is important to outline that the secret key in a pre-shared key implementation (WPA/WPA2-PSK) is only really leveraged to provided authentication; encryption is derived from the client and AP possessing the same key, but the pre-shared key is not directly involved. This is in contrast to WEP's use of a pre-shared key directly in the encryption process by leveraging it for authentication and including it in the RC4 stream cipher. Through the use of EAP and RADIUS, WPA and WPA2 Enterprise editions fully support key management through directory-based structures. These directories manage authentication keys outside of the WPA/WPA2 implementation and add increased

SPECIFIC WIRELESS SECURITY ATTACKS 33 · 31

security through the combination of directory-based authentication key management (LDAP/Active Directory) and TKIP/CCMP encryption key management.

Tools used to exploit WEP vulnerabilities:

- **Airmon-ng**³⁰: Used to place the wireless adapter into monitor mode.
- **Airodump-ng**: Used to capture WEP IVs using monitor mode.
- **Aireplay-ng-1**: Fake authentication to the AP; this will allow you to associate to the AP but not continue the authentication process. This is necessary for the AP to receive data from your wireless adapter.
- **Aireplay-ng-3**: ARP request replay attack mode, this will listen for encrypted gratuitous ARP packets sent by a legitimate connected client. This mode identifies ARP packets by size and other factors despite being encrypted. Once an ARP packet is found, Aireplay will bitflip and turn it into a request which is replied thousands of times. The purpose of this is to generate lots of encrypted packets from the victim and start exhausting IV space.
- **Aircrack-ng-b**: Performs statistical analysis and brute force of captured WEP traffic to recover the key. The process is to use statistical analysis to zone in on potential byte values at each location of the key, then use brute force to complete the process.

33.5.6 WPA/WPA2 Pre-Shared Key. WPA and later WPA2 aimed to address many of the deficiencies in WEP and WEP advancements. The following are attack vectors and mitigating controls implemented in the WPA standard:

- **48-bit Initialization Vectors:** Addresses WEP's 24-bit vulnerability and poor RC4 implementation through the use of TKIP/CCMP to significantly reduce IV reuse and improve the ability to “seed” the encryption function.
- **Key Management:** TKIP and CCMP/AES provide per-frame unique encryption keys in contrast to stale WEP key storage within client supplicants and reliance on the IV as a primary differentiating factor. WPA/WPA2 also natively supports both PSK and enterprise based authentication mechanisms for third-party authentication key management (e.g., LDAP directory).
- **Key Length:** WPA/WPA2 provide a minimum 256-bit pre-shared key length.
- **Message Integrity Code:** WPA correctly provides message integrity through a hashing mechanism known as an MIC or Message Integrity Code. This is contrast to WEP's Integrity Check Value (ICV), which is comprised of an encrypted CRC-32 error-check. WEP's ICV is plagued by bit-flipping attacks due to its incorrect use as an integrity mechanism despite CRC-32's intended use as a transmission error detection mechanism
- **Mutual Authentication:** WPA/WPA2 provides the ability for an AP to be authenticated by a client through the use of certificates and RSN.

The core vulnerability which plagues WPA/WPA2 is not necessarily a vulnerability in the implementation, but an attribute of the necessary 4-way handshake process (EAPoL) for basic pre-shared key authentication. PSK is not designed for enterprise use (but nonetheless finds itself in the enterprise), which is why solutions exist within WPA/WPA2 to eliminate this vulnerability at the enterprise level but not at the SOHO level.

33 · 32 802.11 WIRELESS LAN SECURITY

Pre-shared key authentication is commonly used in absence of a backend directory normally accessed via RADIUS and LDAP protocols. This type of implementation is common for solutions, which are not intended to be used by a wide variety of clients and therefore the cost of implementing such a solution is not effective. Pre-shared key authentication relies on an out-of-band key-management system. WPA-PSK and WPA2-PSK do not provide any key-management or tunneled encryption of the pre-shared key authentication handshake process. The password hash procedure used by WPA/WPA2 combines the stored passphrase with a salt of the destination AP's SSID and sends it over in plaintext.

Reliance on a single attribute for access control means that each client is indiscernible from the next and no amount of accountability can be afforded to each user outside of a DHCP MAC address log.

33.5.6.1 WPA/WPA2 Key Hash Capture and Cracking. The WPA/WPA2 pre-shared key authentication process is a series of handshakes used to eventually transfer a hash of the passphrase to the AP. This process can be captured through monitor mode and used to extract the hashed passphrase. Since this process only occurs when clients are initially authenticated to the AP, an attacker can either wait for a new client or deauthenticate an existing client by spoofing management frames as described earlier in the chapter.

Once a WPA hash is obtained an attacker can attempt to crack it using brute force. Due to the salting of the hash with the SSID of the AP, rainbow tables cannot be used to assist in the cracking process. The use of a salted hash combined with a lengthy hash process makes for slow cracking of WPA hashes on CPU-based systems. An exhaustive dictionary attack on a CPU-based system can typically process between 5,000–15,000 pairwise master keys per second.

GPU-based cracking systems are becoming more popular, as programming interfaces such as Nvidia's CUDA opened the door for developers to access many components of a GPU that would normally require advanced programming experience. GPUs can be leveraged to greatly increase the probability of successfully cracking a captured WPA handshake to the tune of 50,000+ pairwise master keys per second for an average set of CUDA-enabled GPU graphics cards.

Tools used to capture and crack a WPA/WPA2 pre-shared key:

- **Airmon-ng³¹:** Used to place the wireless adapter into monitor mode.
- **Airodump-ng:** Used to capture traffic using monitor mode.
- **Aireplay-ng-0:** Deauthenticates a client from the AP and forces them to re-authenticate and perform the pre-shared key handshake. This tool is necessary to allow Airodump-ng to capture the pre-shared key during the re-authentication process.

33.5.7 WPA/WPA2: Wi-Fi Protected Setup (WPS). Wi-Fi Protected Setup allows users with little understanding of wireless security configurations to easily configure a pre-shared key configuration without necessarily logging into the router. WPS describes the protocol and *personal identification number* (PIN) sequence used to negotiate a WPA/WPA2 pre-shared key without the need to manually configure it on the client and AP. WPS is only provided on SOHO-grade APs and not enterprise-grade hardware. This vulnerability may apply to the enterprise environment through its

SPECIFIC WIRELESS SECURITY ATTACKS 33 · 33

existence as a nonmalicious rogue AP, which an internal user has securely configured but is vulnerable to WPS.

Two types of WPS exist:

- Push-Button-Connect (PBC)
- 8-digit PIN

A PIN may be provided by the AP for input into the client WPS agent or the client WPS agent may provide a PIN for input into AP WPS prompt. WPS authentication is based on 802.1X EAP and a vendor-specific implementation of EAP, which uses a Diffie-Hellman Key Exchange to prove possession of the first half of the PIN between the AP and client and eventually the second half.

Because WPS divides the PIN into two separate authentication steps, this process severely limits the number of combinations for each half of the PIN. An attacker can derive the first and the second part of the PIN through abusing an EAP message, “EAP-NACK.” This message is issued by the AP on failing to provide the correct first or second portion of the PIN. Brute forcing each *separate* portion of the PIN requires only 10^4 tries at most, so there are only $10^4 + 10^4 = 20,000$ possible combinations instead of the original keyspace of 10^8 or 1,00,000,000.³²

On average, WPS cracking tools can recover the WPA/WPA2 plaintext password in 4–10 hours.

Tools used to perform a WPS attack:

- **Airmon-ng:** To detect target AP and place adapter into monitor mode.
- **Reaver³³:** Brute forces the WPS EAP process to recover the WPA/WPA2 key.
- **Wpscrack.py³⁴:** Alternative tool to brute force the WPS EAP process to recover the WPA/WPA2 key.

33.5.8 MS-CHAP-v2—Divide and Conquer. EAP implementations that leverage MS-CHAP-v2 have already been discussed as vulnerable to dictionary attacks with the use of EAP-LEAP or EAP-PEAP-MSCHAP with poor supplicant configurations. Traditionally, the challenge-response hash-cracking process where a challenge-response hash is provided to a cracking program along with a dictionary was the only feasible approach to exploiting this vulnerability.

Thanks to research by Moxie Marlinspike (also the author of SSLStrip), this process was recently expedited through a new understanding of the seemingly painfully obvious and much overlooked functionality within the MS-CHAP-v2 protocol. During MS-CHAP-v2, the MD4 hash of a user’s password is bitwise divided to construct three individual DES keys (seven bytes each). A typical penetration-testing tool such as *Asleep* will input the MD4 hash, break it into its three DES keys and use those keys to encrypt a given plaintext dictionary value for comparison to the original ciphertext hash.

The traditional cracking approach would result in a total keyspace of $2^{(56+56+56)}$ or about 10^{50} —a very large number that’s infeasible for incremental brute-force cracking. This is exactly the reason why dictionaries were traditionally used as a “best guess” approach, because incremental (character-by-character) brute forcing is just not an option for any reasonable amount of time due to dictionaries’ salting through the use of a challenge.

33 · 34 802.11 WIRELESS LAN SECURITY

The use of three DES keys is prime for exploitation. An MD4 hash is only 16 bytes, whereas three seven-byte DES keys brings us to a total of 21 bytes. How is the difference made up? Padding, which reduces the effective keylength of the third DES key to two bytes. If each of the two seven-byte and singular two-byte keys are cracked individually, the effective keyspace is reduced to $2^{56} + 2^{56} + 2^{16}$, or about 10^{17} —a decline in keyspace by a factor of 10^{32} . The interesting part about the MS-CHAP-v2 routine is that each DES key is used to encrypt the same plaintext; therefore, during the cracking process we can use the same key input for each of the two DES functions, effectively reducing the amount of DES functions for each iteration to one, leaving the only dual operation as the comparison of the DES output to each of the two previously determined ciphertexts. This is a final effective keyspace of 2^{56} (10^{16}), which in 1998 could be cracked in 4.5 days on average, and on the massively parallel architectures commercially available by the early 2010s, in about half a day on average.³⁵

Tools used to crack MS-CHAP-v2:

- **Airmon-ng:** To capture MS-CHAP-v2 keying material during the challenge-response process.
- **Aireplay-ng-0:** Used to deauthenticate a connected MS-CHAP-v2 client.
- **Chapcrack.py:** Used to parse/extract the keying material from MS-CHAP-v2 handshake.
- **Super Computer:** A hardware-accelerated cracking machine designed to crack a single DES routine (available through cloud services).

33.6 MITIGATING CONTROLS. This section provides enterprise-grade technical mitigating controls that should be evaluated in every wireless solution regardless of the vendor-specific implementation. These controls respond to the fundamental wireless-security issues discussed earlier and provide a layered defense as part of secure enterprise design.

33.6.1 Improvements. Wired access controls have remained relatively steady with the inception of 802.1X port-based authentication, which is now standard in most modern switches. 802.1X port-based authentication is fundamentally the same authentication function baked into the WPA security standard and combined with the *Temporal Key Integrity Protocol* (TKIP). WPA2 builds on WPA by adding even stronger encryption, which requires dedicated hardware for AES encryption. The 802.1X wired standard is fully implemented in WPA and WPA2 but reliant on the configuration as to whether it is leveraged, such as *WPA-Personal*, which uses a pre-shared key, or *WPA-Enterprise*, which uses a central authentication repository.

The use of wireless for endpoint connectivity can provide security, which rivals traditional wired architecture through the ease of administration and management of a thin-client AP environment and network choke-points for monitoring in close proximity to the initiating clients at the outermost edge of the network.

33.6.1.1 RSN. Complete RSN support in WPA2 and 802.11i draft support in the WPA standard is essential to ensure wireless confidentiality and integrity. To deploy RSN, any AP which can perform 802.1X authentication will suffice. Support for 802.1X on an AP is usually just a matter of understanding the EAPoL protocol and accepting a RADIUS server as a configuration parameter. The AP does not need to be configured

MITIGATING CONTROLS 33 · 35

for any particular EAP type, as this aspect of the network is completely transparent to the AP.

In a typical wireless environment, there are two fundamental management functions at play: user authentication management and key management for encryption, initial keying material, and key distribution. In a nonenterprise implementation, the user and key-management functions can be combined. The fact that a client possesses a correct pre-shared key authenticates them to the AP, and the bits which make up the pre-shared key may be used in the encryption process to derive a pairwise transient key for the encryption of the 802.11 payload.

33.6.1.2 Authentication and Access-Control Key-Management Considerations

- Solid authentication key management is necessary to deliver secure communication of keys in contrast to a pre-shared key environment where keys must be distributed through an out-of-band mechanism. A key can be in the form of a user/password combination, certificate, or other identifying material such an RSA SecurID number and PIN.
- Leverage an existing directory-based infrastructure used to authenticate users to computer systems on the wireless LAN or rollout an endpoint PKI infrastructure. The use of EAP and a backend user repository for authentication through the RADIUS protocol is the silver standard method of authentication. A PKI infrastructure, which provides client public/private keys to every wireless endpoint effectively eliminates the possibility for an attacker to capture a password hash challenge response.
- Strong EAP implementations include a TLS tunnel established between the endpoint client and backend authentication server (e.g., RADIUS) before any identifiable information is exchanged.

33.6.1.3 Encryption Key Management Considerations. As described in this chapter, encryption key management is the Achilles' heel of confidentiality. Modern encryption algorithms have matured to the point of perfection, as with AES, and are only as weak as their weakest link—that is, input. The best key-management functions are available only with the latest evolutions of wireless technologies. WPA2-Enterprise's AES/CCMP is the *de facto* standard and should be the minimum requirement for any new implementation. WEP, WPA, and WPA-2 present multiple methods of key management with improvement in each generation of 802.11 authentication/encryption schemes.

33.6.2 Authentication Server and Client Certificates. Certificate-based authentication between a client and backend authentication server is the *de facto* standard of wireless authentication. X.509 Certificates are used for the authentication of server to the client (RADIUS as a broker), client to the authentication server, or both.

To leverage endpoint-user and computer-based certificates for authentication in place of directory-based password authentication requires the implementation of a full *Public Key Infrastructure* (PKI). PKI is necessary to provide and regularly update/associate certificates with endpoint devices and sign them using an enterprise root Certificate Authority (CA) or outside root CA such as VeriSign.

33 · 36 802.11 WIRELESS LAN SECURITY

For details of PKI, see Chapter 37 in this *Handbook*.

The following are three of the most common uses for AP or client certificates within modern enterprise wireless implementations:

- Authentication Server Certificates Only:
 - A barebones RSN EAP implementation should require at a minimum a signed authentication server certificate. This ensures to connecting clients that they are performing authentication with a trusted AP and authentication server. In absence of a server certificate, clients are susceptible to MITM attacks, which can be used to capture EAP challenge-response hashes from the client. The RSN protocol ensures that a client must authenticate the authentication server before authenticating itself to the server. This sequence protects against RADIUS-impersonation attacks and rogue APs if configured properly.
 - Wireless implementations which leverage authentication-server certificates but do not enforce them at the supplicant level forgo any security provided by this control. An attacker can masquerade as the legitimate backend authentication server through the use of same-SSID and hijack the authentication process.
- **Client and server mutual certificate-based authentication:** Each endpoint client is provided a certificate signed by the enterprise root authority and installed in the Operating System's certificate store. EAP-TLS is commonly used with this type of PKI infrastructure to authenticate the server's certificate (as with directory-based authentication) and the client's certificate by validating the certificate's signature. This client certificate also contains field information to identify the user and any other applicable information for use by the authentication server.
- **Client- and server-based authentication with external client certificates:** A PKI infrastructure can provide even greater security in the form of "something I have that others don't." A client certificate (private key) can be stored on a smart card, requiring the user to possess the certificate separate from the operating system in the event the laptop is lost or stolen. The smart card must be inserted into the computer prior to attempting authentication.

33.6.3 Endpoint Supplicant Configurations. The endpoint-client supplicant configuration is a crucial and commonly overlooked component of a secure wireless deployment. Deploying a *secure* wireless network within the enterprise protects the internal network, but most deployments neglect to assess the security of the endpoints that connect to the new wireless network.

Enabling wireless connectivity at the enterprise levels means that laptop computers, which are taken out of the office, will almost always have their wireless radios on. A wireless radio is not an inherent risk, but when combined with default operating-system and driver configurations, such unprotected devices can easily provide all kinds of connectivity information to listening parties.

Most default wireless-client supplicants lend themselves to some sort of attack, which can potentially compromise the endpoint at a remote location, such as a coffee shop, and re-establish connectivity when the client goes to work and connects the infected machine to the enterprise wireless network.

Employees may be flying, sitting on a bus, or just grabbing something to eat—and if they sit in a physical location long enough, an attacker could glean a wealth of information from an insecure supplicant configuration. This information can range from a pre-shared key hash for a home network, remembered SSIDs (e.g., the

MITIGATING CONTROLS 33 · 37

enterprise wireless SSID), or even password hashes for enterprise directory-based authentication.

The following are recommendations when for securing an endpoint supplicant configuration:

- **Disable SSID Probe Requests:** Many operating systems and third-party supplicants will continuously probe for the last known SSIDs. This process is performed to ensure the client can connect to an AP, which may not be sending beacon frames to identify itself or if the client has missed the AP's beacon interval. A beacon interval is defined on the AP as the period of time between sending beacon frames to identify itself to all nonassociated clients within range. Disabling this feature prevents an attacker from passively monitoring 802.11 traffic and acquiring a list of SSIDs to which a particular client is attempting to connect. This information can then be leveraged to coerce the client into connecting to a rogue AP masquerading as one of the identified APs.³⁶
- **Wireless Profiles:** Wireless profiles should be configured to separate corporate AP profiles from hotspot or home AP profiles. Third-party supplicants commonly provide this capability to prevent the situation where a corporate employee uses their computer at a coffee shop and their wireless driver unknowingly broadcasts the corporate SSID, therefore inviting an attacker to masquerade and potentially capture a WPA/WPA2 pre-shared key hash or challenge response directory-based password hash. The use of different profiles limits the ability for the client to know how to connect to APs on the corporate network when not located at the office.
- **Secure Authentication Server Configuration**
 - **Certificate Authority:** A supplicant should only accept authentication server certificates, which were signed from a designated/single root CA. This situation prevents the ability to connect to an AP and underlying authentication server, which is displaying a valid/matching Common Name signed by a legitimate but different root CA.
 - **Certificate Validation:** Users should not have the ability to overwrite and enable the use of an untrusted server certificate during the authentication process. By default, many supplicants will prompt the user if an invalid certificate is detected and provide the option to ignore. The supplicant should be configured to not prompt the user and remove any possibility for successful authentication attempts to an AP with an invalid or untrusted certificate.
 - **Common Name (CN):** Ensure the backend authentication server's common name or CN is specified in the client. This will prevent a connection to an AP and authentication server which possess a trusted certificate (assuming a nonenterprise-specific root CA) but uses an incorrect hostname/FQDN.

33.6.4 Wireless Intrusion-Detection Systems. The wireless intrusion-detection system (WIDS) has come a long way to become a standard network-security monitoring appliance, becoming an everyday security-monitoring tool, like a traditional wired IDS. A WIDS is capable of monitoring the 802.11 airwaves at the data link and MAC layer regardless of encryption or authentication. Although a WIDS cannot be used in place of a wired IDS, it can provide visibility into a cloudy segment of the network edge. A WIDS is also not exclusive to enterprises that deploy a wireless network. A WIDS can also provide value for enterprises, which don't

33 · 38 802.11 WIRELESS LAN SECURITY

deploy wireless networks by ensuring against the unauthorized use of wireless connections by continuously monitoring the airwaves and alerting appropriately for internal rogue APs.

Multiple large breaches have stemmed from poor wireless management and visibility; this is a large problem for organizations with many physical retail or corporate outlets. Due to the widespread use and ease of installation of an SOHO AP obtained from a consumer retail outlet, it is becoming increasingly important to detect unauthorized extensions of the corporate network. In addition, the increased reliance on wireless communications for various business cases (such as corporate access, vendor devices requiring their own AP, or retail devices requiring dedicated APs with low-grade encryption due to low process power) has resulted in increased difficulty inventorying a network's authorized APs and confirming the proper segmentation or absence of internal connectivity for noncorporate access APs. In the event an enterprise can manage and inventory all known APs, there remains no native solution to quickly and accurately audit the implementation of enterprise standard authentication and encryption levels from a remote location.

33.6.4.1 Types of WIDS. Most modern WIDS infrastructures can be configured to leverage existing APs or use dedicated sensors to monitor 802.11 traffic. Dedicated sensors provide much more functionality at the expense of additional cost. A dedicated sensor has the time and CPU to solicit useful data from neighboring APs and clients, whereas an existing enterprise-thin AP may allocate only a fraction of processing toward these functions.

33.6.4.2 Use Cases. Due to the encryption, which takes place at the data link layer and above, a WIDS cannot inspect data-layer traffic. The majority of wireless threats involve the exploitation of a vulnerability, which exists in the data-link layer such as the 802.11 protocol itself or following authentication and encryption processes. The core value of a WIDS lies in the ability to provide assurance for wireless controls by verifying their existence and continuously monitoring for unauthorized changes.

A dedicated WIDS can provide the following services:

- Identify unauthorized APs across all physical locations (rogue APs).
- Detect and block the use of various wireless hacking tools with identifiable layer 1/2 signatures.
- Identify and continuously monitor for below enterprise-grade encryption or authentication mechanisms on authorized APs.
- Detect and block unauthorized connections to internal APs from clients with unauthorized wireless cards.
- Determine client supplicant configuration parameters as they relate to the transmitting of known SSID beacons.

33.6.4.3 Detailed WIDS Benefits

Regulatory Compliance. Regular wireless penetration testing and rogue AP audits are a necessary function of many regulatory certifications (such as the Payment Card Industry [PCI]). Leveraging a remote dedicated WIDS sensor may be a cost-effective alternative to deploying contractors or internal employees to each physical location to survey and assess controls.

MITIGATING CONTROLS 33 · 39

Rogue-AP Detection. Rogue APs are the number-one issue threatening wireless-network security.³⁷ Even in the face of 802.1X switchport access controls, Port Address Translation (PAT) can provide access to the wired network for multiple wireless clients by sharing a single authorized MAC address. A single unmanaged AP could provide unlimited access to the internal network for an extended period.

Methods of rogue AP detection:

- Corporate AP white-listing by:
 - Vendor: Rogue APs can be detected based on a nonstandard make/model by performing a lookup on the base clients OUI portion of the MAC address.
 - SSID: Rogue APs can be detected through the identification of transmitted SSID beacons containing a nonstandard enterprise SSID (such as “Linksys” or “Netgear”).
 - MAC Address: A static list of authorized AP MAC addresses can be obtained through SNMP or other means and provided to a WIDS system as “known good” APs.
 - Encryption/Authentication Grade: The WIDS can attempt to connect to potential rogue APs and enumerate available/requested encryption/authentication levels. A nonstandard level or method would be indicative of a rogue AP.
- Associated clients
 - A WIDS can perform analysis on a potential rogue AP’s connected clients to determine if clients are enterprise-provided based on MAC address OUIs.
- Signal strength
 - A WIDS is constantly monitoring the airwaves and can alert on anomaly detection of a new AP in close proximity to the physical sensor/location by determining the signal strength of the potential rogue AP measured in dBm. If the dBm strength falls within a user configured threshold, an alert can be triggered.

Confirming internal LAN connectivity:

- Internal Device Connectivity: The remote WIDS sensor can attempt to connect to the rogue wireless network and ping a known internal device to confirm connectivity to the enterprise network (assumed no authentication).
- Simple Network Management Protocol (SNMP) content-addressable memory (CAM) table³⁸ enumeration: Internal network switch CAM tables possess all MAC addresses a device knows about and which port they reside on. This detection technique is based on the fact that wireless APs possess two network cards. One network card resides on the wired side of the AP, and the other on the wireless side/radio. These cards both possess MAC addresses, which are one or two values apart. A WIDS sensor can identify the MAC address in use by the wireless radio through passive monitoring of the data link layer and search all known switch CAM tables via SNMP for a similar MAC address. The WIDS controller can be provided an SNMP read string to the enterprise-wide switch environment used by the network management team for various monitoring/polling processes.

Counter-Action. If a rogue AP is detected, the WIDS can be configured to disable the switch port it resides on through SNMP write strings or to trap an alert to the central

33 · 40 802.11 WIRELESS LAN SECURITY

console. A WIDS sensor can also be placed into a *blocking mode*, which will leverage spoofed 802.11 deauthentication packets to prevent clients from successfully establishing a connection to the rogue AP. Spoofed deauthentication packets were covered previously in Section 33.5.3, Management Frames, and in the correct circumstance can be used to *increase* security. These actions should be carefully analyzed prior to execution to minimize the probability of unknowingly attacking a legitimate neighboring business or household.

Defense in Depth—Rudimentary Controls. The previously described controls are best practice that can make or break the security of an enterprise wireless network. In addition to these controls, there are some controls which by themselves do not afford much protection, but together can reduce an attacker's success rate.

Security through obscurity may be a bad idea if it's the only security measure, but combined with solid security controls, one can potentially reduce legitimate attacks against core security control (such as EAP-MS-CHAP-v2). These controls are what we refer to as *script-kiddie* controls; they offer little resistance to an experienced attacker but cause an adversary to take multiple steps before even exposing their target objective.

- **MAC-address based restrictions:** Depending on the capabilities of each vendor-specific AP, it may be possible to restrict association to an AP based on MAC address, even in an enterprise environment. Methods exist to dynamically update MAC address mappings within a backend repository, similar to wired MAC-based VLAN assignment.
- **SSID Cloaking:** SSID cloaking is a perfect example of security through obscurity. The presence of an AP and BSSID can be identified through various penetration testing tools capable of analyzing raw 802.11 traffic, but an everyday operating-system default or third-party client supplicant will not display the presence of an AP with nonbroadcasting SSID.
- **Client Isolation:** Client isolation is an AP-specific control that monitors the destination MAC address of the incoming traffic. If traffic is destined for another client on the same AP, or SSID in thin client environment, it will be denied access. The purpose of this control is to limit the ability of a malicious user to communicate with other devices connected to the AP. This feature is common on public APs where the end user is not necessarily well vetted or authenticated. Depending on the enterprise use for a particular SSID/VLAN, this feature may be enabled to prevent access to other endpoints in an environment where services are provided over wireless by other wireless clients, such as a printer; analysis should be performed to allow only a subset of communication with those devices.

33.7 SECURE ENTERPRISE DESIGN. Earlier we discussed fundamental and specific vulnerabilities in detail to expose the underlying inherent risk of wireless network mediums and provide a background on some of the latest and greatest threats that should be addressed. Understanding an organization's risks and associated threats as they pertain to their unique network environment is vital to ensuring they are properly managed. The proper management of technical security risk starts with a secure design and ends with the implementation of secure technical controls to further support a secure design. The phrase "crunchy on the outside, soft and chewy on the inside" plays directly into this particular type of governance. Mitigating controls can

SECURE ENTERPRISE DESIGN 33 · 41

provide a crunchy perimeter, but secure design prevents a chewy interior in the event of perimeter-control failure.

In this section, we discuss priority design considerations which foster a secure enterprise wireless network. These design considerations should be accompanied by the technical mitigating controls discussed earlier to provide a wholesome approach. Integrating and allowing for various types of technical mitigating controls such as EAP methods, SSIDs, PKI infrastructures, AP certificates, or supplicant configurations during the initial wireless design phase can increase their effectiveness as they function as a whole in contrast to after-the-fact add-ons.

33.7.1 Benefits of Wireless Controller Architecture. Enterprise wireless network hardware comes in two primary varieties, thin (lightweight) and thick (autonomous) APs. These two types of designs will have a large impact on the way the wireless network is operated and managed, and the extent of available controls to secure it.

- A thick AP possesses a lot of intelligence and features outside of the core 802.11 or RSN functions. These features include helper services such as DHCP, SNMP, QOS, Firewalls, etc. Although these types of services definitely have a place, from a management standpoint it can become difficult to synchronize configurations across thick clients or varying vendors.
- Thin client APs are bare-bone APs, which implement the necessary 802.11/RSN functions only. Enterprises should leverage thin-client APs to lower cost and increase secure management of the wireless network. These APs leverage connected switch hardware or dedicated wireless controllers for intelligence and configuration settings. Upon boot and initialization a thin client will discover, contact, and download the latest configuration settings.

The security afforded by a thin-client environment reaches beyond ease of device management by providing the following additional benefits:

- **Physical safeguard of configuration data:** A thick client's configuration could be compromised through physical theft and analysis. A thin client holds the bare bone configuration to allow 802.11 connectivity, RADIUS shared secrets, and SNMP strings; other stored configuration information is held on the backend switch or wireless controller
- **Mesh Security Functions:** A common form of communication across APs such as the *Lightweight Access Point Protocol* (LWAPP)³⁹ enables widespread communication between APs and controllers across an entire enterprise. This communication mechanism combined with the physical presence of APs across primary and remote locations enables the ability to leverage them for security functions. Leveraging client APs for security monitoring can be beneficial from a cost-saving standpoint in favor of a dedicated WIDS, albeit with a lack of functionality.
 - Detect rogue APs using AP infrastructure rogue-discovery techniques. Vendor-specific solutions exist to leverage client-servicing APs to momentarily scan for rogue AP beacons, which are not part of the LWAPP network.
 - An AP can be placed into a detection role and connected to a switched port analyzer (SPAN) port to track down potential rogue AP beacons for internal network connectivity.⁴⁰

33 · 42 802.11 WIRELESS LAN SECURITY

33.7.2 Network Segmentation. A basic wireless network should be segmented in two ways. First, the wireless network should be segmented from the internal LAN network, and second, wireless networks should be segmented between wireless networks with different business uses. In addition, during the segmentation process consideration should be provided for WIDS sensor placement with the ability to physically span all areas of the wired network.

33.7.2.1 Wired versus Wireless Endpoint Segmentation. Wired networks commonly assign VLANs on a port or MAC address basis. In the case of port-based VLAN assignment, management relies heavily on proper access-layer switch configuration, trunk-port configuration for uplinks, and pruning the propagation of sensitive VLANs for different physical locations. Operational changes may dictate the updating of numerous switches for relatively minor moves. In the event of wired MAC address based VLAN assignment, MAC-authentication-aware switch hardware and backend RADIUS servers for MAC->VLAN mappings must be used. MAC-based VLAN assignment can improve management capability and oversight but poses a significant security risk as they are susceptible to MAC-address spoofing, which may allow an attacker access to sensitive VLANs by spoofing an authenticated device's MAC address.

In a wireless environment, VLANs can be assigned on a per-SSID or per-user-group basis, and administration efforts to provision a user to a VLAN can be greatly reduced or completely removed. A user is no longer confined to particular physical location to gain access to a required VLAN, as the necessary SSID may be available anywhere on campus. In the event of expansion, deploying a new thin AP in close proximate will deliver the necessary VLANs, which a client supplicant profile is configured for. Since VLAN assignment can be controlled through minor client supplicant configurations and backend RADIUS privileges, segmentation can be determined during the asset provisioning process and require relatively minor management thereafter.

33.7.2.2 Segmentation from the Wired LAN. Due to public accessibility of wireless airwave transmission (encrypted or not), segmentation within the existing LAN should be addressed to mitigate the risk of complete internal access due to a lapse of wireless security. A solid network security posture assumes a particularly risky segment may be compromised and is designed to mitigate the damage from such a disaster.

The value in segmenting the wireless LAN from the wired LAN thrives on the ability to hedge against the following events, which could lead to LAN compromise from a remote location, such as the corporate parking lot.

- Malicious use of compromised employee username/password-based authentication credentials through spear phishing emails or other theft.
- Misconfiguration of wireless security controls, such as a weak EAP implementation (EAP-LEAP or absence of AP certificate verification in EAP-PEAP-MS-CHAP-v2) lending itself to dictionary or RADIUS impersonation attacks.
- Rogue AP “hotspotting” of employee laptops to establish connectivity and compromise the endpoint through exploitation of services or HTTP traffic MITM. Once the endpoint is connected to the legitimate enterprise wireless an attacker could gain access to the LAN.

SECURE ENTERPRISE DESIGN 33 · 43

The need for segmentation of wireless traffic is due to the inherent lack of physical controls for 802.11 traffic and range of accessibility. An external attack at a DMZ or other Internet-accessible device could lead to the compromise of DMZ assets, but controls typically exist in the form of multileg or sandwich-style firewall environments which limit the exposure from the DMZ network and prevent direct LAN connectivity. This type of segmentation should be applied to the wireless environment for similar reasons, albeit different attack vectors.

33.7.2.3 VPN over Wireless. Depending on the purpose of the wireless network, users may require the same access on the wireless network that is afforded on the wired LAN. In favor of exposing the wireless network to the complete wired LAN a VPN solution may be used. This architecture would allow wireless network access to the Internet, just the VPN gateway, or other basic internal assets (such as an intranet site), which force the use of a VPN gateway for any other resources.

An additional strength provided by a VPN on top of wireless authentication is defense in depth. To gain internal LAN access from a wireless network, an attacker must compromise and gain access to the VPN in addition to any layer 2 wireless authentication/access control mechanisms. An attacker who gains access to the wireless medium would likely never suspect the use of a VPN and assume the state of limited internal access or mistake the current state for the complete LAN.

An assessment should be performed to understand the business case where wireless assets need complete access to the corporate LAN, and explicitly allow access on an IP/port basis. Fine-grain segmentation can be obtained by leveraging SSID-based VLAN assignments, explained in detail in the following sections, to assign restrictions on a per-VLAN basis with supporting ACLs at higher points within the network.

33.7.2.4 Inter-Wireless Network Segmentation. Inter-wireless segmentation enables different user types to connect to different SSIDs provided by the same physical AP and be afforded different network-based access privileges. The process of segmenting each SSID into its own network usually involves VLAN tagging at the AP level to distinguish layer-2 traffic by forcing inter-VLAN communication to flow over a firewall or upstream router with VLAN/IP-based ACLs. Backend directory systems such as LDAP can be leveraged by modern APs to provide a decision on the VLAN assignment a particular user will receive—regardless of the SSID on which they currently reside. Dynamic backend directory-based systems are not completely secure due to the reliance on MAC addresses as a unique identifying agent.

This type of segmentation can be enforced at the network layer, through aforementioned VLAN firewalling, or through VPN access lists if employee wireless access is forced through a VPN tunnel to communicate with LAN assets.

33.7.2.5 Benefits and Drawbacks of SSID/VLAN-Based Segmentation. Benefits of SSID/VLAN-based segmentation include:

- Differential treatment for wireless users by group or SSID; for example,
- An SSID of “corporate-sales” and “corporate-hr” may provide varying levels of access. It would be best practice to disguise and cloak the SSID names to obfuscate their purpose to an outsider.
- Ease of network administration by removing the need to assign VLANs on a per-port basis. SSID-based VLAN assignment only requires endpoint configuration

33 · 44 802.11 WIRELESS LAN SECURITY

to ensure the endpoint provides the credentials necessary to authenticate to a particular SSID.

A drawback of SSID-/VLAN-based segmentation is increased wireless-network overhead and reduced performance due to multiple SSID memory requirements.

Summary of Wireless Segmentation Controls

- LAN segmentation: protects internal LAN in the event of wireless compromise through VLAN and wireless IP-based firewall ACLs.
 - VPNs can be used over a wireless medium to allow users to access LAN assets.
- Intra-wireless network segmentation:
 - SSID based: The use of separate SSIDs for different business cases can easily be mapped to a VLAN and provide varying levels of segmentation.
 - Group based: RADIUS-based VLAN assignment for wireless user's current group privileges.

33.7.3 User Segmentation. A RADIUS server can make advanced authorization decisions for a valid username/password match. Leveraging designated user groups in the authorization process can provide fine-grained control over which users are able to use a wireless network segmented in large sweeps. If it's known that domain admins, power users, or general service accounts should never be able to authenticate over a wireless connection, they can be implicitly denied through nonmembership in a wireless-users group. Any user group that requires the ability to authenticate over wireless should be a member of this group.

A default RSN implementation that can and will authorize any user within the provided directory can become a victim of attacks on default, inactive, guest, or test accounts within the designated directory. Test accounts may be vulnerable to weak canonical passwords (e.g., “test” or “password”) due to their creation during early stages of software or directory development or through other means that bypass corporate password policy. When a username/password authentication mechanism is in play, an attacker may attempt basic brute-force attacks to gain access to the network through standard test or guest accounts. Many organizations don't feel this is a serious problem, but when a directory contains thousands of accounts, and organizations don't audit on a regular basis, it's easy for things to slip through the cracks.

The use of a designated wireless group:

- Decreases the surface area of an attack, even if an adversary is able to obtain a list of user accounts and leverage them to brute force a wireless network through the use of basic passwords.
- Decreases the number of endpoints with saved or remembered network profiles that can be attacked and used to capture username/password challenge hashes in a hotspot or RADIUS impersonation-style attack.

33.8 SECURITY AUDITING TOOLS.⁴¹ This section looks at the open source and commercial tools available to help audit a wireless environment and to understand the tools that war drivers have to attack a network. Exhibit 33.11 lists the main programs available, including older software still mentioned in discussions, and the platforms they run (or ran) under.

SECURITY AUDITING TOOLS 33 · 45**EXHIBIT 33.11 802.11 Security-Auditing Software**

Tool	Description and Home Page	Operating System
Aircrack	A program that breaks WEP encryption, injects packets (Aireplay), and performs a dictionary attack against WPA PSK. http://www.aircrack-ng.org	Linux Windows
Auditor	An obsolete security collection that had all the required drivers and security tools to audit or hack wireless networks. Replaced by BackTrack (see below).	Linux
AirMagnet	A commercial wireless LAN analyzer that runs on laptop PCs and pocket PCs. http://www.flukenetworks.com/enterprise-network/wlan-security-and-analysis	Windows Pocket PC
Airsnort	Outdated tool. Broke WEP encryption using the RC4 Key Scheduling Vulnerability and Korek optimizations. Last update was 2005. http://sourceforge.net/projects/airsnort/	Linux
Airopeek NX	A commercial wireless LAN analyzer that ran on laptop and tablet PCs. Replaced by OmniPeek (see below).	Windows
BackTrack	A Linux security distribution and is the result of the merging of the Whax and Auditor security collections. http://www.backtrack-linux.org/ http://www.remote-exploit.org/articles/backtrack/index.html	Linux
CoWPAtty	A WPA PSK cracker last updated in 2008. http://www.willhackforsushi.com/Cowpatty.html	Linux
Ethereal	A program that understood the format of packets and dumps out their contents. Could capture packets as well as process a log file created by other programs. Replaced by Wireshark (see below) in 2006 when original programmer hired by a new company.	Linux Windows
Kismet	Wireless auditing tool still being actively developed and improved. Latest updates at time of writing were on 2012-12-06. http://www.kismetwireless.net/	Linux, FreeBSD, NetBSD, OpenBSD, Mac OS X
Linux-wlan	The Linux drivers for Intersil Prism 2, 2.5, and 3-based cards used by open source software such as Kismet and Wellenreiter. Last updated 2004. http://www.linux-wlan.com/linux-wlan/	GNU/Linux
Sniffer Global Analyzer	High-end software for use with portable (hardware) analyzer tools. http://www.netscout.com/products/enterprise/Pages/default.aspx	Hardware modules & client software
Netstumbler	Detects broadcasting 802.11 networks. Netstumbler is not suitable as a network auditing tool because access points can be configured to avoid detection. Last update (v0.4.0) was in 2004. http://stumbler.net/	Windows

(continued)

33 · 46 802.11 WIRELESS LAN SECURITY

EXHIBIT 33.11 (Continued)

Tool	Description and Home Page	Operating System
OmniPeek	Commercial, sophisticated network-analysis tool from WildPackets with wireless capabilities. Replaced Airopeek. http://www.wildpackets.com/products/omnipeek_network_analyzer/wireless_network_analysis	Mac, Windows, hardware unit
Wellenreiter	Provides similar functionality to Netstumbler in that it detects access points. However, Wellenreiter has a major advantage over Netstumbler in that it detects access points that Netstumbler cannot see. At time of writing, the last update was 2012-07-17. http://sourceforge.net/projects/wellenreiter/	Linux
Windows wireless client	A built-in wireless subsystem that displays a list of available networks.	Windows
Wireshark	Highly regarded network protocol analyzer. Version 1.8.4 released 2012-11-28. http://www.wireshark.org/download.html	Windows, Mac

33.9 CONCLUDING REMARKS. The convenience and increasing bandwidth of wireless networks will inevitably result in continued growth of both technological improvements and implementation. Network and security managers must continue to monitor these developments and take appropriate action to secure all of their systems against constantly evolving threats to information security.

33.10 ABBREVIATIONS AND DEFINITIONS. Most of these abbreviations and definitions originate from the ANSI/IEEE 802.11-2012 standard.

Abbreviation/ Term	Meaning/Definition
802.11	The ANSI/IEEE standard 802.11, 1999 edition. Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications.
802.1X	IEEE Std 802.1X-2001 defines port-based network access controls. It provides the means to authenticate and authorize a network device before granting access to network resources. ⁴²
Access Controller	The CAPWAP term for a wireless switch that controls and manages Lightweight Access Points (LWAP).
Access point (AP)	A specialized wireless client that bridges other wireless clients to the physical wired network or distribution system.
Authenticator	An 802.1X component that authenticates a network device before allowing it to access network resources.
Authentication Server (AS)	An authentication server is an 802.1X component that performs the actual authentication of the supplicant on behalf of the authenticator.

ABBREVIATIONS AND DEFINITIONS 33 · 47

Abbreviation/ Term	Meaning/Definition
BSS	Basic Service Set. A set of clients controlled by a single coordination function (an AP).
BSSID	Basic Service Set Identification. Due to the ability for a physical AP to provide multiple SSIDs on one interface, a BSSID reference to the MAC address associated with that interface to uniquely identify it as the provider of the SSIDs.
CAPWAP	Control and Provisioning of Wireless Access Points. An IETF initiative to establish standards enabling interoperability between different vendors APs and wireless switches.
CBC-MAC	Cipher Block Chain–Message Authentication Code. A message authentication algorithm based on a symmetric block cipher, where each block chains into the next data block.
CCMP	Counter mode with CBC-MAC protocol. A symmetric block cipher mode providing data privacy using counter mode, and data origin authentication using CBC-MAC. This is the protocol required for compliance with the WPA2 standard.
CRC	Cyclic redundancy check.
CRC-32	A 32 bit CRC that is used in the 802.11 standard.
Distribution System (DS)	A system used to interconnect a set of basic service sets (BSS) to create an extended service set (ESS).
DMZ	Demilitarized zone. This refers to a network segment that separates and protects a trusted network from an untrusted one.
DSSS	Direct sequence spread spectrum. DSSS is an 802.11 modulation technique that spreads signal power over a wide band of frequencies.
DTLS	Datagram Transport Layer Security. Essentially TLS over UDP.
EAP	Extensible authentication protocol defined in RFC 5247. ⁴³
EAP-PEAP	Protected EAP. A draft IETF standard by RSA, Microsoft, and CISCO for doing EAP authentication protected by TLS. ⁴⁴
EAP-TLS	EAP–Transport Layer Security defined in RFC 2716. ¹³
EAP-TTLS	EAP—Tunneled TLS authentication protocol. A draft IETF standard that tunnels EAP authentication through TLS. ⁴⁵
EAPoL	EAP over LANs. The protocol defined in 802.1X for encapsulating EAP messages exchanged between a wireless client and the AP.
ESS	Extended Service Set. A set of one or more interconnected basic service sets (BSSs) and integrated LANs that appears as a single BSS to wireless clients.
Exclusive-Or (XOR)	An add-without-carry logical operation. XOR is used in cryptographic algorithms partly because it is reversible by repeating an operation. A XOR B XOR A will result in B; A XOR B XOR B will result in A.
FHSS	Frequency hopping spread spectrum. A transmission method that sends bursts of data over a number of frequencies, with the system hopping between frequencies in a defined way.

33 · 48 802.11 WIRELESS LAN SECURITY

Abbreviation/ Term	Meaning/Definition
IBSS	Independent Basic Service Set. An ad hoc wireless network where clients communicate directly with one another rather than via an AP.
ICV	Integrity check value. A CRC-32 value used by the WEP algorithm to detect changes in the ciphertext and to ensure that the packet has decrypted correctly.
IV	Initialization vector. A 24-bit value prepended to the WEP key, which is used as the seed for the RC4 stream cipher as part of the WEP algorithm.
KSA	The key scheduling algorithm of the RC4 stream cipher.
LAN	An IEEE 802 local area network.
LWAP	Lightweight Access Point. A type of AP that is essentially an 802.11 wireless radio that passes the received packets to a wireless switch for processing.
LWAPP	Lightweight Access Point Protocol. The base protocol chosen for the CAPWAP protocol and used by wireless switches to manage and communicate with access points.
MAC	Media access control. This is the data link layer in a wireless LAN that enables clients to share a common transmission medium.
Mb/s	Megabits per second.
MIC	Message integrity code. A cryptographic integrity code to detect changes in a message.
Michael	The message integrity code (MIC) for TKIP.
MITM	Man-in-the-middle. An MITM attack is one where an attacker can sit between two communicating parties and monitor what is being exchanged, without being noticed.
MPDU	MAC protocol data unit. The packets/frames exchanged between clients via the physical layer.
MSDU	MAC service data unit. Information that is to be delivered as a unit between MAC service APs (SAP). If an MSDU is too large to fit into one MPDU then it will be fragmented into multiple MPDUs and reassembled into one MSDU by the receiver.
OFDM	Orthogonal frequency division multiplexing. A modulation technique used in the 802.11a and 802.11g standards that works by dividing data into several pieces and simultaneously sending the pieces on many different subchannels. This mechanism enables throughput of up to 54Mb/s.
PAE	Port Access Entity. An 802.1X object that operates the authentication mechanism in the participants.
PKI	Public Key Infrastructure.
PMK	Pairwise master key. The PMK is a master session key derived from the overall master key. For 802.1X authentication mechanisms (e.g., EAP-TLS), the key is generated during authentication. For Pre-Shared Key authentication, the PMK is the pre-shared key (PSK).

ABBREVIATIONS AND DEFINITIONS 33 · 49

Abbreviation/ Term	Meaning/Definition
Port	Network access port.
PRF	Pseudorandom function. An RSNA-defined algorithm that is part of the key-generation function.
PSK	Pre-shared key. A static cryptographic key that is distributed out-of-band to all clients in the wireless network.
PTK	Pairwise Transient Key. An RSNA cryptographic key that is the source of working temporal keys that protect wireless messages.
RADIUS	Remote authentication dial-in user service, defined in RFC 2865. ⁴⁶
RSN	Robust Security Network. The name for the new security system implemented in the RSN draft standard. An RSN is one that only allows the creation of RSN Associations.
RSNA	A Robust Security Network Association is where two clients have authenticated each other and associated using the 4-way handshake protocol.
RSN IE	Robust Security Network Information Element. The RSN IE is contained in Beacon and probe frames, and lists the supported authentication mechanisms and cipher suites.
Supplicant	The supplicant is the 802.1X term for the network device that wants to connect to the network.
“The Standard”	In this chapter, the term refers to the ANSI/IEEE standard 802.11-1999 edition.
TKIP	Temporal Key Integrity Protocol. TKIP is a suite of algorithms enhancing the WEP protocol to provide secure operations for legacy 802.11 equipment.
TLS	Transport Layer Security defined in RFC 2246.
TSN	Transient security network. A TSN can form associations with both RSNA capable and pre-RSNA clients.
UDP	User Datagram Protocol. An Internet protocol used to send short messages to other systems. UDP does not provide any reliability or ordering guarantees.
WEP	Wired Equivalent Privacy. The protocol defined in the 802.11-1999 standard to encrypt data on a wireless LAN.
Wireless LAN	In this chapter, the term “wireless LAN” refers to 802.11 wireless local area networks using a radio physical layer.
WLAN	Wireless local area network.
WTP	Wireless termination point. The CAPWAP term for an access point
WPA	Wi-Fi Protected Access. WPA was the first 802.11 wireless security standard from the Wi-Fi Alliance, using the TKIP security mechanism in the RSN standard.
WPA2	Wi-Fi Protected Access v2. WPA2 is the Wi-Fi Alliance’s name for the AES-based CCMP wireless security mechanisms in the RSN standard.
XOR	Exclusive-OR.

33 · 50 802.11 WIRELESS LAN SECURITY

33.11 FURTHER READING

- Anjum, F., and P. Mouchtaris. *Security for Wireless Ad Hoc Networks*. Wiley, 2007.
- Edney, J., and W. A. Arbaugh. *Real 802.11 Security: Wi-Fi Protected Access and 802.11i*. Addison-Wesley, 2003.
- Coleman, David D., and David A. Westcott. *CWNA: Certified Wireless Network Administrator Official Study Guide: Exam PW0-105*, 3rd ed. Sybex, 2012.
- Gast, M. *802.11 Wireless Networks: The Definitive Guide*, 2nd ed. O'Reilly & Associates, 2005.
- Geier, Jim. *Designing and Deploying 802.11n Wireless Networks*. Cisco Press, 2010.
- Makki, S., P. Reiher, K. Makki, N. Pissinou, and S. Makki, eds. *Mobile and Wireless Network Security and Privacy*. Springer, 2007.
- Vacca, J. R. *Guide to Wireless Network Security*. Springer, 2006.
- Wrightson, Tyler. *Wireless Network Security A Beginner's Guide*. McGraw-Hill Osborne Media, 2012.
- Yang, Xiao, Xuemin Shen, and Ding-Zhu Du, eds. *Wireless Network Security*. Springer, 2007.

33.12 NOTES

1. IEEE LAN/MAN Standards Committee, “IEEE 802.11™: Wireless LANs,” IEEE Standards Association, 2013, <http://standards.ieee.org/about/get/802/802.11.html>
2. David Chernicoff, “BYTE Guide To In-Flight Wi-Fi,” *Information Week: Byte Reviews*, December 14, 2012, www.informationweek.com/byte/personal-tech/wireless/byte-guide-to-in-flight-wi-fi/240144399
3. IEEE LAN/MAN Standards Committee, “IEEE Std 802.11-2012 Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications,” IEEE Computer Society, March 29, 2012, <http://standards.ieee.org/getieee802/download/802.11-2012.pdf>
4. See also Stephen McQuerry, “Understanding WLAN Security,” Chap. 3 in *Interconnecting Cisco Network Devices, Part 1 (ICND1)*: CCNA Exam 640-802 and ICND1 Exam 640-822, 528, Cisco Press, 2007; an extract is available at [www.ciscopress.com/articles/article.asp?p=1156068&seqNum=3](http://ciscopress.com/articles/article.asp?p=1156068&seqNum=3)
5. Microsoft, “Shared Key Authentication,” MSDN—Microsoft Developer Network, April 8, 2010, <http://msdn.microsoft.com/en-us/library/aa916565.aspx>
6. Sheila Frankel, Bernard Eydt, Les Owens, and Karen Scarfone, “Establishing Wireless Robust Security Networks: A Guide to IEEE 802.11i (SP 800-97),” National Institute of Standards and Technology Special Publications, February 2007, <http://csrc.nist.gov/publications/nistpubs/800-97/SP800-97.pdf>
7. Eric Griffith, “WEP: Cracked in 60 Seconds,” Wi-Fi Planet, April 9, 2007, www.wifiplanet.com/news/article.php/3670601
8. Margaret Rouse, “TKIP (Temporal Key Integrity Protocol),” SearchMobileComputing, March 2006, <http://searchmobilecomputing.techtarget.com/definition/TKIP>
9. Margaret Rouse, “CCMP (Counter Mode with Cipher Block Chaining Message Authentication Code Protocol),” SearchMobileComputing, June 2008, <http://searchmobilecomputing.techtarget.com/definition/CCMP>

NOTES 33 · 51

10. IEEE LAN/MAN Standards Committee, “802.1X-2010—Port Based Network Access Control,” IEEE 802.1 Working Group, February 5, 2010, <http://standards.ieee.org/getieee802/download/802.1X-2010.pdf>
11. B. Aboba, L. Blunk, J. Vollbrecht, J. Carlson, and H. Levkowetz, “Extensible Authentication Protocol (EAP): Request for Comments 3748,” Internet Engineering Task Force, June 2004, <http://tools.ietf.org/html/rfc3748>
12. VOCAL Technologies, “EAPoL—Extensible Authentication Protocol over LAN,” VOCALSecure Communication, 2012, www.vocal.com/secure-communication/eapol-extensible-authentication-protocol-over-lan
13. D. Simon, B. Aboba, and R. Hurst. “The EAP-TLS Authentication Protocol: Request for Comment 5216,” Internet Engineering Task Force, March 2008, <http://tools.ietf.org/html/rfc5216>
14. P. Funk and S. Blake-Wilson, “Extensible Authentication Protocol Tunneled Transport Layer Security Authenticated Protocol Version 0 (EAP-TTLSv0): Request for Comments 5281,” Internet Engineering Task Force, August 2008, <http://tools.ietf.org/html/rfc5281>
15. Matthew Gast, “TTLS and PEAP Comparison,” Opus One, April 13, 2004, www.opus1.com/www/whitepapers/ttlsandpeap.pdf
16. Microsoft, “PEAP-MS-CHAP v2,” Microsoft TechNet Windows Server, March 31, 2005, [http://technet.microsoft.com/en-us/library/cc779326\(v=ws.10\).aspx](http://technet.microsoft.com/en-us/library/cc779326(v=ws.10).aspx)
17. Moxie Marlinspike, “Divide and Conquer: Cracking MS-CHAPv2 with a 100% Success Rate,” CloudCracker: Blog, July 29, 2012, <https://www.cloudcracker.com/blog/2012/07/29/cracking-ms-chap-v2>
18. Cisco, “PEAP/EAP-TLS Configuration Scenario,” Configuration Guide for Cisco Secure ACS 4.1, n.d., accessed January 10, 2013, www.cisco.com/en/US/docs/net_mgmt/cisco_secure_access_control_server_for_windows/4.1/configuration/guide/peap_tls.html
19. Krishna Sankar, Andrew Balinsky, Darrin Miller, and Sri Sundaralingam, “EAP Authentication Protocols for WLANs,” Cisco Press, February 18, 2005, www.ciscopress.com/articles/article.asp?p=369223&seqNum=4
20. Sankar et al., “EAP Authentication Protocols for WLANs,” p. 5, www.ciscopress.com/articles/article.asp?p=369223&seqNum=5
21. INTELLIGRAPHICS Device Drivers, “Introduction to IEEE 802.11,” INTELLIGRAPHICS | Developer Resources | White Papers, 2012, www.intelligraphics.com/introduction-ieee-80211
22. IEEE LAN/MAN Standards Committee, “Part 11: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications—Amendment 2: Fast Basic Service Set (BSS) Transition,” IEEE Standards Get IEEE 802 Program, July 15, 2008, <http://standards.ieee.org/getieee802/download/802.11r-2008.pdf>
23. Niels Ferguson, “Michael: An Improved MIC for 802.11 WEP | document IEEE 802.11-02/020r0,” IEEE 802.11 Documents | Documents | Wiki, January 2002, <https://mentor.ieee.org/802.11/dcn/02/11-02-0020-00-000i-michael-an-improved-mic-for-802-11-wep.doc>
24. Jim Geier, “Understanding 802.11 Frame Types,” Wi-Fi Planet | Tutorials. August 15, 2002, www.wi-fiplanet.com/tutorials/article.php/1447501
25. Aircrack-ng, Home Page v1.1, April 24, 2010, www.aircrack-ng.org

33 · 52 802.11 WIRELESS LAN SECURITY

26. Paul Rubens, “Secure Your WLAN With Aircrack-ng,” Wi-Fi Planet | Wireless Security, December 27, 2007, www.wi-fiplanet.com/tutorials/article.php/3718671/Secure-Your-WLAN-With-Aircrack-ng.htm
27. Dino A. Dai Zovi and Shane A. Macaulay, “Attacking Automatic Wireless Network Selection,” Dino A. Dai Zovi’s Website, March 18, 2005, www.theta44.org/karma/aawns.pdf
28. Lisa Phifer, “The Caffe Latte Attack: How It Works—and How To Block It,” Wi-Fi Planet | Tutorials, December 12, 2007, www.wi-fiplanet.com/tutorials/article.php/10724_3716241_2
29. Martin Beck and Erik Tews, “Practical Attacks against WEP and WPA,” aircrack-ng Website, November 8, 2008, <http://dl.aircrack-ng.org/breakingwepandwpa.pdf>
30. “mister_x,” “Airmon-ng,” Aircrack-ng Website, October 31, 2010, www.aircrack-ng.org/doku.php?id=airmon-ng
31. Beck and Tews, “Practical Attacks against WEP and WPA.”
32. Stefan Viehböck, “Brute Forcing Wi-Fi Protected Setup: When Poor Design Meets Poor Implementation, v3,” .braindump—RE and stuff, December 26, 2011, http://sviehb.files.wordpress.com/2011/12/viehboeck_wps.pdf
33. reaver-wps, “reaver-wps: Brute Force Attack against Wifi Protected Setup,” Google Code Project Hosting, n.d., accessed January 10, 2013, <http://code.google.com/p/reaver-wps>
34. Admin, “How to Hack WPA WiFi Passwords by Cracking the WPS PIN,” Live Technology Guide | Null Byte | Linux Security, 2013, <http://livetechnoguide.com/how-to-hack-wpa-wifi-passwords-by-cracking-the-wps-pin> (URL inactive).
35. Moxie Marlinspike, “Divide and Conquer.”
36. Microsoft, “Description of the Wireless Client Update for Windows XP with Service Pack 2: Article ID 917021 Rev. 8,” Microsoft Support, October 9, 2011, <http://support.microsoft.com/kb/917021>
37. Karen J. Bannan, “Rogue Mobile Devices Threaten Enterprise Security,” DefenseSystems Website, 2011, <http://defensesystems.com/microsites/2011/mobile-wireless/personal-devices-connect-corporate-infrastructures.aspx>
38. T. S. Adams, “What Is a CAM Table?” wiseGeek Website, 2013, www.wisegeek.com/what-is-a-cam-table.htm
39. P. Calhoun et al., “Lightweight Access Point Protocol: Request for Comments 5412,” Internet Engineering Task Force, February 2010, <http://tools.ietf.org/html/rfc5412>
40. Cisco, “Rogue AP Detection under Unified Wireless Networks—Document ID 70987,” Cisco | Wireless, LAN (WLAN), September 25, 2007, www.cisco.com/en/US/tech/tk722/tk809/technologies_white_paper09186a0080722d8_c.shtml
41. This section is an updated extract of a table from Chapter 33 (“802.11 Wireless LAN Security”) by Gary L. Tagg in the 5th edition of this *Handbook*.
42. IEEE LAN/MAN Standards Committee, “802.1X-2010—Port Based Network Access Control.”
43. Aboba et al., “Extensible Authentication Protocol (EAP): Request for Comments 3748.”
44. A. Palekar, Dan Simon, Glen Zorn, and S. Josefsson, “Protected EAP Protocol (PEAP)—Internet Draft,” Internet Engineering Task Force, March 22, 2003, <http://tools.ietf.org/html/draft-josefsson-pppext-eap-tls-eap-06>

NOTES 33 · 53

45. Funk and Blake-Wilson, "Extensible Authentication Protocol Tunneled Transport Layer Security Authenticated Protocol Version 0 (EAP-TTLSv0): Request for Comments 5281."
46. C. Rigney, S. Willens, A. Rubens, and W. Simpson, "Remote Authentication Dial In User Service (RADIUS): Request for Comments 2865," Internet Engineering Task Force, June 2000, <http://tools.ietf.org/html/rfc2865>

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 34

SECURING VoIP

Christopher Dantos and John Mason

34.1 INTRODUCTION	34·1		
34.2 REGULATORY COMPLIANCE	34·2	34.4.3 Logical Separation of Voice and Data	34·7
34.2.1 Enhanced 911	34·2	34.4.4 Quality of Service	34·7
34.2.2 Communications Assistance for Law Enforcement Act (CALEA)	34·3	34.4.5 Network Monitoring Tools	34·7
34.2.3 State Laws and Regulations	34·3	34.4.6 Device Authentication	34·8
34.2.4 International Laws and Considerations	34·3	34.4.7 User Authentication	34·8
		34.4.8 Network Address Translation and NAT-Traversal	34·8
34.3 TECHNICAL ASPECTS OF VoIP SECURITY	34·3	34.5 ENCRYPTION	34·9
34.3.1 Protocol Basics	34·4	34.5.1 Secure SIP	34·9
34.3.2 VoIP Threats	34·4	34.5.2 Secure Real-Time Protocol	34·9
		34.5.3 Session Border Control	34·9
34.4 PROTECTING THE INFRASTRUCTURE	34·6	34.6 CONCLUDING REMARKS	34·10
34.4.1 Real-Time Antivirus Scanning	34·7	34.7 FURTHER READING	34·10
34.4.2 Application Layer Gateways and Firewalls	34·7	34.8 NOTES	34·11

34.1 INTRODUCTION. Whether it is referred to as Voice over Internet Protocol (VoIP) or Internet Protocol Telephony (IPT), the digitization of voice messaging has had and will continue to have an impact on society. Voice messaging is part of a shift that some are calling the Unified Messaging System (UMS).¹ Instant messaging, Short Message System (SMS) text messaging to mobile phones, voice communications, video conferencing, email, file sharing, and network presence are already available in applications that are shared by both the home user and large corporations. Users are empowered by freeing our communications from geographically stationary limits. For example, users can decide to work from home and have their office telephones ring into their laptops. Colleagues can view each other's computer screens from anywhere in the world with Internet access.

34 · 2 SECURING VoIP

Skype, for example, reported about 50 million concurrent online users at the time of writing in mid-January 2013—a phenomenal growth in participation in a UMS created in 2003 (about 590 percent per year over a decade).²

Aside from convenience and capability, there are also great financial justifications for VoIP/UMS systems. They replace formerly expensive equipment with common microcomputers, resulting in substantial savings. Some individuals and small offices are finding that they no longer need to buy landline phones: mobile phones and VoIP software with excellent wireless headsets or earpieces are even better than traditional landline units. An August 2012 report indicated that “New York state has seen a 55% drop in landline phones over the past decade. Some estimate that landlines will be completely obsolete by 2025” and argued that VoIP was “probably the most widely used landline alternative, especially among businesses.”³ By December 2012, U.S. Government reports indicated that more than half of U.S. residents did not have or use landlines in their homes.⁴ Offices are turning to VoIP in the form of Hosted Private Branch Exchanges (PBX) to replace conventional landline PBXs.⁵

Just as there are security issues with our current operating systems and applications, the popularity of VoIP has already drawn the attention of criminals. Most disturbing of all is that many current exploits can be used against VoIP.

34.2 REGULATORY COMPLIANCE. Since VoIP affects key areas of an organization, such as systems, applications, privacy, networks, transmissions, and end users, VoIP must be included in all systematic risk analysis and compliance verification.

Among other U.S. laws, the Gramm-Leach-Bliley (GLB), Sarbanes-Oxley (SOX), and Health Insurance Portability and Accountability Act (HIPAA) are particularly significant in any evaluation of VoIP implementation.

For more information on GLB and SOX, see Chapter 64 in this *Handbook*. For more information about HIPAA, see Chapter 71.

34.2.1 Enhanced 911. Enhanced 911 (E911) mandates a location technology advanced by the Federal Communications Commission (FCC) that enables mobile or cellular phones to process 911 emergency calls and enable emergency services to locate the geographic position of the caller. According to the FCC’s Website:

The FCC has divided its wireless E9-1-1 program into two parts—Phase I and Phase II. Under Phase I, the FCC requires carriers, within six months of a valid request by a local Public Safety Answering Point (PSAP), to provide the PSAP with the telephone number of the originator of a wireless 9-1-1 call and the location of the cell site or base station transmitting the call.

Under Phase II, the FCC requires wireless carriers, within six months of a valid request by a PSAP, to begin providing information that is more precise to PSAPs, specifically, the latitude and longitude of the caller. This information must meet FCC accuracy standards, generally to within 50 to 300 meters, depending on the type of technology used. The deployment of E9-1-1 requires the development of new technologies and upgrades to local 9-1-1 PSAPs, as well as coordination among public safety agencies, wireless carriers, technology vendors, equipment manufacturers, and local wireline carriers.

In a simple solution, a VoIP provider takes the stored location information, passes that through its own call center, and routes it to the local 911. In a more sophisticated solution, some gateways can map MAC addresses to locations, and the 911 calls then can be passed to a PSAP.

Currently, E911 is not required for companies that are using VoIP for internal purposes only. This law is empowered by the FCC, but the Department of Homeland

TECHNICAL ASPECTS OF VoIP SECURITY 34 · 3

Security (DHS) and the Federal Bureau of Investigation (FBI) appear to want it for surveillance and location tracking.

34.2.2 Communications Assistance for Law Enforcement Act (CALEA).

Electronic surveillance consists of either the interception of call content (commonly referred to as wiretaps) or the interception of call-identifying information (commonly referred to as dialed-number extraction) through the use of pen registers and trap and trace devices.⁶

The standards for surveillance of VoIP networks have been a joint effort of industry and law enforcement for some time.

In 2007, the Packet Technologies and Systems Committee (PTSC) published the *Lawfully Authorized Electronic Surveillance (LAES) for Voice over Packet Technologies in Wireline Telecommunications Networks*, Version 2 (Revision of T1.678-2004), which is a surveillance standard for basic VoIP.⁷

A key part of the U.S. Federal Communications Assistance for Law Enforcement Act (CALEA) enforcement and debate is determining who is responsible for compliance. According to the law, “all entities engaged in the transmission or switching of wire or electronic communications as a common carrier for hire.”⁸ However, a question arises when an organization builds its own VoIP solution; that is, does the organization then become a *telecommunications carrier* subject to the law? It would seem so, according to CALEA:

[A] person or entity engaged in providing wire or electronic communication switching or transmission service to the extent that the Commission finds that such service is a replacement for a substantial portion of the local telephone exchange service and that it is in the public interest to deem such a person or entity to be a telecommunications carrier for purposes of this title.⁹

However, the law specifically excludes “(i) persons or entities insofar as they are engaged in providing information services; and (ii) any class or category of telecommunications carriers that the Commission exempts by rule after consultation with the Attorney General.” Given the ambiguity of this exclusionary language, an organization should seek advice from appropriate legal counsel prior to creating its own VoIP network.

34.2.3 State Laws and Regulations.

All states have laws concerning surveillance; 31 states specifically address computers and 14 refer to cell phones. Because a detailed discussion of this diverse environment is beyond the scope of this text, an organization should seek appropriate legal counsel not only concerning the state it is domiciled in, but also the states of any branches, divisions, subsidiaries, or affiliates. The National Conference of State Legislators provides links to the applicable laws of each state and a summary of the coverage (e.g., cell phones, computers, video, photos, etc.).¹⁰

34.2.4 International Laws and Considerations.

Similar to the diversity in the United States, many countries have laws concerning surveillance and intercepts. As such, appropriate legal guidance and advice should be sought before installing or using VoIP networks in international locations. In particular, standards such as ISO/IEC 27002:2005 and privacy laws may be applicable, and may affect VoIP implementation and usage.¹¹

34.3 TECHNICAL ASPECTS OF VoIP SECURITY.

The next section is designed to provide a technical overview of VoIP and related security issues. It starts with an

34 · 4 SECURING VoIP

introduction to the protocols used and then progresses to associated threats. Following that is a discussion of best practices and then encryption.

34.3.1 Protocol Basics

34.3.1.1 Audio Stream Protocols: RTP and UDP. Real-time Transport Protocol (RTP) is the packet-based communication protocol that provides the base for virtually all VoIP architectures. Part of RTP is timing and packet sequence information that can be used to reconstruct audio streams. Like TCP, User Datagram Protocol (UDP) is a layer 4 network communication protocol. Both of these protocols provide the basic packet addressing needed to get from one network address to another. When dealing with VoIP, packet delivery time is of the essence. Both TCP and UDP add overhead to the organization's network communication. This overhead results in greater time delays to the communication. UDP adds less overhead as it provides comparatively little packet error checking. The compromise is that UDP will help deliver packets quicker but will also result in more lost packets. It has been observed that a packet loss of 10 percent spread over a VoIP call may be virtually undetectable to most users. However, all of the features that can make VoIP secure also add latency. Added latency means more lost or discarded packets, which then results in an unpleasant or unacceptable user experience. In the end, a slow packet is a lost packet.

34.3.1.2 Signaling Protocols: SIP and H.323. Session Initiation Protocol (SIP) is a protocol that is used to establish interactive multimedia sessions between users. In addition to VoIP, it is also used for video conferencing and online gaming. SIP appears to be the most commonly accepted form of establishing VoIP calls. Secure SIP (SSIP) is discussed in Section 34.5.1. Like SIP, H.323 can be used for video conferencing or VoIP call setup. While SIP appears to be the standard for new installations. One may find H.323 in use at enterprise scale installations that have large investments in older, analog communication equipment.

Briefly, the basics of a VoIP call include these details. A VoIP client application, whether on a personal computer or a dedicated handset, uses SIP or H.323 to set up the call. This call setup is an exchange of control parameters that may include encryption and compression algorithms to be used. This is called "signaling." Once the call is set up, the VoIP client uses RTP to start packetizing the voice data. The RTP packets are incorporated into a UDP packet that adds addressing and sequencing information. The UDP packets are collected and sorted by sequence number at the receiving station. Some systems use a "jitter buffer" to assemble and store the packets. The endpoint VoIP client then reads the jitter buffer and turns the RTP packets back into voice.

34.3.2 VoIP Threats. Although not an exhaustive list, these hacks represent some of the vulnerabilities likely to be encountered when using VoIP:

- SPam over Internet Telephony (SPIT)
- Eavesdropping
- Theft of service
- Man-in-the-middle attacks

34.3.2.1 SPIT. Most users have become accustomed to finding 10 or 20 (or 100 or 200) spam email messages in an inbox on daily basis. The thought of receiving an equal number of voicemails on a daily basis leads to the unappetizing acronym

TECHNICAL ASPECTS OF VoIP SECURITY 34 · 5

SPam over Internet Telephony (SPIT). Would-be SPIT spammers may have been both encouraged and then disheartened when, in 2004, Qovia,

a company that sells enterprise tools for VoIP monitoring and management, ... applied for a patent on technology to broadcast messages via VoIP—and another one for a method of blocking such broadcasts. The broadcast methodology only works on a pure VoIP network, while most of today's services are hybrids of IP and traditional telephone lines.

The author explains that Qovia realized that broadcasting VoIP messages could be useful for agencies such as Homeland Security but could also be abused by spammers. Therefore, Qovia pledged to “incorporate its SPIT-blocking technology in future releases of its security products, while enforcement of its patent on broadcasting, if granted, could be used to shut down VoIP spammers.”¹²

Sam Rozenfeld, writing in March 2012, summarized some of the measures being used to interfere with SPIT.¹³

- VoIP providers can collaborate with business clients to use network firewalls to block inbound SPIT using blacklisting of sources known to the provider as large-scale SPIT senders.
- It may be possible for the VoIP provider to identify robocalls—automated large-scale SPIT floods—and block them.
- Authenticating the presence of a human being before allowing the caller to complete a call can block SPIT.
- Blacklisting specific individual numbers can help individual users suffering repeated messages from the same sources.
- Since some SPIT sources use area codes from which no legitimate traffic would normally occur, blocking entire area codes may work.
- Whitelisting might work for individuals who know exactly who should be using their phone number, but it is unlikely to be acceptable for businesses who must accept calls from potential customers.

34.3.2.2 Eavesdropping. In an unsecured VoIP environment, eavesdropping is reduced to a task that is quite straightforward. Earlier, it was mentioned that RTP is the de facto protocol for VoIP communication. RTP adds unique sequencing information to its packets. Because of this, one could collect a number of RTP packets and then assemble them in a consecutive order, as a receiving station would collect these packets in a jitter buffer.

Knowing this, the first step in eavesdropping is to obtain a packet-sniffing tool, such as Ethereal from www.wireshark.org. One may then use Ethereal to perform a packet capture, or one may obtain sample capture files from the Ethereal Website. If the user chooses to perform packet capture, care must be exercised not to violate any privacy laws applicable to the user's network. Once obtained, Ethereal is used to sort out any RTP packets. Exhibit 34.1 depicts a packet capture filtered to show RTP packets. Notice that the first column at the left contains a list of sequential numbers. These numbers indicate the packet placement within the capture. However, farther to the right is the actual conversation sequence number. Once the RTP packets are assembled by sequence number, they can be saved as an “.au” file, which can be played on most computers. As shown with this example, eavesdropping on unencrypted VoIP traffic is not a complicated or expensive process.

34 · 6 SECURING VoIP

Filter: rtp.payload					▼	Expression...	Clear	Apply
No.	Time	Source	Destination	Protocol	Info			
624	1444.509099	192.168.1.2	212.242.33.36	RTP	Payload type=ITU-T G.711 PCMA, SSRC=932629361, Seq=28590, Time=1240			
625	1444.579046	192.168.1.2	212.242.33.36	RTP	Payload type=ITU-T G.711 PCMA, SSRC=932629361, Seq=28591, Time=1400			
626	1444.582379	192.168.1.2	212.242.33.36	RTP	Payload type=ITU-T G.711 PCMA, SSRC=932629361, Seq=28592, Time=1560			
627	1444.588245	192.168.1.2	212.242.33.36	RTP	Payload type=ITU-T G.711 PCMA, SSRC=932629361, Seq=28593, Time=1720			
628	1444.590352	192.168.1.2	212.242.33.36	RTP	Payload type=ITU-T G.711 PCMA, SSRC=932629361, Seq=28594, Time=1880			
629	1444.625165	192.168.1.2	212.242.33.36	RTP	Payload type=ITU-T G.711 PCMA, SSRC=932629361, Seq=28595, Time=2040			
630	1444.627060	192.168.1.2	212.242.33.36	RTP	Payload type=ITU-T G.711 PCMA, SSRC=932629361, Seq=28596, Time=2200			
631	1444.664688	192.168.1.2	212.242.33.36	RTP	Payload type=ITU-T G.711 PCMA, SSRC=932629361, Seq=28597, Time=2360			
632	1444.671724	192.168.1.2	212.242.33.36	RTP	Payload type=ITU-T G.711 PCMA, SSRC=932629361, Seq=28598, Time=2520			

EXHIBIT 34.1 Packet Capture Showing RTP Packets

34.3.2.3 Theft of Service. One of the classic theft-of-service attacks that has occurred was reported in June 2006. In this case, a brute-force attack yielded special access codes allowing attackers entry into the provider network. Once in, the attackers obtained router passwords and login credentials; they then programmed the network transport devices to implicitly accept and route VoIP messages from the attackers' server. The attackers then sold VoIP access to public providers, who then sold it to the public. In the end, 15 exploited companies were left to pay the bill for the calls that were routed out of their network.¹⁴

An organization's own employees could exploit its VoIP infrastructure as these attackers did.

34.3.2.4 Man-in-the-Middle Attacks. As with any technology that digitizes communications, VoIP that sends data without encryption is inherently open to manipulation if an attacker can intercept traffic, alter its content, and send it on its way to the recipient in a classic man-in-the-middle (MITM) attack.¹⁵ Once an attacker has control of VoIP traffic, the attacker can not only eavesdrop but can also:

- Initiate calls to a third party and impersonate a caller by sending data that appear to come from a legitimate phone belonging to someone else
- Deflect calls to the wrong destination
- Intercept traffic in real time and generate simulated voice content to create misleading impressions or cause operational errors

This last point warrants expansion. An eavesdropper could collect the digital patterns corresponding to words or phonemes generated by a particular user during normal conversation. Using simple programs, it would be possible to generate data streams corresponding to any spoken sequence including those words or phonemes in real time, allowing the MITM to feed the recipient (or even both sides of a conversation) with distorted or invented information and responses. The potential for mayhem is enormous, especially if the conversation involved, say, emergency response.

34.4 PROTECTING THE INFRASTRUCTURE. This section focuses entirely on VoIP networks. For general information on infrastructure protection, see Chapters 22 and 23 in this *Handbook*.

PROTECTING THE INFRASTRUCTURE 34 · 7

34.4.1 Real-Time Antivirus Scanning. Protecting the VoIP infrastructure would appear to be a routine decision. Remember that as with VoIP, a slow packet is a lost packet. Many of the routine measures used to protect a typical server will introduce latency into the voice system, leading to jitter and sporadic communication. One of the first routine measures to be sacrificed is real-time antivirus protection of the VoIP server. Some vendors will suggest that real-time scanning of the organization's entire VoIP server is next to impossible, and unsupportable. This is an unacceptable myth. The VoIP server must be scanned in a fashion that is at least consistent with other production servers. Requests from system administrators asking to disable real-time scanning is a common first step in creating issues with the VoIP system.

34.4.2 Application Layer Gateways and Firewalls. An organization's VoIP infrastructure is much more than a series of systems that are digitizing voice and forwarding packets. Its servers will have real-time contact with email servers and central authentication systems using Remote Authentication Dial-in User Service (RADIUS) or perhaps Active Directory. It may also have database systems devoted to logging call information or even recording the calls themselves. An attacker who gains access to part of the organization's VoIP infrastructure may be able to access the most sensitive parts of its network. Consider the use of application layer gateways (ALGs) or SIP/VoIP-aware firewalls to segregate the VoIP systems.

34.4.3 Logical Separation of Voice and Data. It would be ideal to have a separate network for the organization's VoIP system. Bandwidth issues would be minimized and troubleshooting simplified. In most instances, it may not be possible to make a business case justifying the installation of a separate set of cables and network gear. The VoIP handset on the user's desk or the VoIP softphone in the user's PC will share the same wire as the workstation. The logical separation of voice and data begins with assigning the organization's VoIP devices to a network subnet separate from the data devices. This is initiated by a Dynamic Host Configuration Protocol (DHCP) request from the user's handset. Part of this request allows the DHCP servers to distribute addresses based on hardware identification parameters. For example, connecting a laptop to the organization's network will result in the assignment of an address that is on a subnet different from the Cisco handset that is plugged into the same connection. This logical separation of voice and data allows the organization's firewalls to protect its VoIP infrastructure by screening out protocols and requests that are not voice related.

34.4.4 Quality of Service. The term "quality of service" (QoS) refers to a set of configurable parameters that can be used to control and/or prioritize communication through the network. Again, with respect to VoIP, a slow packet is a lost packet. QoS can be used to prioritize the VoIP packets so they can be delivered in a timely fashion. Most vendors will provide a choice of default QoS configurations to be used according to the organization's needs. Some even provide VoIP firewalls that are capable of buffering VoIP messages and retransmitting packets. Regarding QoS, it is not just one parameter that can be turned on to make the VoIP installation work properly; it is a series of parameters that may need to be tuned to fit the organization's exact needs. For more technical detail, see IEEE standards 802.1p and 802.1q.

34.4.5 Network Monitoring Tools. Best practices require a dedicated security operation center (SOC) watching the organization's networks for attacks. Whether

34 · 8 SECURING VoIP

it is a 7×24 SOC or just one network administrator who does everything, it is absolutely critical to provide the tools and training necessary to detect attacks and to troubleshoot performance issues. With VoIP, the staff will be facing attacks using a new set of attack vectors and protocols. If the network administration is outsourced, review of existing service-level agreements to guarantee that the provider is capable of supporting VoIP is strongly recommended.

34.4.6 Device Authentication. Device authentication can be accomplished in a variety of ways. The simplest is to store a list of device MAC addresses on the organization's VoIP server and to authenticate all SIP requests through that list of addresses. The standard is to deploy devices to desktops without configuration. Upon connection, a technician enters a setup utility that connects to the VoIP server and then downloads a preconfigured image.

34.4.7 User Authentication. The organization's finance people likely will demand some type of call tracking and usage so departments can be charged or reviewed appropriately. Users can also be placed into different groups that the organization can configure. This group-level access can be used to limit services, such as long distance or international calling. It is common for a VoIP infrastructure to have ties to a central authentication server, such as Lightweight Directory Access Protocol (LDAP) or Active Directory. This central authentication will ease functions such as forwarding voicemail to a personal computer or cell phone. There are two common problems to watch for:

1. The *authentication interval* will be set by user. The organization can configure a demand for authentication to be hours, days, or months; 24 hours is suggested. Some users complain, and demand that their central credentials be stored in the handset for months so they only need to log in infrequently; this practice is generally undesirable.
2. Most handsets will come with *default accounts and passwords*. These accounts must be disabled or, at the least, strong password discipline must be maintained over them.

34.4.8 Network Address Translation and NAT-Traversal. Network address translation (NAT) is a technique commonly used by firewalls and routers to allow multiple devices on an internal network to share one IP address on the Internet. A user's internal address should be known only to systems on that user's own network. When connecting to an external network, the organization's router or firewall forwards the user's communication to an external address but replaces the user's private address with its public address. When the communication is returned, the firewall routes the message back to the correct private address. Similarly, consider the broadband router in a home. The user may connect a series of systems to this device, each with a unique internal address distributed by the router, but the ISP sees it as only one address. These internal or private addresses are stored in what is commonly called a "translation table." At a higher level, the process of getting a packet through a NAT device is called NAT-Traversal (NAT-T). Understanding how NAT-T issues affect VoIP is vital to understanding the danger of sending a voice call to the Internet.

Section 34.3.1.2 outlined how SIP is commonly used to set up a call. Once the call is set up, the actual audio stream typically is relayed via RTP/UDP. This is where

ENCRYPTION 34 · 9

the issues start. The firewall does not have a problem with passing SIP traffic back and forth to the Internet, as the internal address of the VoIP device is stored in the translation table. However, the SIP signaling has passed on the private address of the user's VoIP device. This means that a device outside of the local network is trying to send the RTP/UDP audio stream to a fictional IP address. Essentially, NAT prevents VoIP from functioning.

Several work-arounds are commonly used. Sometimes the NAT device can be configured to provide VoIP support. Sometimes VoIP devices can be configured to work over otherwise open ports, overloading a common protocol such as HTTP, unfortunately, often with unintended side effects. Or VoIP proxy servers can be used on either side of the NAT in order to facilitate the traversal. Each of these solutions opens up its own security concerns, which should be carefully addressed; these concerns include the consequences of external proxy servers or creating anomalous traffic over other protocols, as in the overloading example.

34.5 ENCRYPTION. Encryption plays a critical role in communications security. For a general introduction to encryption, see Chapter 7 in this *Handbook*; for more details of public key encryption, see Chapter 37.

34.5.1 Secure SIP. Transport Layer Security (TLS) was sponsored by the Internet Engineering Task Force (IETF) to secure and encrypt data communications crossing public networks. It is intended to replace Secure Sockets Layer (SSL) as a widely accepted form of securing data communication. This protocol consists of a “handshake” and a “record.” TLS was designed to be application independent so developers could choose their own way of initiating a TLS session.

Secure SIP is a mechanism designed to send SIP signaling messages over an encrypted TLS channel. A SSIP session is initiated by a SIP client contacting a SIP proxy and requesting a TLS session. The proxy returns a certificate that the SIP client then authenticates. The client and proxy then exchange encryption keys for the session. If the call is destined for another network segment, the SIP proxy will contact that segment and negotiate a sequential TLS session, so the SIP message is protected by TLS the entire time.

34.5.2 Secure Real-Time Protocol. Secure Real-Time Protocol (SRTP) is an enhancement of RTP that provides encryption, authentication, and integrity to the VoIP audio stream. The Advanced Encryption Standard (AES) originally was a block cipher; SRTP incorporates AES into the data stream with an implementation that utilizes it as a stream cipher.

Encryption is good but it does not protect the user or organization against replay attacks. SRTP uses a Hashed Message Authentication Code (HMAC-SHA1) algorithm to provide authentication and integrity checks. The MAC is calculated using a cryptographic hashing function in conjunction with a private key. SRTP uses one of the five Secure Hash Algorithms (SHA) designed by the National Security Agency. All five of these algorithms are compliant with requirements set in the Federal Information Processing Standards (FIPS).

34.5.3 Session Border Control. To this point, a number of issues affecting a VoIP deployment have been identified. Session border control (SBC) is a set of services that address VoIP issues related to security, QOS, NAT traversal, and network

34 · 10 SECURING VoIP

interoperability. SBC collects real-time bandwidth statistics that can be used to allocate the network resources necessary to maintain the QOS desired. SBC will also support a number of NAT-T algorithms that will allow calls to be routed to public networks while maintaining the anonymity of internal resources. At the same time, SBC can accommodate both SIP and H.323. This allows the signaling protocol translation necessary to connect both types of networks.

34.6 CONCLUDING REMARKS. VoIP provides expanded functionality and lower costs for corporate users, but managers must integrate security considerations into the architecture and implementation of all such systems to prevent interception, deception, and denial-of-service attacks. In addition, technologists must monitor developments in this rapidly changing field to keep abreast of new attack methodologies and countermeasures.

34.7 FURTHER READING

- Boyter, B. "Voice-over-IP Sniffing Attack," 2003, www.giac.org/certified-professionals/practicals/gcih/0442.php (URL inactive).
- Collier, M., and D. Endler. *Hacking Exposed Unified Communications & VoIP Security Secrets & Solutions*, 2nd ed. McGraw-Hill Osborne Media, 2013.
- Davidson, J., J. Peters, and B. Gracely. *Voice over IP Fundamentals*. Indianapolis, IN: Cisco Press, 2000.
- Endler, D., and M. Collier. *Hacking Exposed—VoIP: Voice Over IP Security Secrets & Solutions*. New York: McGraw-Hill Osborne Media, 2007. See also the associated Website, www.hackingVoIP.com
- Kuhn, D. R., T. J. Walsh, and S. Fries. "Security Considerations for VoIP Systems." NIST Special Publication 800-58, 2005, <http://csrc.nist.gov/publications/nistpubs/800-58/SP800-58-final.pdf>
- Long, T. "Eavesdropping an IP Telephony Call." GIAC Security Essentials Certification Practical Assignment, 2002, www.sans.org/reading_room/whitepapers/telephone/318.php
- Miller, W. *Wireless Mesh Network Security: An Overview*. Death and Blind Media, 2013
- Molitor, A. "Deploying a Dynamic Voice over IP Firewall with IP Telephony Applications," 2000, http://cnscenter.future.co.kr/resource/rsc-center/vendor-wp/aravox/aravox_deploying_dynamic.pdf (URL inactive).
- Molitor, A. "Securing VoIP Networks with Specific Techniques, Comprehensive Policies and VoIP-Capable Firewalls," 2000, http://cnscenter.future.co.kr/resource/rsc-center/vendor-wp/aravox/aravox_specifictechniques.pdf (URL inactive).
- Porter, T., B. Baskin, L. Chaffin, M. Cross, J. Kanclirz, A. Rosela, C. Shim, and A. Zmolek. *Practical VoIP Security*. Rockland, MA: Syngress, 2006.
- Thalhammer, J. "Security in VoIP Telephony Systems." Master's thesis, Institute for Applied Information Processing and Communications at the Graz University of Technology, Graz, Austria, 2002; www.iaik.tugraz.at/teaching/11_diplomarbeiten/archive/thalhammer.pdf (URL inactive).
- Thermos, P., and A. Takanen. *Securing VoIP Networks: Threats, Vulnerabilities, and Countermeasures*. Boston, MA: Addison-Wesley, 2007.
- Verma, P. K., and L. Wang. *Voice over IP Networks: Quality of Service, Pricing and Security*. Springer, 2013.
- VOIP Security Alliance White Papers: www.voipsa.org/Resources/whitepapers.php

NOTES 34 · 11**34.8 NOTES**

1. J. M. S. Wams, “Unified Messaging Atop a Cloud of Micro-Objects,” (doctoral dissertation, Vrije Universiteit Amsterdam, December 18, 2012), <http://dare.uvbu.vu.nl/handle/1871/39352>
2. Jean Mercier, “50 Million Concurrent Users Online!” *Skype Numerology*, January 21, 2013, <http://skypenumerology.blogspot.com/2013/01/50-million-concurrent-users-online.html>
3. Kate Harrison, “How To Break Up With Your Landline.” *Forbes* Website, August 8, 2012. www.forbes.com/sites/kateharrison/2012/08/08/how-to-break-up-with-your-landline
4. Stacey Higginbotham, “Over Half of American Homes Don’t Have or Use Their Landline.” *GIGaom*, December 26, 2012, <http://gigaom.com/2012/12/26/over-half-of-american-homes-dont-have-or-use-their-landline>
5. Enhanced VoIP Communications Inc., “What is a Hosted PBX?” *Easy Office Phone* Website, 2013, www.easyofficephone.com/resources/hosted-pbx
6. Electronic Frontier Foundation, “CALEA FAQ: What are the different kinds of communications surveillance methods, and how are they relevant to the current debate regarding CALEA?” *Electronic Frontier Foundation: Defending Your Rights in the Digital World*, n.d. <https://www.eff.org/pages/calea-faq#2>
7. IHS, “ATIS Releases LAES Standard for Internet Access, Services—ATIS PP-1000013,” *IHS News & Analysis*, March 30, 2007, www.ihs.com/news/atis-laes-internet.htm
8. Alliance for Telecommunications Industry Solutions. “Lawfully Authorized Electronic Surveillance (LAES) for Voice over Packet Technologies in Wireline Telecommunications Networks, Version 2,” *ATIS Document Center*, May 2006, www.atis.org/docstore/product.aspx?id=22771
9. Cornell University Law School, “CALEA: 47 USC §1001—Definitions (Preliminary Version),” *Legal Information Institute* Website, January 4, 2012, www.law.cornell.edu/uscode/text/47/1001
10. National Conference of State Legislatures, “Electronic Surveillance Laws,” *NCSL* Website, 2013, www.ncsl.org/issues-research/telecom/electronic-surveillance-laws.aspx
11. International Trade Administration, “Worldwide VoIP Regulatory and Market Information,” *U. S. Department of Commerce*, November 3, 2009, [http://web.ita.doc.gov/ITI/itiHome.nsf/f3e8a6b8413b3e5d85256cc40075c5df/cb2a434afea6790485256d020053fef0/\\$FILE/voip%20worldwide%202009.pdf](http://web.ita.doc.gov/ITI/itiHome.nsf/f3e8a6b8413b3e5d85256cc40075c5df/cb2a434afea6790485256d020053fef0/$FILE/voip%20worldwide%202009.pdf) (URL inactive).
12. Susan Kuchinskas, “Don’t SPIT on VOIP,” *Small Business Computing.com*, August 4, 2004, www.smallbusinesscomputing.com/news/article.php/3399011/Dont-SPIT-on-VOIP.htm
13. Sam Rozenfeld, “No more SPIT!” *Business VoIP/VoIP Security*, March 25, 2012, www.telephonyyourway.com/2012/03/25/no-more-spit
14. “VoIP Hacker Arrested on Fraud Charges,” *Technology News Daily*, 2006, www.technologynewsdaily.com/node/3252 (URL inactive).
15. See, for example, P. Thermos, “Two Attacks against VoIP,” *Security Focus*, 2006, www.securityfocus.com/infocus/1862/1 (URL inactive).

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 35

SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

Carl Ness

35.1	INTRODUCTION	35·1	35.5	SECURING SMS	35·14
35.2	GENERAL CONCEPTS AND DEFINITIONS	35·1	35.5.1	Dangers to the Business	35·14
35.2.1	Peer to Peer	35·2	35.5.2	Prevention and Mitigation	35·15
35.2.2	Instant Messaging	35·2	35.5.3	Reaction and Response	35·18
35.2.3	Short Message Service	35·3			
35.2.4	Collaboration Tools	35·3	35.6	SECURING COLLABORATION TOOLS	35·18
			35.6.1	Security versus Openness	35·18
35.3	PEER-TO-PEER NETWORKS	35·4	35.6.2	Dangers of Collaboration Tools	35·19
35.3.1	Dangers to the Business	35·5	35.6.3	Prevention and Mitigation	35·20
35.3.2	Prevention and Mitigation	35·7	35.6.4	Reaction and Response	35·22
35.3.3	Response	35·8			
35.3.4	Case Study	35·9			
35.4	SECURING INSTANT MESSAGING	35·9	35.7	CONCLUSIONS	35·22
35.4.1	Dangers to the Business	35·9	35.8	FURTHER READING	35·23
35.4.2	Prevention and Mitigation	35·11	35.9	NOTES	35·23
35.4.3	Response	35·13			
35.4.4	Safe Messaging	35·13			

35.1 INTRODUCTION. Peer-to-peer (P2P) communications, instant messaging (IM), short message service (SMS), and collaboration tools must be directly addressed in any comprehensive security plan. The dangers are real, as is the probability that at least one of these technologies is in use in almost every organization.

35.2 GENERAL CONCEPTS AND DEFINITIONS. This chapter is designed to present enough information and resources to aid in integrating the defense of each function into the organization's security plan. A list of resources is provided at the end of the chapter to aid in further research.

35 · 2 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

35.2.1 Peer to Peer. Peer-to-peer networking, also referred to as P2P, is not a new concept or technology. The term was contained in some of the original designs and proposals for the Internet as an efficient and logical way to exchange information from one resource, or peer, to another, on a large interconnected network. Today, the term is most associated with applications that transfer multimedia files across the Internet.

Peer-to-peer networks generally consist of different computers, or nodes, that communicate directly with each other, often with little, if any, need for a central computer to control the activity. Often utilizing an application with a client-server appearance, the two computers set up a direct connection between each other for file transfer. A central indexing computer may or may not be needed to help these computers “find” each other, to index and publish their contents, or to facilitate the connection. However, what is most important, the two computers must have a direct, logical connection to transfer the file or files. File transport may take place over a local network (LAN), a wide area network (WAN), a value-added network (VAN), or via the Internet.

Peer-to-peer technologies and applications were much more common in the early days of networking when it was not financially possible for many organizations to have expensive servers and complicated network topologies. This is especially true of personal computer networks that performed simple file sharing from computer to computer in a one-to-one model instead of today’s much more common one-to-many server-to-client setup. However, there are legitimate uses for peer-to-peer technologies. One common example is the sharing of Linux distribution software images. Peer-to-peer sharing of these often large ISO disk images requires much fewer resources for the distributor, because there may be thousands of computers distributing the software among themselves, instead of every user trying to download the file from a single server.

35.2.2 Instant Messaging. Instant messaging, or IM, has become one of the most widely used communication mediums and is on pace to overtake email as the preferred technology to communicate with others. This tool allows users to communicate with each other in a real-time, synchronous, instantaneous fashion via computer, tablet, or mobile device. Today’s IM applications are nowhere near the first generation of IM. The concept of communicating, or chatting, in real time made its appearance on multiuser computer systems, when users could initiate a text-based conversation with each other. The most common example of this type of communication was a host-based system such as a mainframe environment or UNIX system using programs like *talk* or *ytalk*. Initially, users may have been restricted to messaging each other when logged into the same machine; eventually users were able to communicate with each other, either via Internet Relay Chat (IRC) or via an early online service like America Online. The first widespread uses of IM were made possible with the popularity of the PC with a modem and were used mainly for brief, informal, personal conversations.

With time, IM has become a business tool and, in some organizations, a necessity, especially with telecommuters. The need to communicate with colleagues, salespeople, clients, customers, and the like has transformed a gimmick technology into ubiquity. With this change, it is necessary for security management to change and adapt accordingly. Users are able to send messages, files, real-time streaming video and audio, utilize collaborative online whiteboards, and share desktops, almost instantly. Essential to the organization or not, IM can become a dangerous medium for security breaches.

GENERAL CONCEPTS AND DEFINITIONS 35 · 3

35.2.3 Short Message Service. Short message service, or more commonly SMS, is another previously minor technology that has become ubiquitous and a large part of everyday life for many people. Although some mobile phone standards and companies had different ideas for the uses of SMS, a common early use was to notify customers of information one way, from the mobile phone provider to the user. A popular example was alerting the user of a missed call or voicemail message. Many carriers never dreamed customers would actually be able to send text messages from one mobile phone to another, nor did the carriers think users would ever *want* to do such a thing. The name, short message service, also implied a limited amount of text a message could contain. Originally, users were limited to 160 characters or less.

SMS has morphed into something much larger. The commonality of mobile phones has pushed the original concept far beyond its original meaning and function. Today, two-way communication between mobile phone customers, often on different mobile phone carrier networks, between customers and mobile phone providers, and between customers and other information systems has become a way of life. Customers expect instant, always-on, reliable SMS services. Most mobile phones are capable of SMS text messages, taking and sending pictures, instant alerts, and a number of other services that utilize or expand on the original concept of short message service.

35.2.4 Collaboration Tools. People working together have created a need for even more technology to aid them in completing their tasks. There are many products in today's market to facilitate sharing, collaboration, and organization of data. As some information security professionals joke, "Computers and technology are generally safe and secure, until you let a human near them." Humans are inevitable when it comes to collaboration tools and systems. Many collaboration tools and systems are designed to aid workgroups that are physically far apart. Once a system has requirements that contain the words "open," "via the Internet," or "access from anywhere," information security professionals cringe. Securing collaboration tools can be difficult, especially when it comes to balancing functionality versus security. A handbook would not be complete if the cloud were omitted, which is precisely where most of the current collaboration tools reside, including their data. These tools are also creating a tug-of-war between users and security professionals—users want to use these tools; security pros worry about the consequences.

File-sharing services such as Dropbox and Google Docs have put effort into securing the data stored and shared by individuals and groups. For example, Dropbox answers the question "How secure is Dropbox?" as follows:

We have a dedicated security team using the best tools and engineering practices available to build and maintain Dropbox, and you can rest assured that we've implemented multiple levels of security to protect and back up your files. You can also take advantage of two-step verification, a login authentication feature which you can enable to add another layer of security to your account.

Other Dropbox users can't see your files in Dropbox unless you deliberately share links to files or share folders. Dropbox employees are prohibited from viewing the content of files you store in your account. Employees may access file metadata (e.g., file names and locations) when they have a legitimate reason, like providing technical support. Like most online services, we have a small number of employees who must be able to access user data for the reasons stated in our privacy policy (e.g., when legally required to do so). But that's the rare exception, not the rule. We have strict policy and technical access controls that prohibit employee access except in these rare circumstances. In addition, we employ a number of physical, technical, and heuristic security measures to protect user information from unauthorized access.¹

35 · 4 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

Nonetheless, there are serious security concerns about cloud-based file-sharing tools and Dropbox in particular. At a Black Hat EU conference in 2013, the paper “DropSmack: How cloud synchronization services render your corporate firewall worthless” caught the eye of writer Michael Kassner. The paper summary included the following points:

- “... [C]loud-based synchronization solutions in general, and Dropbox in particular, can be used as a vector for delivering malware to an internal network.”
- “... Dropbox synchronization service can be used as a Command and Control (C2) channel.”
- “... [F]unctioning malware is able to use Dropbox to smuggle out data from exploited remote computers.”

The paper’s author, experienced penetration-tester Jacob Williams, warned that if a bad actor has any access to a secured Dropbox folder, it is possible to synchronize a remote-access Trojan he wrote called DropSmack with all the shared Dropbox folders. The tool would allow the infiltration of the entire corporate network. Williams also warned that access to a Dropbox folder by employees using their personal computers raises legal issues:

Many general counsels are more than a little worried about the appearance of authorizing us to pen test what could end up being home machines. That’s becoming a sticky issue with pen-testers these days as people open spear phishing emails delivered to the corporate email addresses on machines that may be privately owned.²

Integrated collaboration tools have another danger, especially for inexperienced users who don’t take daily backups: deleting one or more files (or all of them) in a shared Dropbox folder will propagate the deletion to all users of the shared folder. If any one of the users keeps a daily backup, the entire group of users can be protected against disaster; if none of them do, they may be in serious difficulty. In a related issue, any user who moves the files out of the Dropbox folder into a local folder will wipe the data from all the other users’ Dropbox folders too.³

In a 2012 article, Matthew J. Schwartz urged corporate uses to pay attention to Dropbox use among their employees. His five recommendations (with more details in the original article) are:

1. Monitor Dropbox use
2. Compare cloud service security
3. Beware lackluster security cloud practices
4. Treat Dropbox as a public repository
5. Beware insider data exfiltration⁴

A free tool, Cloudfogger, automatically encrypts data on the client side when it is uploaded to any external collaboration tool using 256-bit AES encryption. The tool then automatically decrypts the data when it is downloaded by an authorized user.⁵

35.3 PEER-TO-PEER NETWORKS. One of the earliest mass applications of P2P was for free file-sharing of music through Napster, LLC. Despite difficulties over copyrights and a subsequent bankruptcy, Napster’s technology, in substantially the same form, is still in widespread use. Practical applications have expanded beyond music

PEER-TO-PEER NETWORKS 35 · 5

downloads into the business world, such as allowing small groups of users to share files without the interaction of a systems administrator and distribution of open source software. Likewise, it may be possible that employees are utilizing the organization's high-speed Internet connection to supplement their at-home movie collection via P2P downloads.

35.3.1 Dangers to the Business. Using P2P technology without proper care and controls, an organization may face serious consequences. There are many threats to an organization that does not properly control P2P networking, as for any other network configuration or protocol. Many problems are discussed in Chapters 21, 25, and 26 in this *Handbook*. However, this section contains several important issues that information security management should consider while performing risk analysis and policy implementation for P2P networking.

35.3.1.1 Abusing Company Resources and Illegal Content. Organizations must have an acceptable usage policy in place, one that limits what employees can do with the technology resources provided to them. The policy should clearly state the kinds of technology and applications that are prohibited or restricted in specific ways. In most cases, P2P technology used to download music or videos for personal use will violate the policy.

P2P technology is a specific danger to company technology resources because the inherent nature of P2P technology is to use every resource to the maximum extent possible. For example, a single P2P application, configured properly, will use every bit of bandwidth that is made available to it. This would include LAN bandwidth, WAN bandwidth, and Internet bandwidth. One of the most popular uses for P2P technology still remains the sharing of extremely large files, especially multimedia files, including full-length movies. These large files can take hours to download in full. This fact can have an extremely negative impact on an organization's network infrastructure—including expensive Internet bandwidth.

In practice, a single P2P application has been demonstrated to completely saturate a 12-megabit Internet connection, virtually denying, or severely limiting, access to all other computers.⁶ In this case, these dangers to the business are common to many areas of information security management:

- Threat to **availability**. If an organization's resources, including network resources, are not available, the business cannot properly function.
- Threat to **integrity**. If the organization's resources are crippled or misused by employees utilizing P2P technology, data may suffer from a breakdown of integrity and usability.
- Threat to the **organization's image**. If the organization's information systems and infrastructure cannot be relied on because of interruptions from P2P abuses, there is a risk of financial or public image degradation. Some organizations are not able to overcome a substantial loss of image, credibility, or both.
- Threat from **litigation**. It is common to see illegal content being shared via P2P technology; illegal music and video sharing is often credited with having made P2P technology popular. An organization may suffer legal troubles, including copyright and intellectual property suits, if its resources are involved with the sharing of illegal materials. Some anti-piracy groups have become extremely aggressive in combating illegal sharing of copyrighted content.

35 · 6 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

35.3.1.2 Loss of Confidentiality. There are many ways an organization may suffer a loss of confidentiality from P2P technology. One common mistake is a misconfigured P2P application. The case study in Section 35.3.4 describes one situation. However, the dangers of a misconfigured P2P application are real—it is quite easy for data to be shared inadvertently. When users are in a hurry or do not understand what they are doing, a P2P application may allow for unauthorized access to information because its restrictions are too lax or missing altogether. A common mistake in a Microsoft Windows environment may be to share the entire “My Documents” folder when a user intended to share only photos. A misconfigured or maliciously altered P2P application may also become a conduit or access point for an attacker to enter an otherwise secure network environment.

Another, less well-known and often overlooked threat involves the amount of data a P2P application can reveal to unauthorized persons. For example, a P2P application may offer detailed information about its host, including:

- Operating system, version, and configuration
- Corporate network address scheme, host naming convention, DNS information
- Detail about the P2P application version or build (useful for attackers to exploit known vulnerabilities in a “buggy” release or version)
- Network routes, privileged access from a certain PC to sensitive networks within the organization (behind a firewall, etc.)
- Open network ports in the organization’s firewall

Although many of these examples may seem rather benign by themselves, the P2P application may be revealing information that an attacker can use as part of a bigger attack. Chapter 19 of this *Handbook* details how small pieces of information can be gathered and used together in an information security breach. The nature and functionality of P2P applications leaks sensitive information that otherwise would not be revealed.

All P2P applications are not created equal; a P2P application may be different from the user’s expectations. Can the P2P application actually be a reliable, malware-free, secure application—especially when the application is a free download from the Internet? It is possible that a backdoor exploit, malware, spyware, or the like may be built into the P2P application, or introduced later. This was especially true in the days of Napster; many applications included unwanted malware that ranged from innocent to downright dangerous.⁷ Similar exploits are still possible.

Many users still do not completely understand the application’s functionality. What may start out as a user trying to download a single movie can turn into an application that never gets uninstalled and is always running, uploading, or “seeding” that single movie for weeks or months on end. Since many of these applications are built to silently operate in the background, the user may think that the movie is downloaded, and the application isn’t used any longer, when in reality, it is still running until it is removed. This is not only a waste of resources, but a likely way to expose the organization to a DMCA (Digital Millennium Copyright Act) complaint from the organization that owns the intellectual property.

35.3.1.3 Consequences. Any organization that does not protect against data loss via P2P networking is at great risk of public disclosure and scrutiny, financial

PEER-TO-PEER NETWORKS 35 · 7

penalties, regulatory penalties, and so on. The functionality and nature of P2P applications may provide an investigator or, worse yet, the press with definitive evidence of the use of P2P technology within an organization. (It is easy to discover an organization that has users utilizing P2P applications because there are several online databases of IP addresses recorded while participating in a P2P “hive.”) A majority of the public may only understand P2P technologies to be used in conjunction with illegal music sharing; even this simple, negative perception can greatly influence public opinion on the organization. It would be difficult to refute packet analysis or screen shots containing an organization’s IP address in which the computer was compromised, used for illegal software or media sharing, or the computer was used by an unauthorized entity to extract data. In the age of P2P applications commonly used to illegally share and distribute the intellectual property of unwilling participants, organizations are taking aggressive steps to find and prosecute offenders. See Chapter 55 in this *Handbook* for a discussion of cyber investigations; see Chapter 61 for guidance on working with law enforcement.

35.3.2 Prevention and Mitigation. Protecting the organization from information security breaches via P2P technologies is one of many important parts to an overall security plan. Depending on an organization’s structure, leadership, function, and similar factors, methods for preventing and mitigating P2P threats can range from simple to complicated. Obviously, each organization must perform a risk analysis and determine its threat threshold when it comes to P2P technology. Chapter 62 provides means for risk assessment. The guidelines that follow can help an organization defend against the threat of P2P technology causing security breaches.

35.3.2.1 Policy. It is important for every organization to address the use of P2P technology in a policy, such as an acceptable use policy, HR policies, or security policies. The relevant policy, along with all other security-related policies, should be clearly stated, clearly communicated to the entire organization, uniformly and equally enforced, and updated as necessary.

35.3.2.2 Complete Ban on Peer-to-Peer Technology. In *most* cases, the organization can ban the use of P2P completely, especially through enforceable policy. Care should be taken to ensure all employees and computers are in compliance with the ban. It should be forbidden or, even better, impossible to install P2P applications on personal computers, servers, and all other information systems that could be used to send and receive P2P-related traffic. Most computer users should be running as a standard user of their computer, not as an administrator. If employees are allowed to install software, or are allowed administrative access to their desktop, regular, automated inventories and audits of the computers should take place. Removal should be immediate and appropriate corrective actions taken.

Several technologies may also aid in disallowing P2P traffic, although no technological solution is completely foolproof. These measures are additional safeguards, not complete solutions. Firewalls should be configured with a default deny policy and should only allow TCP/IP ports to pass that are necessary for normal business operations. While many P2P applications are able to tunnel through TCP/IP ports such as those used by HTTP or other common protocols, this is a necessary first defense. Packet-shaping technologies can also be useful to identify P2P-related traffic and block its communications. Packet-shaping and traffic management devices are often able to detect the signature of P2P traffic, no matter what TCP/IP port the application may

35 · 8 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

be using. Some intrusion-detection and intrusion-prevention systems may also be able to identify and block P2P traffic, as would many Internet filtering devices. Logs and reports should be examined daily and infractions should be quickly remedied. This is also a case where a managed desktop environment is important—automated software install/removal, uniform desktop imaging, and automated software inventory/reporting aids the desktop administrators in the effort of combating P2P software installed against policies.

35.3.2.3 Information Security and Information System Audits. All information systems and components should be audited (both on a schedule and via random audits) to ensure that they are not configured, intentionally or unintentionally, to participate in P2P file sharing. This task should be part of every organization’s regular information system and information security audit processes. If possible, external and neutral resources are most useful to ensure all systems are audited in a uniform, exact, repeatable, and objective fashion.

35.3.2.4 Legitimate Business Use Must Be Managed. There are times when an organization does not wish to completely ban or block the use of P2P technologies. One increasingly common and legitimate example for P2P involves open source software distribution or updates via BitTorrent. BitTorrent is a P2P-based protocol for the distribution of data—often large amounts of data. A widespread use includes the distribution of several distributions of the Linux operating system. Installation of Linux often involves obtaining CD-ROM or DVD images to create install discs. By utilizing BitTorrent technology, software vendors and distributors are able to provide large amounts of data to their clients without carrying the entire burden of distribution, bandwidth, and computing resources. However, the organization must manage how this technology is used to ensure resources are not abused and that the P2P applications are used for only allowed, legal ends.

This can be accomplished through policy, auditing, and various network access and control technologies. Each organization must define its own level of acceptable risk for legitimate P2P technology usage and must find solutions that will match the acceptable level. Some examples include:

- Use of encryption
- “Anonymous” P2P routing technologies such as onion routing (see Chapter 31 in this *Handbook*)
- Network isolation for computers used to obtain software with P2P applications
- Company-acquired DSL or cable modem connections to the Internet, avoiding the use of corporate network resources

35.3.3 Response. It is necessary for all organizations to define exactly how to respond to security breaches and policy violations, including situations where P2P technology is involved. Not only should the process be included in the overall security plan, but also incident response processes should be in place to remove offending systems from the network. In some cases, the rebuilding of a compromised resource may be necessary, but some organizations may choose to remove the compromised machine from production and/or preserve a forensic copy of the machine for legal, forensic, or investigative processes.

SECURING INSTANT MESSAGING 35 · 9

35.3.4 Case Study. Misconfiguration, unintentional use, curiosity, and experimentation with P2P in the workplace do happen, with consequences. Although this case is only one type of specific security incident involving P2P technology, it should serve as an example of how such a situation can occur.

An employee of one organization reported a slow-running computer to the helpdesk. All of the usual helpdesk suggestions and tricks were exhausted with little effect on the performance of the computer. The usual symptoms of a slow computer were present—massive wait times to accomplish simple tasks, random errors and shutdowns, lockups, and other operational problems. However, there was one difference: After a reboot, it would take several minutes for the computer to slow down and become unresponsive. After some time, an employee commented, “I did try installing a music sharing program last week, but I didn’t like it and uninstalled it.” This led the engineer to examine each and every process that was running on the computer.

Although it appeared that the P2P application had been uninstalled, it actually had not been; it was still installed and running in a stealth mode. The uninstaller only masked the P2P application. Not only was the P2P application still running, but it was misconfigured to share the entire contents of the C: drive. There were literally thousands of other P2P users attached to the machine actively searching, uploading, downloading, and altering the contents of the computer’s hard drive. The computer was not only giving away all of its data, it was being used for a server to host thousands of media files. Since the computer was on a network segment that had full TCP/IP 1-to-1 network address translation, it was effectively completely open to the outside world—and the outside world was taking full advantage of the opportunity. The hard drive was virtually full, and files were being added and deleted at will by remote users. The host’s firewall was even modified by the P2P application’s install to open all necessary ports to the world.

It is unknown if the user’s personal data was actually accessed, downloaded, or used for any malicious activity, but the capability was certainly there. Because of a user’s unauthorized download, inadequate network security, and other policy violations, the organization could not be sure of the confidentiality or integrity of the computer or its data. Necessary steps were taken to prevent this incident from occurring again, but this scenario has played out at other organizations, and will continue to do so as long as the P2P risk exists.

35.4 SECURING INSTANT MESSAGING. Instant messaging has become an integral part of communications—both business and personal—for many people. It was fairly easy in the recent past to simply create a policy that banned IM for personal use and only allow internal, business-related IM within the company. However, IM has become an integral part of life. IM is pervasive in most organizations—for personal and business use. As a result from executives to interns, IM may be found on many desktops, but it must be managed and secured on all.

35.4.1 Dangers to the Business. With any technology, especially those that make connections to the Internet, there is a risk to the organization. Instant messaging is not a petty annoyance that should be taken lightly; if the technology is not controlled by the organization, a serious breach of security could occur. Remember, IM applications today go well beyond just txt; they are capable of transmitting much more than just interpersonal casual banter.

35 · 10 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

35.4.1.1 Loss of Information. Information loss and loss of confidentiality, intentionally or unintentionally, is most likely the biggest threat of IM to the business. There are several ways information can be harmfully conveyed via IM⁸:

- Revealing secrets via text chat, especially when a company has taken extensive effort to filter and block such communication via email with email gateways, intellectual property data management, and other technologies
- Copy-and-paste functions (including screen captures) used to transmit confidential or secret information from other secured applications or environments
- File transfers
- Screen sharing and real-time collaboration functions such as shared whiteboards or desktop sharing functions
- Forgetting to turn off a voice or video session and unintentionally relaying sound or images to the original correspondent—especially dangerous when leaving voice mail through the IM client
- Relaying voice, video, or both to another party (unintentional or intentional)
- Use of Webcam technology to relay visual information within a secure facility
- Downloading malware to collect and steal data
- Impersonation (This tactic usually involves stealing a known IM account or creating a fake account to impersonate someone the victim knows.)
- Subpoenas or search warrants executed to collect IM logs, conversations, and so on

Although this is not a complete list, it should serve to aid in security planning and policy development. There are many good resources on the Internet to further explain similar threats and consequences, but the preceding list should encourage thoughtful brainstorming about the ways in which an organization may lose data. Some of the listed methods would be extremely difficult to detect and remedy. With high-speed networks and high-speed Internet links at most organizations, a massive amount of data can be conveyed within a small amount of time.

35.4.1.2 Consequences. Like other security threats, the consequences of not securing IM technology can be serious. Many organizations today have experienced a security breach that involved IM, and there are probably more to come. Instant messaging security breaches can be deadly to an organization by themselves or as part of a much larger attack on the business. Stealing or transmitting information through IM is no less risky than any other form of information theft. One single file, whether sent through an IM file transfer or meticulously cut and pasted, bit by bit over a great period of time, can destroy a company's reputation and standing in the public eye or even benefit a competitor. A breach from a single IM conversation has the potential to depress a corporation's stock price in a matter of hours or days. There have even been cases where a chief executive officer's confidential information was captured and posted on the Internet for all to see.⁹

35.4.1.3 Denial of Service. IM cannot be written off as a tiny application with no real footprint on network resources. Instant messaging can be a tool used to create a denial of service attack on an organization, resulting in a loss of availability. IM

SECURING INSTANT MESSAGING 35 · 11

technology can be a powerful and useful tool for an attacker, including the use of IM clients with a direct connection to the Internet. With the right combination of malware and access, an attacker may be able to exploit one of many vulnerabilities discovered in IM applications, including the ever-popular buffer overflow. The National Vulnerability Database listed 38 current vulnerabilities in instant messaging software as of March 2012.¹⁰

35.4.2 Prevention and Mitigation. Every organization must guard against the threats caused by IM technology. Proper review and analysis of the risks associated with IM must be carried out within the organization, and the organization must determine the amount of risk it is willing to take. It is also necessary to evaluate the costs and efforts associated with the prevention and mitigation of this threat. Organizations will judge the risks and rewards of using IM differently. There is no set standard for every organization or business; there is no universal set of rules that can be applied in all situations. (For a discussion of risk assessment and management, see Chapter 62 in this *Handbook*.) The next sections provide strategies, tactics, and considerations for securing IM.

35.4.2.1 Policy. Policy must come first, especially with the popularity and widespread use of IM. Without adequate policies, the organization has no chance of actually protecting itself. Policy must be the foundation that all other considerations rest on. Clearly defined, well-communicated, and equally enforced policy is one of the most important fundamentals information security relies on. No matter what the organization decides when it comes to IM rules, it must be stated in a policy.

Instant messaging, while risky, is one of the most visible policy decisions a business will make for employees. While it might be best, and preferred for best security, to completely disallow IM, which could lead to frustrated employees, unable to use the facility for personal use, business use, or both. Every management team should be conscious of the potential ramifications of an overly strict policy. Conversely, allowing unfettered IM is certainly not the best solution.

Some organizations, depending on their software platform, may be able to provide employees with a managed IM environment where the employee can connect to the organization's official IM platform, as well as some of the popular IM services used for personal communications. At the same time, the organization can centrally deploy settings such as which services are allowed, how they connect (encrypted or not), and what logging options are in use—including if conversations can be recorded or archived.

Compliance and governmental regulations must be taken into consideration. If IM communications are to be allowed, they may become part of the organization's business communications, and therefore may be subject to subpoena, open records requests, and archiving and document retention requirements. This can be especially dangerous if the chosen IM platform is integrated into the organization's voice over IP system. New regulations and legal rules may greatly affect policy decisions. It is important to remember that instant messaging logs and conversations *may* be subject to legal discovery, search, and seizure. Consult counsel for proper legal advice.

35.4.2.2 Effects of a Complete Ban. A ban on all IM technology would be the best way to ensure better enterprise security. However, this will only produce dissatisfied users, without being effective. Users can become technology-savvy in a hurry if they are determined to circumvent a policy. A block on IM communications often

35 · 12 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

causes users to do just about anything they can to accomplish their goal of unobstructed IM. Some clients are designed with this in mind and will happily tunnel out of commonly allowed TCP/IP ports such as 80, usually open for HTTP communications. The software may be configured to bypass firewall rules, detect and avert packet-shaping technology, and tunnel its way to the Internet via encrypted connections. Many IM services also provide Web-only interfaces that do not require software to be installed while communicating via HTTP. Unfortunately, this technology can be a difficult and frustrating one to ban within the organization; a complete ban is probably not a practical solution.

35.4.2.3 Prevent Installation of Instant Messaging Software. Although a complete ban may not be possible, or even desirable, one step that a more secure organization can take is to prevent users from installing IM software. This tactic is not going to solve the whole problem, but it certainly will help. Controlling the installation of IM software should be part of the organization's overall software installation policy. In general, installing software without permission should be denied. If it is feasible, local workstation administrator rights should be denied for most employees. In many user or system management solutions, it is also possible to block software installation through individual, group, or workstation policy templates and procedures. Universal software installation prevention is much easier than trying to define policies or templates for every possible IM peer, application, group, or tool.

35.4.2.4 Fight Technology with Technology. Technology on its own will not provide the organization with an all-in-one solution for securing IM. However, there are a number of network devices, appliances, traffic monitoring software, and other technologies to help an organization minimize IM use. Do not believe marketing claims that any device or technology can guarantee IM blocking; few can deliver on this promise. The only *true* way to guarantee an IM-free company is to block access to the Internet completely—which is not realistic. An organization can also leverage existing security infrastructure such as IDS or IPS.

35.4.2.5 Limit Risk and Exposure. For most organizations, *limiting* IM through policy and technology is the solution to the threats that IM introduces. Combining those two approaches will help to reduce the possibility of data loss. Security managers and administrators should agree on what can and what cannot be allowed within the organization. An organization may choose to block file transfers, Webcam functions, or screen-sharing functions for IM communications. These types of actions will not prevent IM security breaches, but they could limit data loss. As with any policy and risk management, proper audit, reporting, and compliance controls must be in place.

35.4.2.6 Providing Secure Instant Messaging. In environments where IM is needed to run the business, the best strategy is to provide secure, managed IM services to the employees. Of course, the needs will vary among different organizations for different levels of IM connectivity, functions, and software. Many of today's popular corporate email and collaboration systems have built-in or optional IM services. When properly deployed, these IM systems can meet many of these secure IM best practices:

- Encrypt IM communications wherever possible: client to server, server to Internet, and so on.

SECURING INSTANT MESSAGING 35 · 13

- Encrypt logs and chat conversations at the workstation and server.
- Ensure that all logs, chat conversations, file transfers, and archives meet data transmission, retention, and destruction policies.
- Ensure that “presence awareness” features (software features that allow the user to communicate his or her presence or availability, such as “online,” “away,” or “out to lunch” to all users) comply with corporate personnel policies.
- Administratively disable features that cannot be encrypted or properly managed (screen sharing, file transfer, whiteboard, etc.).
- Where possible, lock or force configuration settings to ensure policy compliance.
- Establish procedures for periodic monitoring and auditing of IM systems; do not ignore logs.
- Enforce prudent password policies for IM systems.
- Properly secure IM communication systems with Internet connectivity; consider using proxies or intermediary gateways that protect internal corporate IM systems; ensure appropriate server lock-down policies and procedures.

Corporate-owned and -managed IM systems may not be possible in all situations. In those cases, the organization must form policies and procedures to limit risk and exposure with commercial IM systems. Some systems do provide “secure” IM, but one must be skeptical of exactly how much protection they provide. Consider limiting commercial IM needs to nonessential computers with limited network access, limiting or restricting users to specific IM applications or services, and monitoring instant message network traffic and usage. Some commercial IM services also provide “corporate” or “business” IM services, often for a fee. These premium offerings may provide the organization with the necessary or acceptable level of functionality and security. Access to public IM services should also *always* be blocked from machines with privileged access to critical data assets.

35.4.3 Response. Instant messaging breaches and compromised systems generally do not require special handling after a security incident. In general, normal policies and procedures can be followed to properly investigate, document, and respond to security breaches. There are many commercial tools, including forensic software, to aid in incident response. Infected or compromised systems, if no longer needed for investigation should be reimaged before redeployment to an employee; never allow a machine to be “cleaned” from an incident and returned to production.

35.4.4 Safe Messaging. Although most users at the organization are generally satisfied with mainstream IM systems, clients, and services, there are dangers to be considered. There seems to be almost an unlimited number of open source IM clients, Web-based IM and chat providers, social networking Websites providing IM, and the like. When considering policy and management of IM within the organization, it is important to analyze the source and intentions of all of the possible services. All IM software and services are not created equal; some may originate from untrusted sources and may contain malware and other security risks such as password stealers or keyloggers. Instant messaging software or providers may also be remotely logging information without the user’s knowledge or consent.

Also, if the organization will utilize commercial IM software and services, it is critical to carefully examine the provider’s terms of use and license agreements. The

35 · 14 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

responsibilities and liabilities of both parties should be carefully weighed by information security managers, company executives, and legal counsel before allowing use of the software and associated services.

35.5 SECURING SMS. Few technologies are more ubiquitous than SMS. Virtually all mobile phones are capable of sending and receiving SMS communications. Since mobile phones are virtually everywhere, security considerations must be in place to guard against the threats that they present. A technology with a relatively minor footprint can cause a world of destruction when used as a weapon. Today, SMS technology, and its associated complementary services, has grown exponentially. Securing and defending against SMS must be included in every organization's comprehensive security plan.

To understand SMS security, it is also important to look beyond the cell phone. SMS does not require a cell phone to utilize the technology. Many phone carriers allow SMS messages to be generated and sent from an unsecured, public Website. SMS messages may also originate from email messages, instant messaging services, and the like. SMS messages can even come from subscription services such as a daily horoscope or critical systems like emergency notification services. Beyond this, today's mobile phones, including smart phones, are more powerful and contain many more features and show no sign of slowing down. Phones are increasingly gaining processing power, memory, complex operating systems, and other features that essentially could redefine the device as a personal computer. Phones are able to access the Internet, install applications, communicate from phone to phone, and even access corporate data networks. Information security managers and professionals should never underestimate the power or versatility of a mobile phone. They are a threat to all of an organization's information security.

35.5.1 Dangers to the Business. SMS can introduce many types of security threats into an organization. SMS can cause a data breach by innocent mistakes or by deliberate attacks. This technology can be used as a criminal tool to deliberately steal information, to extract data, to extort information, and to deceive. It may also be a conduit for inadvertent data loss. The consequences of data lost via SMS are relatively the same as any other data breach: loss of confidence in the organization, loss of image, bad public relations, financial penalties, and so on. A serious or even minor data breach may appear to communicate to the world that the organization does not have a comprehensive security plan in effect, or the company does not abide by such a plan—whether true or not. Some investors, customers, or people in the general community may look at a breach of such a simple technology and ask, “How could the company not have proper security for something as simple as a cell phone?” A missing laptop with confidential data is a serious security breach, but a mobile phone, with all its capabilities, must be treated as nearly the same type of critical infraction.

35.5.1.1 SMS as a Tool for Deliberate Data Loss. One danger an organization may face involves an individual or group of people utilizing SMS technology to ferry critical data to unauthorized persons, usually outside the organization. This action would replicate an age-old tactic of stealing information piece by piece from within the organization to someone who should not possess the information. Consider classic tricks of criminals, such as copying information in tiny pieces over great amounts of time to avoid causing suspicion. Any number of technologies can be used to move data, including flash or thumb drives, iPods, scraps of paper, photographs, screen printouts,

SECURING SMS 35 · 15

embedded code, or even memorization. Disgruntled employees may use SMS to send confidential information to an accomplice or even to themselves for later use, such as selling the data, extortion, and the like. It would be virtually impossible to know that an employee is slowly leaking data outside the business from a mobile phone, especially when that phone is not owned or controlled by the organization. What may appear to coworkers as a serious text-messaging addiction may actually be a serious data breach.

Another fact information security management must consider is that SMS service, whether exactly true to the original definition or not, has expanded well beyond messages of only 160 characters. Mobile phone users are able to send real-time video streams, recorded video, photographs, substantially longer text messages beyond 160 characters, Web page links, and just about anything else the phone carriers can implement. If the mobile phone industry considers all of these features to be synonymous with SMS, the organization's security plan should as well. Business risk has increased greatly with every new technology addition.

35.5.1.2 Inadvertent Data Loss via SMS. Data loss can occur by mistake, bad luck, stupidity, misinformed user, or misunderstanding of features as well as by theft of the data device itself. Both deliberate data theft and inadvertent data losses are extremely dangerous, with potentially serious consequences. Search engines reveal many different tactics and war stories of data loss from a mobile phone as well as other SMS-specific security issues. These are scenarios and techniques to consider:

- SMS via email or the Internet
- SMS snooping or sniffing
- Recovery of improperly deleted data
- Stolen, mixed up, or lost phones
- Misdialed numbers
- Wi-Fi connectivity
- Unattended phone with no password
- Malware installed on phone (keyloggers)
- Recipient's phone is lost, stolen, or borrowed
- Impersonation

35.5.2 Prevention and Mitigation. SMS technology is not going to disappear anytime soon, so every organization must come up with a plan to prevent data loss and protect itself from this risk. Once again, the organization's leadership, security management, and security professionals must evaluate the risk of SMS versus the need to operate the business and maintain an amiable group of employees. Every organization must decide for itself exactly what kind of practices to put into place for SMS security and the cost/benefit of each practice. Everything must be considered, from policy and procedures, to deployment of security technologies and mobile phone company-provided services. With today's varying needs, newly emerging technologies and an array of mobile phones, it is difficult for any two organizations to adopt the same prevention and mitigation strategies. However, the next suggestions can be used to begin, update, or enhance the organization's security plan when addressing SMS technology.

35 · 16 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

35.5.2.1 Policy. An SMS policy should be developed for all cell phones that are company owned or company sponsored. The policy, when written on a solid foundation, is key for enforcement by human resources, especially when something may end up in litigation. Vague acceptable use policy will not be enough. A clearly stated position must be written, adopted, and communicated to all employees. The policy should apply to every employee, executive or trainee, with no exceptions. The policy should regularly be reviewed, updated, and redistributed, with recurrent training as necessary, especially in a rapidly changing world such as mobile technology.

A good policy must also address an important distinction common to mobile phone use in the organization: personal phones versus company-provided/sponsored phones. The policy must address: what is acceptable for employee conduct on the job; whether personal phones are allowed on the premises; what type of phone is allowed (usually refers to whether employees are allowed to have smartphones); where, when, and for what purposes can they use personal mobile phones; what is allowed on business-provided phones; and the like.

35.5.2.2 Mobile Phone Ban. In some cases, security needs may necessitate prohibiting the use or even possession of mobile phones on company grounds or in certain areas. This type of action should be included in company policy and should be clearly communicated. It may be necessary to remind employees with signs and repeated communication as well. This is a common practice to prevent data loss from any mobile phone function, including SMS. The organization should be sure to make distinctions for employee-owned phones and emergencies. If an area requires a high amount of security, err on the side of caution, and forbid mobile phones completely. The policy must extend to visitors, vendors, contractors, and other outside entities as well as to every employee—regardless of rank.

35.5.2.3 Providing Secure SMS. Providing “secure” SMS can prove to be difficult, and it can be easy to fall into a false sense of security. Information security managers must know exactly how their mobile phone infrastructure works before declaring the system secure. Although one component of a phone’s connection may be secure—for example, from the phone to its messaging infrastructure—the entire path of an SMS message may not be secure. Some devices, such as the BlackBerry from Research in Motion, provide encrypted transport for messaging. However, email or messages to users on different phone networks or other messaging servers may not be encrypted. Information security professionals must clearly understand the technology they are deploying, and they must test for proper installation and configuration as well as working with the vendor to verify marketing claims. However, it is important to remember that if a solution is not as secure as the organization’s policies and needs require, SMS, mobile phones, or both should be banned. Some phones or smart phone solutions allow administrators to “block out” services such as SMS or to install secure communications software. Carefully consider and evaluate all options and solutions.

35.5.2.4 Bring Your Own Device. There was a time when the corporate standard was the BlackBerry, provided by Research in Motion. While still popular and the most manageable platform for smartphones, a much more recent problem has emerged in today’s organization. Tax law changes have increasingly made company-provided cell phones a thing of the past. Today, the rule is bring-your-own device; more

SECURING SMS 35 · 17

and more often, that device is an unmanaged Apple iPhone or Google Android variant smartphone.

Employees will want to connect their smartphone to the organization's email and calendaring system, as well as utilize the Wi-Fi infrastructure. They will want to stay connected for professional reasons (corporate email keeps them connected to the office and more productive), connected for personal reasons (personal email, social networking), and utilize battery-saving Wi-Fi instead of the cellular data networks. On the flipside, they may want to use their phone to tether another wireless device via Wi-Fi where they cannot or are not allowed to connect to corporate Wi-Fi.

This, combined with the knowledge these devices are unmanaged, have almost unlimited unknown applications, or "apps," and basically are small mobile computers, is enough to make most security pros wince. It is increasingly difficult to tell employees "no" when it comes to them utilizing these devices in the workplace, especially when they are sometimes required to do so for their job. This new challenge must be addressed in the organization's policy for mobile phones/mobile computing. It should be strictly laid out what is allowed, what is not, and how the rules will be enforced. Device management features should be utilized when allowed by the organization's platform to control things like encryption, email sync, data retention, device passwords, and device lock-out.

35.5.2.5 Other Considerations. The next list provides points to consider when planning for SMS security, many of which are from NIST Special Publication 800-48, "Wireless Network Security."¹¹

1. Create policies and procedures to deal with lost mobile phones. The phone may contain sensitive data, including stored and deleted SMS messages.
2. If cell phones are banned from the organization's premises, ensure that physical security has procedures and rules for checking visitors and employees for mobile phones.
3. Many mobile phones have the capability to back up and synchronize their contents to the desktop. Ensure proper procedures to secure data and data leakage.
4. Policies and procedures should be in place to limit and manage the acquisition of mobile phones by employees—information security may not be aware of the existence of new phones in the environment.
5. Mobile phones are not easily audited, nor is there much software to aid in the auditing process.
6. Despite proper labeling of a company-owned device, if lost it will rarely be returned to the organization. Plan to mitigate damage caused by a lost mobile phone; utilize security features such as remote wiping the device after loss via "poison pill" features or "auto-destruct" features after several invalid password attempts.
7. If a mobile device supports screen-lock and power-on passwords, use these simple protections wherever possible.
8. Through policy and education, prevent as much sensitive and private information on the organization's mobile phones as possible, including SMS messages.
9. Utilize Public Key Infrastructure (PKI) technology where possible.
10. Install antivirus software where possible.

35 · 18 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

11. Utilize VPN and firewall technology for safer data communications.
12. If a phone is to be carried on international travel, SMS messaging should be prohibited if at all possible. The risks associated with taking a mobile phone to international destinations increase exponentially.

35.5.3 Reaction and Response. When a security incident involving a mobile phone and SMS does happen, unless the organization has in-house staff trained in mobile incident response or forensics, it may be best to work with the mobile phone provider, possibly also the manufacturer. Procedures and correct processes associated with data retrieval, preservation, investigation, and so on are best handled by those most qualified. If necessary, involve law enforcement. This is an area where a long-standing good relationship with local, state, or federal law enforcement is extremely beneficial—even if the investigation would not necessarily require law enforcement investigation. For more information on this subject, see Chapter 61 of this *Handbook*.

Investigating SMS issues, including tracking messages, tracking phone location, and tracking the path an SMS message took, can often be accomplished with the help of the mobile phone carrier. Law enforcement and court-ordered subpoenas may be necessary, depending on the situation.

Compromised devices should be carefully reviewed before returning them to regular use. Mobile phone providers can assist in “wiping” the device clean of all software, including malware if the device doesn’t have such capability built-in. Specific practices and procedures vary by phone and provider, but some organizations may also choose to archive or destroy devices involved with a security breach of any kind.

35.6 SECURING COLLABORATION TOOLS. Information systems that provide online facilities for collaboration are increasingly valuable business tools. Although these tools provide excellent conduits for increased information sharing, they also have the potential to increase security threats. Even the Internet itself, with many Websites dedicated to information sharing, groupware, shared tools, and data storage, has become a collaboration workspace. New features and movements such as Google Apps, “Web 2.0,” “The Cloud,” and even free or online conference calling services must be taken into account in any organization’s security plan. The nature of collaboration and the need to get critical business done efficiently is critical to most of today’s organizations. Many companies and organizations are trying to get more work done with less people. Technology has become an important partner to allow employees to work together and to accomplish more in less time. Collaboration tools have become even more critical as businesses expand to include people working together from different geographical locations.

35.6.1 Security versus Openness. One of the longtime battles for security managers is security versus openness or functionality. The nature of collaboration requires uninhibited data and information sharing, which can be difficult to secure. Organizations have to find the right balance between allowing users free and open information exchange and providing the required level of security. Finding this balance takes cooperation and respect between the two groups: those who use the tools and those charged with securing the organization. The two groups must fully understand each other’s position; without this understanding, finding a middle ground and negotiating compromise cannot take place. The goal of the organization surely must be efficient, uninterrupted business, but not at the expense of good security. The only way to work

SECURING COLLABORATION TOOLS 35 · 19

through this complication is with good-natured, open, goal-oriented communication. This is *not* an information technology–only problem or process. Finding that optimum balance of security and functionality will require all types of management and staff to work together. Although this may be true of all information security domains, it is especially true of collaboration tools security. Without this important balance, the tools are essentially worthless: too secure and they will not be used, too open and the business could suffer catastrophic data and integrity loss. Some businesses are not able to recover from such a loss.

35.6.2 Dangers of Collaboration Tools. Collaboration tools are becoming powerful, and they must be given full security considerations. These tools should not be installed, accessed online, or integrated into the business without the proper planning, risk analysis, security configurations, and testing; ad hoc, unmanaged systems, installed without the knowledge of security personnel, must be prohibited, and violations corrected. Collaboration tools can easily become a nightmare for security management, especially if securing these tools is not a primary consideration from the beginning. Designing and implementing security measures on an already-deployed production system is invariably a frustrating exercise in futility for both the users and the information security personnel. Likewise, finding out after-the-fact that a cloud-based collaboration tool has become mission critical to the organization is a terrifying thought for most security professionals.

Some of the features and general dangers associated with many of these systems include loss of confidentiality, integrity, or availability. These dangers can occur due to any of these problems:

- **Lack of authentication requirements, rules, or procedures.** A wide-open system or one with poor authentication would allow for unauthorized persons to gain access.
- **Data snooping or capture.** Transmission of data to and from the system could be intercepted by unknown persons.
- **Impersonation.** Proving exactly who the user is may be difficult if not well managed, especially with weak authentication and authorization methods.
- **Unauthorized posting** of confidential information into insecure or public areas.
- **Misconfiguration.** A simple mistake in configuration could reveal private information.
- **Search engines.** Documents or other information may be subject to search engine crawlers/agents/spiders if proper security is not established.
- **Rogue collaboration systems.** If a department or group deploys its own tools, privately or publicly, without the knowledge of the security group, proper security cannot be guaranteed.
- **Internal threats.** One cannot be concerned only with external threats. One department's collaboration system may be another department's limitless temptation.
- **Users.** Users may not always have security in mind. Small mistakes or shortcuts could lead to major security breaches.

When deploying or evaluating collaboration tools, risk analysis must be performed to determine if the organization is able and willing to accept the associated risks. Security groups should thoroughly brainstorm and research as many possible security

35 · 20 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

threats to the collaboration system as possible. It may be beneficial to work with the solution provider's support group to minimize or eliminate as many security risks as possible.

Workgroups utilizing collaboration tools place a great amount of trust in the application and the tools. Many of the applications available today are light on security and heavy on marketable features. Although many online companies have become much more serious about data security, they are not the owners or protectors of the organization's data; that is still up to the organization.

35.6.3 Prevention and Mitigation. Collaboration tools and systems should receive the same security care as any other information system. Although the nature of collaboration may be somewhat open, the same policies, procedures, and careful controls should apply. The goal must still be the confidentiality, integrity, and availability of the data and the information system. Collaboration tools must still benefit the business while ensuring the business will not be harmed by a security incident. By taking the necessary steps to prevent and mitigate security issues, collaboration tools can be invaluable to the organization.

The next suggestions can be utilized to aid an organization in securing collaboration tools and systems.

35.6.3.1 Policy. It is debatable whether collaboration tools necessitate specific, separate policies. What is more important is that a complete, well-written, and well-communicated policy exists, one that includes provisions for collaboration tools, systems, and associated technology. Clear understanding and communication of collaboration tool security must be well researched, well written, concise, well communicated, and updated regularly. It is critical that the policy remain valid as new and more complex collaboration tools are developed and deployed.

Policies should also include security options that may otherwise be out of the control of the organization. For example, if a company forbids using public file-sharing services, the policy should cover users attempting to use the service from outside the organization as well as within it. Employees should not be able to use services or systems that do not comply with the policy, no matter where or how the service is to be used.

A good policy should be inclusive, especially when defining exactly what the organization considers a collaboration tool. It would be easy to forget applications such as email, IM, online meetings, blogs, social networking, shared network resources, remote access software, peer-to-peer file sharing, and the like. Many technologies have collaboration components that must be considered to ensure security.

35.6.3.2 Prevent Access or Use. Another option, in conjunction with policy, is to block the use of collaboration tools, depending on the organization's needs. This may involve deploying technology to accomplish this goal, including content blocking, firewalls, or both. This should disallow installation or use of rogue collaboration tools. Periodic review of networked systems and network traffic should be conducted to ensure compliance with prevention or limitation of collaboration tools.

35.6.3.3 Limit Access. Many collaboration tools can be deployed as an internal-only system, a cloud-based system, or both. Organizations will want to choose how users will access these systems. For example, disallowing unsecured communications from the Internet may help increase security. Likewise, it may be necessary to block

SECURING COLLABORATION TOOLS 35 · 21

access to public services from within the organization’s network. Or technical solutions, such as VPN connections from outside the organization’s network, may be used to meet communication needs.

35.6.3.4 Deploy or Enhance Security Frameworks and Technologies.

Wherever possible, install solutions that will increase collaboration tool security and that can be integrated into existing security frameworks. If the organization has a high-security, single sign-on solution, integrate the collaboration systems into it. Another example would be to integrate the collaboration systems into a new or existing PKI infrastructure. Utilize well-known and reliable solutions such as Secure Sockets Layer (SSL) and encryption for the host and all participants. This greatly reduces the risk of security breaches during data transmission.

35.6.3.5 Audit. No matter what level of policy, procedures, or preventions are put into place, every organization *must* audit for compliance. Procedures for auditing collaboration tools and their use should be included in the organization’s regular, structured, information security auditing functions. Any deviations from the policies and procedures mandated for collaboration tools must be acted on in a timely manner.

35.6.3.6 Monitoring. Any organization that deploys collaboration tools must monitor and report on the system’s usage, audit results, and data contents. (The organization must examine the actual data contents to ensure compliance with protections such as protected health information [PHI] or Social Security number [SSN]. Many new products have rules written for this reason.) Monitoring and reporting work to ensure that collaboration tools and systems are being used for their intended purposes. Monitoring and reporting of active projects should look for unusual patterns of use, policy violations, inactive users, inactive or outdated systems, and the like. Proper system management should already be in place, but it is important to check the systems periodically. For example, if a group is utilizing a collaboration system for a project, once the project has been completed, all project materials and users should be removed from the system. Reports from system monitoring and auditing should be acted on at once.

35.6.3.7 Consider Outsourcing and the Cloud Carefully. Some organizations are tempted to use commercial online-only collaboration systems or hosted solutions. This decision should not be made lightly; consider the risks versus the returns. The organization should carefully review all terms of service, license agreements, service-level agreements, and legal responsibilities carefully. Legal counsel must be involved to ensure the organization is protected, especially in the area of data ownership, possession, legal discovery, and subpoena power.¹²

The cloud has gotten a *lot* of attention lately. While cloud security is far beyond the scope of this chapter, the cloud should be treated no differently than any other outsourced application. The cloud often contains the same risks as any other service where the data is stored outside of the control of the organization. Security should be a primary concern when evaluating solutions in the cloud, with consideration given to:

- Service level agreements
- Data ownership

35 · 22 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

- Data security technology in use while data is at rest (encryption, backup, replication, etc)
- Data recovery if the provider is sold, closes, or goes bankrupt
- Physical location of the data and what national, state or local laws apply to that location and data
- Administrator access (for the customer and which administrators at the service provider also have access)
- Liability of data breaches at the service provider
- Key escrow for encrypted data

While this is not an exhaustive list, this should help an organization in brainstorming and researching the risks of cloud-based collaboration tools. There are a lot of online resources from reputable organizations that have already provided guidance for organizations wading into the cloud. As usual, consult legal counsel.

35.6.3.8 Audits and Penetration Testing. As with most information systems, providing necessary security should involve regular, external, third-party penetration testing and audits. Collaboration tools and associated systems should be tested and evaluated for their security fitness. Any problems discovered should be documented and swiftly remedied. Allowing a neutral, external entity to test the system independently is superior to internal testing, so that bias can be ruled out.

35.6.3.9 Keep Collaboration Tools Current. Keeping collaboration tools and their associated information systems up to date is critically important. Applying patches for vulnerabilities is good information technology and information security best practice. After thoroughly testing patches in a test environment, they should be applied to production environments as soon as possible. Do not ignore software vendor patches, especially those for known vulnerabilities.

35.6.4 Reaction and Response. Once a security breach has been discovered involving collaboration tools, the organization's usual policies and procedures should be followed. Procedures for compromised information systems should be well-formed, repeatable processes to preserve evidence, provide for rapid discovery and investigation, and meet necessary regulatory guidelines. When necessary, law enforcement, legal advisors, or both should be utilized to ensure proper evidence collection and documentation. Policies should also dictate the procedures for post-investigation tasks as well, such as requiring compromised systems to be copied, archived, destroyed, reimaged, or reinstalled. It is generally not advisable to try simply to "clean up" a compromised system. It can be difficult to guarantee that a compromised system once again has integrity.

35.7 CONCLUSIONS. This chapter introduces security managers and professionals to securing peer-to-peer technologies, instant messaging, short messaging services, and collaboration tools. The suggestions and information in this chapter are meant to aid in making decisions regarding these tools within the organization's overall security plan. Many of the examples and concepts are meant to aid in the planning, policy development, and review of the organization's exposure to these technologies and their dangers. This chapter should serve as only a starting point for the

NOTES 35 · 23

organization's research on each topic and to ensure that information security managers at least have a brief understanding of each concept, its risks, prevention and mitigation strategies, and suggestions for response. It is difficult to recommend solutions for every type of business, so each organization must make its own judgment for securing these technologies. The popularity and ubiquity of P2P, IM, SMS, and collaboration tools ensures that they will be part of every security plan for many years to come.

35.8 FURTHER READING

- Covill, Jared J. *Going Google: Powerful Tools for 21st Century Learning*. Corwin, 2012.
- Flynn, Nancy. *Instant Messaging Rules: A Business Guide to Managing Policies, Security, and Legal Issues for Safe IM Communication*. AMACON, 2004.
- Kunz, T., and S. S. Ravi, eds. "Ad-Hoc, Mobile, and Wireless Networks." 5th International Conference, ADHOC-NOW 2006, Ottawa, Canada, August 17–19; 2006 Proceedings. New York: Springer, 2007.
- Lamont, Ian. *Google Drive & Docs in 30 Minutes*. Digital Media Machine, 2012.
- Piccard, P., B. Baskin, G. Spillman, and M. Sachs. *Securing IM and P2P Applications for the Enterprise*. Norwell, MA: Syngress, 2005.
- Rittinghouse, J., and J. F. Ransome. *IM Instant Messaging Security*. Burlington, MA: Elsevier/Digital Press, 2005.
- Rittinghouse, J., and James Ransome. *Cloud Computing: Implementation, Management, and Security*. CRC Press, 2009.
- Rhoton, John, Jan De Clercq, and David Graves. *Cloud Computing Protected: Security Assessment Handbook*. Recursive, 2013.
- Taylor, I. J., and A. Harrison. *From P2P to Web Services and Grids: Peers in a Client/Server World*. New York: Springer, 2004.
- Winkler, J. R. *Securing the Cloud: Cloud Computer Security Techniques and Tactics*. Norwell, MA: Syngress, 2011.

35.9 NOTES

1. Dropbox, "How Secure is Dropbox?" Dropbox Website. 2013, <https://www.dropbox.com/help/27/en>
2. Michael Kassner, "DropSmack: Using Dropbox to Steal Files and Deliver Malware," *TechRepublic | Security*, April 15, 2013, www.techrepublic.com/blog/security/dropsmack-using-dropbox-to-steal-files-and-deliver-malware/9332
3. M. E. Kabay, "Dropping the Ball on Dropbox," *InfoSec Perception*, April 19, 2013, <http://resources.infosecskills.com/perception/dropping-the-ball-on-dropbox/>
4. Matthew J. Schwartz, "5 Dropbox Security Warnings for Businesses." *InformationWeek/Security*, August 14, 2012, www.informationweek.com/security/management/5-dropbox-security-warnings-for-business/240005413?pgno=1
5. Cloudfogger, "Protect Your Privacy on SkyDrive, Dropbox, Google Drive in the Cloud," *cloudfogger*, 2012, www.cloudfogger.com/en
6. It should be noted, that at the time this chapter was updated, it is not uncommon for a house or small organization to have an Internet connection of 40–50 mb/sec. While the resource-hogging nature of P2P is becoming less and less of an issue, that should not negate the need to have a policy to limit this technology in the organization.

35 · 24 SECURING P2P, IM, SMS, AND COLLABORATION TOOLS

7. J. Borland, “‘Spyware’ Piggybacks on Napster Rivals,” *CNET News.com*, May 14, 2001, <http://news.cnet.com/2100-1023-257592.html>
8. N. Hindocha, “Instant Insecurity: Security Issues of Instant Messaging,” *Security-Focus*, January 13, 2003, www.securityfocus.com/infocus/1657
9. P. Festa, “ICQ Logs Spark Corporate Nightmare,” *CNET News.com*, March 15, 2001, <http://news.cnet.com/2100-1023-254173.html?legacy=cnet>
10. National Vulnerability Database, <http://nvd.nist.gov/nvd.cfm>
11. K. Scarfone and D. Dicoi, *Wireless Network Security for IEEE 802.11 a/b/g and Bluetooth (DRAFT)*, NIST Special Publication 800-48 Revision 1 (Draft), 2007, <http://csrc.nist.gov/publications/drafts/800-48-rev1/Draft-SP800-48r1.pdf> (URL inactive).
12. M. Rasch, “Don’t Be Evil,” *SecurityFocus*, 2007, www.securityfocus.com/print/columnists/447 (URL inactive).

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 36

SECURING STORED DATA

David J. Johnson, Nicholas Takacs, Jennifer Hadley, and M. E. Kabay

36.1 INTRODUCTION TO SECURING STORED DATA	36·1	36.3.2 Trusted Hosts	36·9
36.1.1 Security Basics for Storage Administrators	36·2	36.3.3 Buffer Overflows	36·9
36.1.2 Best Practices	36·2		
36.1.3 DAS, NAS, and SAN	36·4		
36.1.4 Out-of-Band and In-Band Storage Management	36·4	36.4 CIFS EXPLOITS	36·9
36.1.5 File System Access Controls	36·5	36.4.1 Authentication	36·10
36.1.6 Backup and Restore System Controls	36·5	36.4.2 Rogue or Counterfeit Hosts	36·10
36.1.7 Protecting Management Interfaces	36·6		
		36.5 ENCRYPTION	36·10
		36.5.1 Recoverability	36·11
		36.5.2 File Encryption	36·11
		36.5.3 Volume Encryption and Encrypting File Systems	36·11
		36.5.4 Full Disk Encryption	36·12
		36.5.5 Vulnerability of Volume, File System, and Full Disk Encryption	36·12
		36.5.6 Database Encryption	36·13
		36.5.7 Smart Phone Encryption	36·15
36.2 FIBRE CHANNEL WEAKNESS AND EXPLOITS	36·7	36.6 DATA DISPOSAL	36·15
36.2.1 Man-in-the-Middle Attacks	36·7		
36.2.2 Session Hijacking	36·8		
36.2.3 Name-Server Corruption	36·8	36.7 CONCLUDING REMARKS	36·16
36.2.4 Fibre-Channel Security	36·8		
		36.8 FURTHER READING	36·16
36.3 NFS WEAKNESSES AND EXPLOITS	36·8		
36.3.1 User and File Permissions	36·9	36.9 NOTES	36·17

36.1 INTRODUCTION TO SECURING STORED DATA. This chapter reviews methods of securing data stored on nonvolatile media. Nonvolatile media include magnetic disks and their (hard) drives, compact discs (CDs), and digital video disks (DVDs), with their optical drives), and flash drives (also known as *USB drives*, *flash*

36 · 2 SECURING STORED DATA

disks, and memory keys). Volatile storage devices, which are not covered in this chapter, include random access memory (RAM) and other storage that loses its contents with a power loss.

One of the most important issues for security management is the frequent loss or theft of laptop computers and other mobile computing devices. For example,

- A 2008 report suggested that in the United States alone, over 600,000 laptop computers were left behind at U.S. airports every year, with about two-thirds of them left unclaimed.¹
- In 2010, “... 275 businesses in Europe... lost a combined 72,000 laptops, costing those organizations a total of \$1.8 billion. Of the 265 laptops lost by the average organization per year, a company will typically only recover 12.”²
- In 2012, the rate of loss was still 12,000 laptops lost per week at U.S. airports; another detail of reports at that time was that the unclaimed laptop computers were being auctioned to anyone able to pay for them.³

Unencrypted files on these stolen or repurchased laptop computers are an invitation to violations of confidentiality and control. For example, an audit of discarded computers from the New Jersey state government revealed that “Taxpayers’ Social Security numbers, confidential child abuse reports, and personnel reviews of New Jersey workers nearly went to the highest bidder after the state sent surplus computers out for auction.” In addition, “Nearly 80 percent of discarded computers in a comptroller’s office sample had not been scrubbed of data before being shipped to a warehouse. ...”⁴

Since the last edition of this *Handbook* was published in 2008, cloud-based storage solutions are increasingly being used by individuals and by organizations, so this chapter includes a short discussion of security for such repositories.

36.1.1 Security Basics for Storage Administrators. Storage systems have developed outside the security umbrella of other organizational assets. Because the storage arena is one of the most strategic parts of the infrastructure, professionals should take the same care in developing comprehensive security controls as those that are addressed for the remainder of the network. A number of vendors have focused on the development of secure storage environments that are scalable yet flexible; most address both the logical and physical aspects of security. However, any appropriate strategy for data storage protection includes a balance between protecting the confidentiality and integrity of the information while also ensuring its availability and utility to the system and to authorized users. Ultimately, those with a responsibility for data storage will also be tasked with maintaining this balance at a reasonable cost.

Knowing how and where data will be stored on a network and addressing the known risks to the data is the best option, as it is often more efficient in terms of resource usage (time and money). An organization need not be a high-profile entity to suffer from a compromised pool of data. A single backup can contain enough concentrated personal or sensitive corporate information to experience a loss of credibility, to lose revenue, and possibly to bring the organization to its knees. Worse, a backup that is copied illicitly may show no signs of having led to loss of control over confidential data.

36.1.2 Best Practices. Every organization must keep its applications, servers, and end-user systems up and running to make use of information and to maintain the highest degree of information availability and integrity. A tiered data protection model

INTRODUCTION TO SECURING STORED DATA 36 · 3

works best; it includes a layered defense, due diligence, and restricted management. Implemented correctly, security should be transparent. Best practices for providing a secure data storage environment include the following:

- Performing an audit and risk assessment on the storage infrastructure, looking for risks and vulnerabilities.
- Implementing authentication across the storage network that could coordinate authorization, password maintenance, and encryption.
- Implementing strong role-based access controls and assigning access rights to parties on a need-to-know basis.
- Adopting and enforcing data-encryption and data-classification policies. Based on the classification level assigned to the data, the organization's policy may require the encryption of the data at rest throughout the lifecycle of the data. There may also be requirements to encrypt the data in flight (across the network) as well.
- Requiring strong security features and practices from storage-system vendors and offsite-storage providers.
- Remembering to secure the Storage Area Network (SAN) at the switch (or fabric) level. Carving up the fabric by zones is one technique that limits access to various parts of the SAN.
- Including any data-replication or storage-replication technologies in the overall storage security plan, where such replication traffic may include transaction journals and other temporary locations with partial or full copies of sensitive data.
- Evaluating and implementing data-loss prevention (DLP) technologies to maintain the security of stored data outside the enterprise. Data tagging, or labeling, along with pattern matching (e.g. 999-99-9999 for a Social Security number) are key for assisting with DLP.
- Creating a policy for discarding old devices and media and data-storage services, to include routinely performing tasks such as secure wiping or physically destroying all data-storage devices and media.
- Evaluating and implementing retention and data-destruction policies, ensuring compliance with any applicable organizational or government regulatory issues.
- Isolating the storage management network from the organization's primary network. By not isolating the network, every employee potentially has access to the stored data.
- Establishing access-log monitoring.
- Performing employee and contractor background checks as part of the Human Resources (HR) hiring procedures.
- Establishing facilities controls in the organization to restrict physical access to data centers, locking storage cabinets and server racks, using locks built into some servers, and ensuring the reliability of the perimeter and building(s) and verifying that such controls are in place for remote storage locations.
- Treating the security of backups as warranting monitoring and alarms at a high level of importance. Adopting secure media management tracking and handling policies that include backup requirements for financial information, employee data, and intellectual property.⁵ Chapter 57 in this *Handbook* contains much information about data backup.

36 · 4 SECURING STORED DATA

36.1.3 DAS, NAS, and SAN. There are three primary methods for storing data: Direct Attached Storage (DAS), Network Attached Storage (NAS), and Storage Area Networks (SANs). In addition, remote data storage using cloud computing services are growing in popularity, especially with individual users but even with organizations.

Direct Attached Storage drives are those that are connected directly to the computer. DAS can either be internal, contained within the computer's case, or external and attached via a peripheral component interconnect, PCI, eSATA, or other bus channel. The risks to DAS devices are either their physical theft or access through the computer system that they service.

Network Attached Storage devices are specialized servers that run minimized operating systems and file systems designed specifically to support input/output (I/O) from other servers. The servers that attach to the NAS devices have DAS that contains their operating systems, applications, and other components but, normally, write all data to the NAS device via TCP/IP over Ethernet connections. NAS is utilized via file sharing protocol such as Network File System (NFS) for UNIX systems and Server Message Block (SMB) or Common Internet File System (CIFS) for Microsoft systems. (As CIFS grew out of SMB, the two are often noted as SMB/CIFS or CIFS/SMB.) As with any Ethernet connection, connections between a system and the NAS server that it uses are subject to being sniffed, to eavesdropping, and to packet capture. NFS and CIFS threats are discussed later in this chapter.

Storage Area Networks are collections of centralized disks that can be accessed by numerous servers. Using SANs can facilitate company growth, as most SANs options allow additional disks to be added to the pool as data storage needs increase, without having to take the attached systems offline, as would need to occur if new DAS were being added to individual systems. Data backups can also be easier to control, as a single storage resource could potentially be backed-up instead of each individual system. With the implementation of RAID or other disk redundancy implementations, writing data to multiple disks can be accomplished without impacting the application servers' performance. Systems can be attached to the SAN by various methods including TCP/IP and fibre channels. Fibre channels are discussed later in this chapter. Using IP for SANs connections enables servers to connect over the Internet—but this option must be used with caution, due to the security concerns of transferring data over the Internet.

Cloud backups use Internet connectivity to store backup data on remote systems out of the control of the data owners.

36.1.4 Out-of-Band and In-Band Storage Management. Managers may have to control storage locations remotely, that is, from a location other than a directly attached console. There are two approaches to such control communications, each with its own particular security issues.

In-band management uses the same network as the data transfers; *out-of-band management* uses a separate network.⁶

While in-band storage management uses the same channels as the data itself traverses for storage, out-of-band management uses alternative methods. For example, an out-of-band solution might have a storage administrator working from a desk and connecting to the storage system over the primary network used by all employees, while the data traverses a dedicated channel between the application server and the storage system.

With out-of-band management, consideration must be given to how to ensure that only authorized systems, such as the administrator's, are connecting to the storage system. This is especially necessary if the storage system does not require authenticated

INTRODUCTION TO SECURING STORED DATA 36 · 5

connections being established before a command is accepted. Without authentication, any system able to communicate with the storage system could issue commands that would negatively impact the storage system. Another risk is due to the interface used by management communications and the commands being sent across the wire without being encrypted. For storage systems that are managed by HTTP interfaces by default, it may be possible to use HTTPS instead, in order to mitigate the risk of commands and logins from network packet capture.

In-band management also has concerns. Commands sent in-band are normally sent in clear text. Other threats of in-band storage management include⁷:

- Management interfaces being subjected to denial-of-service (DoS) attacks,
- Commands providing information on other devices and controllers, and
- Set and reset commands being issued inappropriately.

36.1.5 File System Access Controls. File systems provide access control to data. UNIX file systems provide controls based on the user owner, the group owner, and “other,” or those that are not the user or a member of the group that owns the data. Microsoft® Windows systems allow for data owners to be specified and access granted by either individual usernames or group accounts. Access control lists (ACLs) can also be used to provide access exceptions to the normal access permissions of the data files.

When correctly applied, these access controls can be effective in preventing unauthorized access to data through normal usage. However, the file system trusts the computer’s operating-system access controls to have correctly authenticated and authorized the user. If the access controls for the operating system are circumvented, then file system access controls lose their effectiveness.

For more details of operating system and local area network access controls see Chapters 24 and 25 in this *Handbook*.

36.1.6 Backup and Restore System Controls. Systems used for the backup and restoration of data need extra security consideration because the data contained on them are often critical to disaster recovery and business continuity. There are several threats to backup data that are not faced by attacks against other data stores. While most systems only have their data stored on disk (DAS, NAS, or SAN), backup systems often write data to tape cartridges or other removable media specifically intended for off-site storage. Another option is to backup data electronically to a remote system, either at another of the organization’s data centers or with an off-site backup vendor. As with any transmission of data to remote facilities, the data must be protected during transit and at the destination facility.

Regardless of the media used, all data backups need to be stored at a secure, environmentally protected facility sufficiently distant from the originating location to mitigate the risk of losing both primary and backup data from a single major event, such as an earthquake, flood, or volcanic eruption. Additionally, the backup media storage needs to be secured by restricting access to authorized personnel only. Media used for backups needs to be evaluated to ensure that it meets or exceeds the longevity needs of the data, as defined by the organization’s backup policies and standards.

Additional risk from system imitation also needs to be considered. If an attacker is able to insert a system that impersonates a backup system, all of the data intended to be backed-up may instead be written directly to the attacker’s system. Conversely, if an attacker is able to insert a system that can masquerade as one or more data storage

36 · 6 SECURING STORED DATA

systems, then the attacker could request a restoration of data from existing backups, gaining unauthorized access to the information.

To mitigate such risks, interactions between data storage systems and backup systems should be authenticated. For manually controlled backups, systems could have backup accounts created that require an interactive login to authenticate the backup request. For automated backups, other options may be available, including the use of client and server certificates to authenticate both systems involved in the backup connection.

In recent years, developments in cloud computing technology and cloud-based storage have enabled organizations to leverage hybrid data-backup solutions. Instead of backing up data at regularly scheduled intervals, changes to systems replicate to a cloud-based service in near real-time, effectively acting as a replicated copy. The addition of de-duplication technology, sending only those blocks of data with changes, further improves the efficiency of the solution. However, as with all remote backup solutions, continual due diligence with the service provider is necessary to ensure that appropriate security controls exist for the transmission, storage, access, and retrieval of the data. This is especially important given the “on demand” nature of a cloud or hybrid-based solution for data or system recovery. Any use of cloud services should be done only after:

- A thorough assessment of the cloud service provider’s (CSP) security practices,
- Validation that storage of the data in the solution is in compliance with all legal, regulatory, and company requirements, and
- Employing data encryption that maintains in-house management of the data encryption keys.

Off-site storage of data, whether written directly to a storage vendor’s server or on removable media such as tape, deserves its own security considerations. For example, how does the vendor secure physical access to its site? How do they vet their employees? Most risks such as these can be mitigated by performing due diligence of the vendor prior to entering a contract, and by eliminating prospective vendors that do not meet the organization’s security needs.

Another mitigation strategy is to encrypt data backups before they are sent offsite. Many backup applications offer encryption methods. The use of encryption for files and other data storage systems is discussed later in this chapter. Similar techniques can be applied to backup data by encrypting it as it is written to the backup media.

For more details of backup strategies and security, including cloud backups, see Chapter 57 in this *Handbook*.

36.1.7 Protecting Management Interfaces. Management interfaces pose one of the greatest threats to the security of stored data. These interfaces provide administrative access to the data stores, allowing individuals with appropriate access the ability to manipulate data elements, update account security, and perform other housekeeping activities. Therefore, care is required when implementing a storage solution to ensure that a well-defined defense in depth exists. While two-factor authentication is more secure, it is not always practical. At a bare minimum, each administrative user should have a set of credentials, and complex password requirements, with a regular password change frequency. For situations where the data storage occurs via an Internet-based source, certificate or token-based authentication can provide an extra

FIBRE CHANNEL WEAKNESS AND EXPLOITS 36 · 7

layer of authentication security. In addition, the individuals responsible for administering security should not be the same individuals responsible for managing the storage environment. This separation of duties limits the ability of any one individual to circumvent security controls, without some type of conspiracy between the storage and security administrators.

Another important component of the defense in depth strategy involves the use of audit logs. Logging should be enabled to detect violations of policy and of prescribed procedures. However, logs are valueless unless subjected to regular and random review, with follow-up if anomalies are detected. It is unrealistic to expect an individual to pore over voluminous log files on a daily basis. However, log aggregation and correlation technology can be employed to provide an additional layer of confidence. Regardless of the final implementation, the use of audit logs, and restrictions on the ability to access and modify those logs, plays an important part in guaranteeing that no data corruption occurred.

Mitigating these risks to the management interface requires careful monitoring and control of who can access and install these interfaces. With the movement to Web-based interfaces, this discussion comes down to strict user access control and authentication. In any case, it is imperative that only trusted users with a need-to-know be allowed access to the management interface. Once inside, individuals can manipulate the environment as needed to support their goals, whether to further organizational objectives or to perform malicious actions.

36.2 FIBRE CHANNEL WEAKNESS AND EXPLOITS. Fibre channels, while very economical, present unique challenges to the storage environment. The term *Fibre channel* refers not just to an optical-fiber-based communication pathway, but rather to a complex communication protocol. There are a number of inherent weaknesses in the technology, some of which are fairly straightforward and manageable, but others introduce questions of the viability of the technology in larger storage environments.

One of the most serious security weaknesses with Fibre channel is that all communications occur in cleartext. Fortunately, when Fibre channel implementations occur completely within a data center or other secured area, this is not a major concern, as the ability of unauthorized individuals to intercept traffic *on the wire* is limited unless they have physical access to the cabling and can avoid detection as they attempt to intercept traffic or to install interception devices. In early Fibre-channel implementation, native authentication and encryption was not available. However, newer fabric switches and SAN devices provide built-in capabilities for authentication and encryption to significantly reduce the risk posed by cleartext data transmission.

From a vulnerability perspective, attackers can use the Internet Protocol (IP) to craft exploits against Fibre channel, since both protocols use a frame-based communication scheme. Unfortunately, based on the cleartext issue discussed above, an attacker could sniff frames from the Fibre channel connection, and gain information needed to craft an attack. This section focuses on three types of common attacks: *man in the middle*, *session hijacking*, and *name-server corruption*. Most attacks on a storage network require physical access to that environment, or access to the appropriate sniffing hardware, increasing the difficulty of successful, undetected attacks.

36.2.1 Man-in-the-Middle Attacks. Man-in-the-middle attacks (MIMAs) take advantage of weaknesses in the frame-based communications through Fibre channel. Like IP-based attacks, MIMAs involve an attacker intercepting communications, stealing or changing data, and passing that frame on to the intended destination. Fibre

36 · 8 SECURING STORED DATA

channel includes a *Sequence ID* and *Sequence Count*, both intended to ensure consistent communication from sender to receiver. Much like IP, the Sequence Count is an easily predictable sequential number, allowing an attacker to anticipate the next sequence number and forward a packet ahead of the sending system. The introduced packet allows the attacker to intercept the stream without the authorization of either party. Mitigating this issue requires data-integrity checking to guarantee reception of the correct information and rejection of the fraudulent packets.

36.2.2 Session Hijacking. Session hijacking presents the same type of problem as MIMAs, and occurs in much the same way. However, this type of attack focuses on the lack of authentication in Fibre-channel environments. Instead manipulating data in each frame and passing it on, the hacker uses knowledge of the Sequence ID and Sequence Count to intercept and control the session, making the recipient believe that the attacker is really the original sender. The newly controlled session can then be used to extract whatever data or other information the attacker wishes. Mitigating this issue requires strong authentication to provide a guarantee that the original sender is still the same system throughout the length of the connection.

36.2.3 Name-Server Corruption. The last type of attack in this section involves address spoofing, similar to DNS spoofing in the IP world. Each Fibre-channel connection registers its name with World Wide Name (WWN) service, through two processes, a *Fabric Login* (FLOGI) and a *Port Login* (PLOGI). Typically, name-server corruption occurs during the PLOGI process by allowing an incorrect host to register itself with the Fibre channel switch (which contains the WWN service) using a spoofed address. The switch registers the host under that address as if it were valid because of the lack of any host authentication. When the real host tries to connect, the switch denies that connection because the incorrect host is already connected. This type of attack requires some timing on the attacker's part, but can be easy to accomplish because of the weaknesses previously discussed. Modern fabric switches provide additional authentication mechanisms to defend against these attacks through validation of the device with the switch at regular intervals.

36.2.4 Fibre-Channel Security. This short examination of Fibre channels has pointed out that many weaknesses exist. However, this does not mean that Fibre channel should be dismissed as a suitable technology. When implementing Fibre channel into an environment, care must be taken to address these vulnerabilities, taking into consideration location, distance, and availability of the implementation to individuals, systems, and other devices on the network. The physical placement of the devices and wiring must also be optimized to mitigate the risks associated with using Fibre channels. Vendors are also answering the call by offering technology to help secure Fibre channel through built-in authentication and encryption capabilities.

36.3 NFS WEAKNESSES AND EXPLOITS. Network File Systems (NFS) provide a service allowing a user on a client machine to access network-based resources as if they were local to that user. This service is built upon the Remote Procedure Call (RPC) service, and although very useful in a networked computing environment, NFS presents a number of security issues that must be addressed prior to implementation. NFS is typically geared toward high-bandwidth environments, such as an LAN, or networks sharing nonsensitive information. Since NFS does not provide encryption

CIFS EXPLOITS 36 · 9

between hosts, using this technology for other networks, especially those exposed to the Internet, introduces additional risk.

This section describes three of the most common weaknesses and exploits for NFS: user and file permissions, trusted hosts, and buffer overflows.

36.3.1 User and File Permissions. Aside from the lack of encryption, NFS allows user access based on the particular host connected to the NFS share. This means that any user connected to that host can access the network resources. Restricting users to read-only access eliminates the potential to use NFS as a collaborative technology, because users can no longer create or update information on the shares.

When mounting shares with read-write access under NFS, any user connected to the host can access another user's files, as the only protection the file has is its permissions. Administrators attempt to mitigate this risk by forcing all users to access the share under a group or common set of read-only credentials, but this approach eliminates some of the benefits that a network share provides. The read-only share then requires administrators to update or edit files, providing a library-style approach rather than a collaborative file management environment.

36.3.2 Trusted Hosts. Problems with NFS specifically concern the authentication of hosts to the NFS environment. Because NFS controls mount requests based on the host connecting, and not on the particular user, a rogue host could request an NFS mount and make changes to resources. An attacker could also compromise a DNS server used by the system exporting the NFS to point that system to an unauthorized machine. Because there are no login credentials shared prior to mounting NFS shares, if the hosts are not trusted, there are no additional checks to validate the integrity of the new host.

36.3.3 Buffer Overflows. In many implementations of NFS, data input checking does not occur before processing a request. This flaw presents an opportunity for a buffer-overflow attack. When a directory-removal request comes to an NFS server from a user with read-write privileges, the server does not check the length of the path name, and the user can include additional instructions beyond what the server should receive. Those instructions, presumably malicious, could then be executed by the server as an administrative account, such as *root* or *administrator*. As a result, unintended or unauthorized data manipulation can occur.

Recent implementations of NFS have included Kerberos authentication to help validate the users and what they are able to do. In addition, the same validation can be used to validate hosts before they connect to the NFS server. However, these improvements are only partial solutions. Criminal hackers continue to develop and apply new buffer overflows, and it would be unwise to assume that NFS, or any other network technology, can be completely secured.

For more detailed discussions of secure programming techniques and software quality assurance to preclude buffer overflows and identify and prevent other software vulnerabilities, see Chapter 38 and 39 in this *Handbook*.

36.4 CIFS EXPLOITS. The Common Internet File System (CIFS), an Internet-enabled Server Message Block (SMB) protocol, builds upon that protocol by including encryption and secure authentication to the existing resource sharing capabilities of SMB. Unfortunately, from a security perspective CIFS blends some of the new with some of the old, and the result includes a number of security issues.

36 · 10 SECURING STORED DATA

36.4.1 Authentication. CIFS implementations provide either a straight password-based authentication scheme or a challenge-response scheme. Both of these approaches occur in cleartext, allowing anyone with wire access to intercept and capture authentication credentials to the network share. Even with the challenge-response approach, attackers could spoof a transaction and gain access to the share. Recent implementations of CIFS rely on Kerberos for authentication, and much like NFS, while providing additional security, introduces Kerberos-based vulnerabilities, which are outside the scope of this section. Some CIFS implementations also provide a “share-level” security model, rather than a “user-level” security model. In essence, rather than each user maintaining individual credentials, the share has only one set of credentials, which all users share. The weaknesses inherent to a share-level model are similar to those found with the use of group or shared accounts on any other system.

In addition to the authentication issues, CIFS is also vulnerable to dictionary and brute-force attacks against a user’s credentials. These generally involve a chosen plain-text attack, helped by intercepting challenge-response pairs during the authentication process. However, both online and offline dictionary attacks are available to the attacker based upon the amount of time available to watch the connection attempts.

36.4.2 Rogue or Counterfeit Hosts. It is important to identify the differences between the attack surface of CIFS and that of NFS. Man-in-the-middle and the trusted host issue both apply to CIFS, with some different concerns. Improperly configured CIFS clients can be fooled into thinking that they should supply a password instead of interacting with a challenge-response scenario, thus supporting man-in-the-middle attacks. Additionally, if a CIFS environment that does not enable session or message authentication, it removes security controls designed specifically to protect against these types of attacks.

CIFS shares many of the same security issues with NFS. However, most of the issues can be avoided by enabling security features included with the protocol. Man-in-the-middle and session hijacking attacks, in addition to replay and spoofing attacks, can be avoided by message and session authentication. This does not suggest that CIFS can be completely secured, but that care must be taken during the configuration step of implementation to make use of all available security controls.

36.5 ENCRYPTION. Many individuals and organizations focus on encrypting data while it is in motion, transiting networks and the Internet. The use of encryption for data while at rest, or stored, is equally important. As previously mentioned, stored data are compromised by security breaches more often than data in transit. With the use of a sufficiently strong algorithm and sufficiently sized key, data that are stored encrypted can be made unusable for those without the ability to decrypt the data. Even if an attack is successful at gaining full control of copy of the data for brute force attempts at decryption, a proper algorithm and key can prevent the data from being decrypted for a reasonably sufficient period—long enough that any personally identifiable information (PII) or other sensitive, confidential, or proprietary information would be of no value except to historians.

When a system is lost or stolen, an attacker has, essentially, an unlimited amount of time to gain access to the data. If the data are not encrypted, they can simply be read once the userID and password have been determined. If cracking the operating-system access controls is not successful, an unencrypted disk drive can simply be read on a different computer using the file system or, if that’s not possible, using

ENCRYPTION 36 · 11

forensic utilities. With encrypted file systems or volume encryption, only a portion of the data is stored encrypted. Data that are not in one of these encrypted locations are vulnerable—including swap spaces and temporary file locations used by the operating system.

See Chapter 7 in this *Handbook* for more details of cryptography; see Chapter 37 for details of the public key cryptosystem.

36.5.1 Recoverability. Key principles of information assurance (IA) include protecting the availability and utility of data. By its nature, encryption can take away these protections in exchange for protecting the authenticity, confidentiality, and integrity of the data while it mitigates the risk of losing possession of the data. When using data encryption it is important to consider the possible need to recover the data in the event that the user or primary key is unable to be located. There are several ways to accomplish this.

Key escrow is one method to facilitate the recovery of encrypted data. By storing the key with a trusted party, a lost key can be removed from escrow and used to decrypt the data. With public key encryption, additional decryption keys (ADKs) can be used with some encryption tools. Corporate ADKs are keys that can be used during the encryption process to automatically encrypt the data with a key that is tightly controlled and only used by designated individuals for the recovery of encrypted data.

36.5.2 File Encryption. The use of encryption on a file-by-file basis is a good method for securing data. With file encryption, a user can pick and choose which files to encrypt. Files containing sensitive data can be encrypted, while nonsensitive files are stored without encryption. This has the least impact on a system, but it puts most responsibility on a user who must determine which files should be encrypted.

Operating system files cannot be encrypted, so configuration files that may contain information on the organization’s systems are left as risk—as are application code files. These code files may be for the organization’s proprietary application that provides a significant competitive advantage. However, such files are rarely encrypted, as doing so can provide a hindrance for the user who must remember to decrypt each and every file when required.

36.5.3 Volume Encryption and Encrypting File Systems. Both volume encryption and encrypting file systems offer protection to data and can be easier for users than per-file encryption. The ease of use comes from the single operation required to access multiple files. Depending on its configuration, a location can remain decrypted and accessible for a few minutes or several hours.

A critical difference between file encryption and volume encryption is that decryption is carried out at the driver level as data are moved from disk into RAM. The entire volume is *not* decrypted. However, input/output (I/O) to and from the mounted volume can be significantly slower than from an unencrypted disk drive.⁸

In addition, dynamic decryption ensures that there is no extensive plaintext version of the original materials stored on the hard disk or in virtual memory for cryptanalysis in case of an aborted shutdown.

With partial encryption using file encryption or encrypted volumes, only a portion of the data stored on the disk and specifically written to the encrypted partition is protected. Usually, operating-system and application files are not encrypted, so a stolen or otherwise compromised disk is vulnerable to attackers who would gain access to any file not saved into the encrypted volume or file system. If a user forgets to encrypt

36 · 12 SECURING STORED DATA

a file or folder containing sensitive data, then those unprotected data are accessible to anyone who can access the operating system or read the disk drive using forensic or other disk utilities.

36.5.4 Full Disk Encryption. Author Ryan Groom lists cogent reasons for using full disk encryption on laptops.⁹ These can be restated for any system. The primary reasons to use full disk encryption are that it

- Protects data if a drive is lost or stolen,
- Is safer and more effective than volume encryption or encrypting file systems,
- Can be transparent to users, and
- Helps comply with legal and regulatory issues.

Full disk encryption secures the file system and operating-system files but leaves a small boot portion of the drive unencrypted. The unencrypted region allows the encryption software to load, to request the password, passphrase, or token needed to initiate dynamic decryption of the drive contents on demand, and to continue loading the operating system.

Depending on the solution chosen, users can see little difference in functionality or performance between a system using full disk encryption and a system that does not do so. The primary visibility to users is that on a system boot, the user has to identify and authenticate to enable decryption and continue to system boot. System performance is almost undetectably reduced, with only minor delays during the boot of the system while significant amounts of operating-system programs and data from the disk are decrypted, and again on shutdown as unencrypted data are cleaned up to prevent readability without authorization. Although I/O may be slower than on unencrypted systems, most file accesses do not require large volumes of I/O, so overall, these minimal effects on users are greatly outweighed by the protection afforded by full disk encryption.

As with volume encryption, full disk encryption involves dynamic decryption of ciphertext as it flows from disk to memory buffers. With modern data-transfer speeds and processor capabilities, any performance delay due to on-the-fly decryption once the operating system has been loaded is negligible in practice.

With full disk encryption, the entire contents of the drive are protected. Even with full physical access to the disk (e.g., by installing it into another computer under the attacker's control) or with a bit-for-bit copy of the encrypted disk, an attacker must break the encryption in order to gain any information—an almost impossible task with the key sizes currently in use, except perhaps by government cryptanalysis labs using massively parallel architectures for brute-force cracking. With strong encryption, management may be able to satisfy the concerns of clients if an organization has to disclose the loss of equipment and must provide an assurance that, even though the disk was lost, the client and organization's data cannot be accessed.

36.5.5 Vulnerability of Volume, File System, and Full Disk Encryption. As strong as the protection provided by volume, file system, and full disk encryption, there is one significant weakness. Once a user is authorized to access the data, and the operating system dynamically decrypts data as required, the system is vulnerable to attacks by any interloper who has physical access to the unlocked, unprotected session. For example, if a system contains sensitive data and is connected to

ENCRYPTION 36 · 13

a network, any attack over the network can potentially compromise the data when the authorized user's session allows access to the decrypted data.

This vulnerability must be stressed to users who may misunderstand the implications of encryption, especially those who insist on storing sensitive data on their laptops. While full disk encryption protects the data when the system is not booted, once the user decrypts the disk at boot the data are available not just to them but to anyone else who gains access to the system. Systems must use defense-in-depth strategies with personal firewalls to prevent unauthorized network access to the system; users must maintain physical possession of the system, especially once they have entered their decryption key. If a user with a Windows operating system on a laptop locks the screen and then walks away, the data are protected only by the strength of the system password.

Especially sensitive data should be encrypted at the file level. Alternatively, volume and file system encryption can be used with reasonable timeouts applied. Both of these options provide increased protection of data while a system is booted. Combined with full disk encryption, the risks to data are greatly reduced. Some application programs offer options for data encryption (e.g., MS-OUTLOOK includes options for encryption of the PST files containing all the user's data). And what could be worse for an attacker who breaks the full disk encryption only to find that the information is super-encrypted one or more levels, if file encryption, volume encryption, and full disk encryption are all in use?

An important point with the discussion of multilevel encryption is that it requires a strong commitment from the organization to support security at that level. Full disk encryption can be enabled at a global level through group policies or other technology-driven solutions. However, as mentioned above, a user must selectively choose to encrypt sensitive data, placing the responsibility and accountability on them. However, with a multi-layered encryption approach, there is still a layer of protection against a user forgetting or deliberately avoiding encrypting a sensitive piece of data.

Many data privacy laws have arisen over the last several years, including differing laws for the majority of U.S. states and the European Union. Many of these laws require notification of impacted users when data is lost. However, many of these laws include "safe harbor" from notification if the data was encrypted on the lost or stolen device or media. (The authors strongly recommend that you consult with your legal counsel for proper interpretation of the laws for your jurisdiction and/or the jurisdiction of any data subjects potentially impacted by a breach or data loss.)

36.5.6 Database Encryption. For many organizations, the databases that are the primary location for storing data are also an excellent target for attackers. By employing database encryption, the time an attacker needs to gain access can be greatly increased.

Databases can be protected by placing them on a system that can only be accessed via a secured connection and one that employs volume or full disk encryption. In addition, there are database encryption functions or external tools designed to protect the data held within these stores even without disk encryption. These tools employ encryption at the field (individual data element), row (a collection of fields that align in a row when seen in a tabular view), and full database encryption. One of the advantages of database-specific encryption is that it can support strict compliance with legal requirements for data protection. For example, PII such as names, Social Security numbers, medical information, and so on can be protected at the highest level of security available in case unauthorized personnel or intruders gain access to the database.

36 · 14 SECURING STORED DATA

James C. Foster, in an article for *SearchSecurity.com*, stated that organizations “often jump into database encryption as a quick fix for compliance without considering several key factors. The greatest among these factors is the speed or performance of the application, because poorly implemented database encryption could impact production applications.”¹⁰ Foster recommends four “simple guidelines that will help you secure your database without impeding the business you’re trying to protect,” which are expanded on below:

- Do not encrypt foreign or super keys. These keys are used for indexing, and encrypting them can negatively impact the utility of the database. Since these keys are not encrypted, the keys should never contain information that should be protected—such as using a customer’s Social Security number or credit card number as a key to link tables.
- As with any use of encryption, symmetric algorithms are faster than asymmetric keys. However, if all data are encrypted with a single key, this key must be well protected or else an attacker finding the key could, quite literally, have the key to the kingdom of an organization’s primary data store.
- “Full database encryption is rarely advised or a feasible option. Security best practices would teach you to encrypt everything with multiple keys and differing algorithms.” However, technical staff should evaluate the effects on performance.
- “Encrypt only sensitive data [columns]. This is typically all that is required or recommended by regulations and, after all, is what needs protection.”

36.5.6.1 Improving Vendor Provided Options. Database vendors such as Oracle Corporation and Microsoft have encryption options that are specifically designed to protect the information in their databases. These options have improved over time and are becoming more robust and mature.

Microsoft SQL Server 2005 offers enhancements to column encryption. Also introduced was “an integrated and hierarchical infrastructure for managing encryption keys.” The product documentation continues, “Built-in encryption functions and application programming interfaces (APIs) make it easier for an organization to create an encryption security framework.”¹¹

Oracle Database 10g Release 2 improves on existing encryption options within Oracle databases by introducing Transparent Data Encryption (TDE). When using TDE, a database administrator is able to specify that a column needs to be encrypted, and the database automatically encrypts data during insert operations and decrypts the data during selects. This can be achieved “without writing a single line of code.”¹² Arup Nanda provides a good overview of this feature in the September/October 2005 issue of *Oracle Magazine*.¹³

36.5.6.2 Implementation Considerations. As with any implementation of encryption, careful consideration must be given to determining the method and process for implementing the solution, as well as to the data that are being encrypted. Encrypting a key field needs to be avoided. If sensitive data are in the key field, there may be significant work required to create new key fields and recreate table linkages, or the organization may decide to accept the performance degradation that could occur with encrypting the key field. That is, assuming that encrypting the key field does not cause the database to become unusable.

DATA DISPOSAL 36 · 15

The costs associated with implementing an encrypted database solution must also be weighed against the business risk. For companies that have little data needing encryption, database encryption may be inappropriate. Legal and regulatory requirements also need to be considered as database encryption may be mandated in order to protect data and to avoid criminal or civil liabilities. The attendant loss of public and customer confidence, in the event of a data breach, is also a powerful incentive.

36.5.7 Smart Phone Encryption. Users may store confidential information on mobile phones and tablets; typical entries include contact entries with phone numbers and sometimes with ancillary sensitive data such as passwords, personally identifiable information (e.g., government-issued identification numbers), and call records. Such devices may also be used as if they were flash drives, with potentially gigabytes of sensitive data downloaded from other sources and carried in a pocket, briefcase, or handbag—and therefore easy to steal or to lose.

Another factor that can be significant for some users is that under U.S. law at the time of writing (in June 2013), a suspect who is questioned, interrogated, or arrested cannot normally be forced to divulge the decryption code.¹⁴

Phones using Android 2.3.4 or later usually come with integrated total encryption; the process typically takes about an hour, ideally starts with a fully charged battery and connection to a power supply, and must not be interrupted. Interruption of this encryption process can damage or delete the data stored on the phone and requires a factory reset that wipes all current data and personal settings from the device.¹⁵

Apple iOS and Microsoft Windows Phone 7 also include encryption functions with varying coverage. Third-party software is available for all the operating systems discussed above.¹⁶

In March 2013, researchers at the Friedrich-Alexander University discovered how to access data encrypted on a version of the Android operating system:

The team froze phones for an hour as a way to get around the encryption system that protects the data on a phone by scrambling it. . . . The attack allowed the researchers to get at contact lists, browsing histories, and photos. . . . [They] put Android phones in a freezer for an hour until the device had cooled to below –10 C. . . . [Q]uickly connecting and disconnecting the battery of a frozen phone forced the handset into a vulnerable mode. This loophole let them start it up with some custom-built software rather than its onboard Android operating system. The researchers dubbed their custom code *Frost—Forensic Recovery of Scrambled Telephones*.¹⁷

36.6 DATA DISPOSAL. A final consideration for securing stored data is the disposal of the media that contains the data.

For sanitizing electronic media, United States Department of Defense standard 5220.22-M can provide guidance.¹⁸ Essentially, the media should be sanitized so that the data originally written to it cannot be recovered. One method for achieving this goal can include the following steps:

1. Erase the data
2. Write random or meaningless data to the media
3. Erase the data
4. Repeat until the desired level of sanitization is met

For information stored on paper, the paper should, at a minimum, be shredded before being discarded. Use of a cross-cut shredder is necessary as it increases the difficulty

36.16 SECURING STORED DATA

of piecing the documents back together. In one notorious incident, confetti examined after it was thrown during a Thanksgiving Day parade in New York City in 2012 proved to be easily read horizontal shreds of confidential files from the Nassau County Police Department. “There were whole sentences, license plate numbers, and police reports.”¹⁹

There are additional steps that can be taken if needed, and facilities are available to accomplish these, including burning the shredded paper and mixing it with water to speed its deterioration. Over the past several years, vendors specializing in onsite paper destruction have become commonplace. These vendors bring trucks to a client’s site and collect paper that is then shredded onsite before being taken to a facility that recycles, burns, or otherwise disposes of the waste in a secure manner.

For more information about data disposal, see Chapter 57 in this *Handbook*.

36.7 CONCLUDING REMARKS. The security of stored data is of critical importance. More data breaches occur against data in unsecured storage locations than is compromised during transit. With proper data storage security there is less risk to data from all kinds of threats, internal and external. Using secure channels for writing data to disk and protecting the disks themselves is of increasing importance, as attackers usually have one of two goals—causing systems or data to be made unavailable, or compromising the confidentiality and integrity of data. By combining secure communication channels, data encryption for data at rest, and physical protection of data storage devices, the security of information can be better assured.

36.8 FURTHER READING

- Curtis, W. C. *Backup & Recovery*, 2007. Sebastopol, CA: O’Reilly, 2007.
- Curtis, W. C. *Using SANs and NAS*, 2002. Sebastopol, CA: O’Reilly, 2002.
- Dwivedi, H. *Securing Storage: A Practical Guide to SAN and NAS Security*. Chapter 2: “SANs: Fibre Channel Security.” Addison-Wesley, 2006. www.awprofessional.com/content/images/0321349954/samplechapter/Dwivedi_ch02.pdf
- Griffin, D. *The Four Pillars of Endpoint Security: Safeguarding Your Network in the Age of Cloud Computing and the Bring-Your-Own-Device Trend*. CreateSpace Independent Publishing Platform, 2013.
- EMC Education Services. *Information Storage and Management: Storing, Managing, and Protecting Digital Information in Classic, Virtualized, and Cloud Environments*, 2nd ed. Wiley, 2012.
- Hitachi Data Systems. *Storage Concepts: Storing And Managing Digital Data*, Vol. 1. HDS Academy, 2012.
- Loftus, A., and C. Kerner, “SAN Lessons Learned,” 2003. http://dms.ncsa.uiuc.edu/set/san/src/San_Lessons_Learned.pdf
- Loshin, P. *Simple Steps to Data Encryption: A Practical Guide to Secure Computing*. Syngress, 2013.
- Shepler, S., B. Callaghan, D. Robinson, R. Thurlow, C. Beame, M. Eisler, and D. Noveck. “RFC 3530 Network File System Version 4 Protocol.” April 2003. Proposed Standard Requests for Comment, Internet Engineering Task Force. <http://tools.ietf.org/html/rfc3530>
- Souppaya, M., and K. Scarfone. *Guidelines for Managing the Security of Mobile Devices in the Enterprise*. NIST Special Publication SP800-124 Revision 1. June 2013. <http://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-124r1.pdf>

NOTES 36 · 17

- Systems and Network Analysis Center. “Securing Fibre Channel Storage Area Networks.” National Security Agency. February 10, 2009. www.nsa.gov/ia/_files/factsheets/securing_fibre_brochure.pdf
- Troppens, U., R. Erkens, and W. Müller, *Storage Networks Explained: Basics and Application of Fibre Channel SAN, NAS iSCSI and InfiniBand*. Wiley, 2004.
- Vacca, John. “The Basics of SAN Security, Part I”. *Enterprise Storage Forum*. July 23, 2002. www.enterprisestorageforum.com/sans/features/article.php/11188-1431341_2
- Verhelst, Wouter. “Securing NFS.” Wouter Verhelst’s Blog, *Free Software Magazine*. November 26, 2006. www.freesoftwaremagazine.com/blogs/securing_nfs

36.9 NOTES

1. Ponemon Institute, “New Study Reveals Up To 12,000 Laptop Computers Lost Weekly and Up To 600,000 Lost Annually in U.S. Airports,” www.ponemon.org/news-2/8
2. M. J. Schwartz, “Lost Laptops Cost \$1.8 Billion Per Year: Only One-Third of Missing Laptops Have Full-Disk Encryption for Preventing Data Breaches, Finds Ponemon Study of European Firms,” *InformationWeek Security*, April 21, 2011, www.informationweek.com/security/mobile/lost-laptops-cost-18-billion-per-year/229402043
3. D. Klugh, “12,000 Laptops Lost in Airports Every Week,” *WMBF News*, October 1, 2012, www.wmbfnews.com/story/19604344/12000-laptops-lost-in-airports-every-week
4. A. Delli Santi, “NJ Audit: Confidential Data on Junked Computers.” ABC Action News, WPVI-TV, March 9, 2011, <http://abclocal.go.com/wpvi/story?section=news/local&id=8002765>
5. M. Kincora, “Strategic Storage: Storage Security—Change Old Habits and Stop Data Theft,” *SearchStorage.com*, November 2005, <http://searchdatabackup.techtarget.com/news/1300734/Strategic-Storage-Storage-security-Change-old-habits-and-stop-data-theft>
6. W. C. Preston, H. Dwivedi, M. E. Kabay, and S. Gordon, *Storage Security Handbook*, Neoscale Systems, Inc., 2002, www.neoscale.com/English/Downloads/Storage_Security_Handbook/ (URL inactive)
7. Preston et al., *Storage Security Handbook*.
8. For example, observations in June 2013 showed that mounting an entire 7 GB partition encrypted by Symantec PGP Desktop v10.20.0 took ~7 seconds. In comparison, the same 7 GB of data encrypted using WinZip and 256-bit AES encryption took about an hour to decrypt *in toto* on a 7200 rpm two-disk performance RAID 0 set running under 64-bit Windows 7 Professional SP 1 on a 3.2 GHz AMD Phenom™ II X4 955 quad-core processor with 12 GB of RAM. Copying a 1.23 GB file took 85 seconds (~15 MB/sec) on the PGP encrypted volume (located on the RAID 0 hard drive) but only 20 seconds (~63 MB/sec) on the RAID 0 itself.
9. R. Groom, “8 Reasons for Full-Disk Encryption,” *Business Security*, About.com <http://bizsecurity.about.com/od/windowsdesktopsecurity/a/top8fulldisk.htm> (URL inactive)
10. J. C. Foster, “Look Before Leaping into Database Encryption,” *Compliance Counselor*, *SearchSecurity.com*, September 29, 2006, http://searchsecurity.techtarget.com/tip/0,289483,sid14_gci1219561,00.html

36 · 18 SECURING STORED DATA

11. Anonymous, "Improving Data Security by Using SQL Server 2005," Microsoft® TechNet. <http://technet.microsoft.com/en-us/library/bb735261.aspx>
12. A. Nanda, "Transparent Data Encryption," Technology:Security, ORACLE® Technology Network, www.oracle.com/technology/oramag/oracle/05-sep/o55_security.html
13. Nanda, "Transparent Data Encryption."
14. R. Radia, "Why You Should Always Encrypt Your Smartphone," *ArsTechnica*, January 16, 2011, <http://arstechnica.com/gadgets/2011/01/why-you-should-always-encrypt-your-smartphone/>
15. How-To Geek, "How to Encrypt Your Android Phone and Why You Might Want To," 2013, www.howtogeek.com/141953/how-to-encrypt-your-android-phone-and-why-you-might-want-to/
16. UNC-Chapel Hill, "Encrypting Cell Phones," University of North Carolina at Chapel Hill | Help & Support, March 19, 2013, <http://help.unc.edu/help/encrypting-cell-phones>
17. BBC, "Frozen Android Phones Give Up Data Secrets: Freezing an Android Phone Can Help Reveal Its Confidential Contents, German Security Researchers Have Found," *BBC News | Technology*, March 7, 2013, www.bbc.co.uk/news/technology-21697704
18. Anonymous, "DoD 5220.22-M National Industrial Security Program Operating Manual (NISPOM)," 1997, www.usaid.gov/policy/ads/500/d522022m.pdf
19. C. Boyette, "Sensitive Documents Found in Macy's Thanksgiving Day Parade Confetti," *CNN US*, November 26, 2012, www.cnn.com/2012/11/26/us/new-york-confidential-confetti

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER **37**

PKI AND CERTIFICATE AUTHORITIES

**Santosh Chokhani, Padgett Peterson,
and Steven Lovaas**

37.1 INTRODUCTION	37·2	37.6.7	Public Key Infrastructure Interoperability	37·14
37.1.1 Symmetric Key Cryptography	37·2			
37.1.2 Public Key Cryptosystem	37·2	37.7 FORMS OF REVOCATION	37·18	
37.1.3 Advantages of Public Key Cryptosystem over Secret Key Cryptosystem	37·3	37.7.1 Types of Revocation-Notification Mechanisms	37·18	
37.1.4 Combination of the Two	37·4	37.7.2 Certificate Revocation Lists and Their Variants	37·19	
		37.7.3 Server-Based Revocation Protocols	37·20	
37.2 NEED FOR PUBLIC KEY INFRASTRUCTURE	37·4	37.7.4 Summary of Recommendations for Revocation Notification	37·21	
37.3 PUBLIC KEY CERTIFICATE	37·5			
37.4 ENTERPRISE PUBLIC KEY INFRASTRUCTURE	37·7	37.8 REKEY	37·22	
		37.9 KEY RECOVERY	37·23	
37.5 CERTIFICATE POLICY	37·8	37.10 PRIVILEGE MANAGEMENT	37·25	
37.6 GLOBAL PUBLIC KEY INFRASTRUCTURE	37·9	37.11 TRUSTED ARCHIVAL SERVICES AND TRUSTED TIME STAMPS	37·25	
37.6.1 Levels of Trust	37·10			
37.6.2 Proofing	37·10			
37.6.3 Trusted Paths	37·10	37.12 COST OF PUBLIC KEY INFRASTRUCTURE	37·27	
37.6.4 Trust Models	37·11			
37.6.5 Choosing a Public Key Infrastructure Architecture	37·13	37.13 FURTHER READING	37·27	
37.6.6 Cross-Certification	37·14	37.14 NOTES	37·28	

37 · 2 PKI AND CERTIFICATE AUTHORITIES

37.1 INTRODUCTION. In the 1990s, the use of encryption across the Internet consisted mainly of individuals with Pretty Good Privacy (PGP) and later the Gnu Privacy Guard (GnuPG, originally abbreviated GPG) exchanging secure email and each maintaining a private *web of trust*, today's use of encryption encompasses a much wider range of elements, including proofing, issuance, revocation, identification, federation, bridging, encryption, digital signing, and a myriad of ancillary processes. In fact, what the user experiences is just the tip of the required support structure. Proper management of information involved in an encryption infrastructure (or cryptosystem) is about trust, what is required to establish that trust, and how much to grant.

Back then, the prime use of encryption was outside the network perimeter, where unprotected data sent and received by the organization over public networks were perceived as most vulnerable to disclosure, modification, insertion, deletion, and replay attacks. To protect the data being transported over untrusted networks, the only practical and cost-effective technology is cryptography. Cryptography is at the heart of both virtual private networks (VPNs) and public key infrastructures. For further information on encryption, see Chapter 7. For a review of cryptography, see Chapter 7 in this *Handbook* and *Applied Cryptography* by Bruce Schneier.¹

37.1.1 Symmetric Key Cryptography. Symmetrical cryptography (also known as *single key* or *secret key*) uses the same key to encrypt cleartext into ciphertext and to decrypt ciphertext back into cleartext. This process is illustrated in Exhibit 37.1.

Although symmetric cryptography (e.g., Advanced Encryption Standard [AES], Data Encryption Standard [DES], Rivest Cipher 4 [RC4]) has specific advantages in that it has a dense keyspace (any integer) and is computationally very fast, it has a fundamental weakness: Both the originator and each recipient must have the key. Key sharing has two management issues:

1. Key exchange must be performed: Somehow each user must have the same key. This exchange is generally accomplished out of band.
2. There is weak confidentiality, control, or authentication: Anyone with the symmetric key could have encrypted the cleartext or decrypt the ciphertext.

37.1.2 Public Key Cryptosystem. In contrast to secret key cryptosystems, public key cryptosystems (PKCs) use pairs of related keys that are generated together. It is accepted that it is infeasible to generate one key from the other. The ciphertext produced by one key can be decrypted only with the other member of the same key

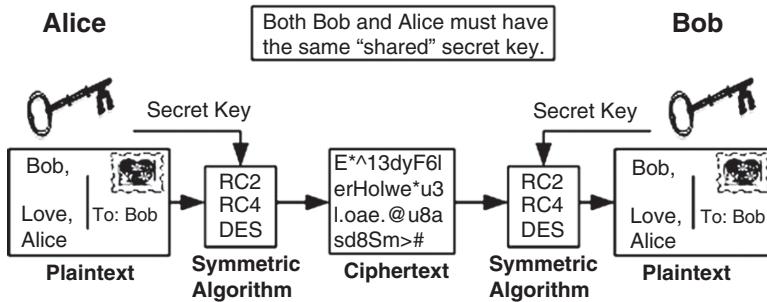
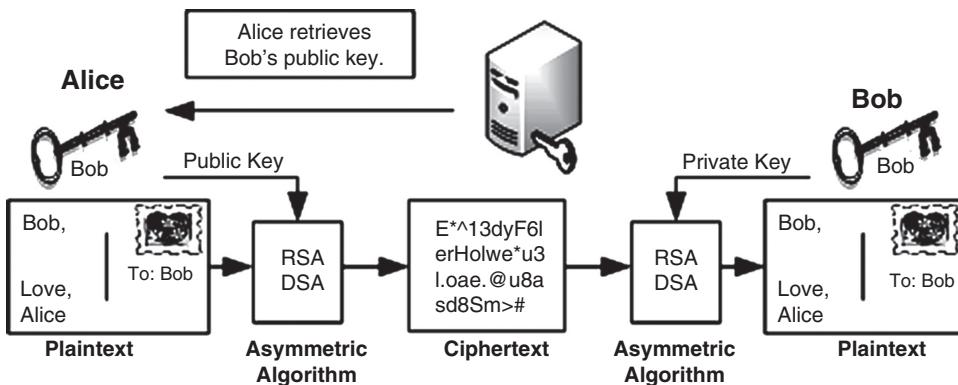


EXHIBIT 37.1 Symmetric (Secret) Key Encryption

INTRODUCTION 37 · 3**EXHIBIT 37.2** Asymmetric or Public Key Encryption

pair. One of these keys is kept secret (the *private key*) and the other is published for all to use (the *public key*). It does not matter which is designated as which, but once designated, the key cannot be changed. In practice, encryption with one key is often faster than with the other. This one is generally selected as the public key.

In the simplest form, to conceal a message in transit so that only the desired recipient may read it, the cleartext is encrypted using the recipient's public key, as shown in Exhibit 37.2. Only the recipient's secret key can decrypt the transmitted ciphertext. Similarly, to verify message integrity and authenticity, it is possible to encrypt information with a sender's private key. This allows anyone with access to that sender's public key to validate the message by decrypting the ciphertext successfully.

37.1.3 Advantages of Public Key Cryptosystem over Secret Key Cryptosystem.

For securing data transmissions, public key cryptosystems are preferred over secret key cryptosystems for these reasons:

- Public key cryptosystems require fewer keys to manage: Each party (n) has a key pair, so the total number of keys is $2n$, instead of being proportional to n^2 as for secret key cryptosystems.
- Because private keys need not be distributed or otherwise managed, public key cryptosystems require only demonstrated integrity and authenticity of the public keys themselves. Users (the *relying parties*) must have assurance that their public keys truly belong to the publishers. This requires signing by a trusted third party or by a mutually trusted source.
- Because no secret keys are transmitted over any network, PKCs are not susceptible to compromise even when public keys have to be changed. PKCs can be used to encrypt temporary keys (*session keys*) that can be used one time for secret key cryptography to obviate the heavier computational load of the PKC.
- To encrypt a message so that multiple PKC users can receive and decipher the ciphertext securely, PKC software can create a session key. This secret key is then encrypted with each recipient's public key separately and sent with the ciphertext to all recipients without compromising confidentiality. Each recipient can then extract the key used to encrypt the file.

37 · 4 PKI AND CERTIFICATE AUTHORITIES

EXHIBIT 37.3 Symmetric versus Asymmetric Encryption

Type	Keyspace	Speed
Asymmetric	Sparse	Slow
Symmetric	Dense	Fast

PKC-based digital signatures can also provide the basis for nonrepudiation in the event of a dispute. Only the possessor of a private key could have sent a message decrypted by its public key. In contrast, because of the use of shared secrets, symmetric secret key cryptosystems alone cannot reasonably support nonrepudiation.

37.1.4 Combination of the Two. At this point, it may be easiest to consider the elements of the different types of cryptography, as shown in Exhibit 37.3.

The difference is similar to that between electronic data interchange and e-commerce. Symmetric encryption is good for a very few exchanges with people you trust; asymmetric encryption is good for many, many small exchanges with people or devices you may have never met.

Today, the most common use is to combine both types: For a document to be emailed securely, a symmetric algorithm is selected and a random key generated. The document is encrypted using the symmetric algorithm and key. The symmetric key is then encrypted with the asymmetric public key of each recipient and is then added as a header to the document.

For digital signing, a similar mechanism is used except that a hashing algorithm (e.g., Secure Hash Algorithm [SHA]) is used to create a hash value of the document, and the hash is then encrypted (signed) with the private key of the originator and accompanies the document. The originator's public key can be used to decrypt (verify) the hash, and the hash is used to verify the integrity of the document.

37.2 NEED FOR PUBLIC KEY INFRASTRUCTURE. The PKC depends on the integrity of each public key and of that public key's binding to a specific entity, such as a person, an institution, or a network component. Without mechanisms for ensuring integrity and authenticity, a relying party is vulnerable to masquerading attacks through public key substitution.

To illustrate, suppose that ABC Company wants to send a confidential message to XYZ Corp. that no one else can read. ABC could use XYZ's public key to encrypt the message, although, for the sake of efficiency, ABC probably would use a symmetric algorithm to encrypt the message and the public key algorithm to encrypt the symmetric key. However, if ABC can be tricked into using an attacker's public key as if it were XYZ's public key, then the attacker would be able to decrypt the message. This technique is known as public key *spoofing*.

Such public key spoofing by the attacker would, however, make it impossible for XYZ to read the message from ABC, as was originally intended. Therefore, such an attack probably would continue with the attacker reencrypting ABC's message, using XYZ's real public key, and sending it on to XYZ. Such interception and reencryption is an example of a *man-in-the-middle* attack. However, if the sender also signed the original message, this is something that the attacker could not duplicate, so the message could not be changed but only intercepted.

PUBLIC KEY CERTIFICATE 37 · 5

Another example of breaching the connection between a public key and its owner involves digital signature verification. Suppose ABC wants to verify XYZ's signature. If the attacker could trick ABC into using the attacker's public key as if it were XYZ's public key, then the attacker would be able to sign messages masquerading as XYZ using the attacker's private key. ABC would unknowingly use the replaced public key and be spoofed into thinking that the message actually was signed by XYZ.

This same problem exists with the core of Internet commerce today, Secure Sockets Layer v2 (SSLv2). SSLv3 and Transport Layer Security (TLS) can prevent this, but then every participant must have a certificate trusted by the other party. This is where the concept of a *trust chain* comes in. If a trusted third party has signed XYZ's public key, then the attacker would also have to have a key signed by the same authority as well. Otherwise, a warning would appear that the key was not recognized.

A serious issue is the sheer number of certificate authorities (over 100) built into most browsers, from a variety of countries. The browser trusts all certificates issued by any of these by default, and not all are necessarily worthy of trust.

In summary, both for digital signature and encryption services, the relying party must use the public key of the correct party in order to maintain security. There are various manual, electronic, and hybrid mechanisms for the distribution of public keys in a trusted manner, so that the relying party can be sure to have the correct public keys of the subscribers. These mechanisms for distribution and binding of public keys are known as a Public Key Infrastructure [PKI]).

The beauty of a signed public key is that the trust is inherent with the signature provided the signer is trusted. Unlike a *web of trust*, such as that used by the original PGP which requires all keys be accepted individually by each user, with a *chain of trust* a single signer (the root certificate authority) can provide trust to millions of certificates. Trust need not be absolute; it may be contextual. For example, a certificate signed by an employer can be trusted for elements relating to employment but not for credit cards. See the discussion of trust levels in Section 37.6.1.

37.3 PUBLIC KEY CERTIFICATE. The technique that is most scalable uses a public key certificate issued by a trusted party called the certification authority (CA). A CA issues public key certificates to the various subscribers by putting together subscriber information and signing the information using the CA's private key. The generally accepted standard for public key certificates is the X.509 version 3,² as defined in RFC 5280 and subsequent text.

Note: Hierarchical structures may use one or more *root* CAs and each Root CA may have multiple *intermediate* CAs which issue the end user and server certificates

X.509 certificates are expressed in Abstract Syntax Notation 1 (ASN.1), which is a complex binary notation. In order to be passed through email, certificates are usually MIME (aka Base 64) encoded, expressing the binary syntax in ASCII characters.³

The advantage to using a CA and a chain of trust is that by trusting the root, the user automatically trusts all keys that it has issued whether or not the user has ever seen them.

Each CA's certificate may contain this key information:

- Version number of certificate standard
- Certificate serial number (unique for every certificate issued by the CA)
- Algorithm and associated parameters used by CA to sign the certificate
- CA name

37 · 6 PKI AND CERTIFICATE AUTHORITIES

- Validity period for the certificate
- Subscriber name
- Subscriber public key, public key algorithm, and associated parameters
- CA unique identifier (optional)
- Subscriber unique identifier (optional)
- Extensions (optional)
- CA's digital signature

The relying parties require the CA's public key so that they can verify the digital signatures on the certificates issued by the CA. The relying party must trust the CA's public key, most likely obtained during the registration process. Once the signatures are verified, relying parties can use the subscriber name and subscriber public key in the certificate with as much confidence in the accuracy of the information as they have in the trustability of the CA.

In some situations, a CA may need to revoke the binding between a subscriber and that subscriber's public key. For example, the subscriber private key may be compromised (i.e., there may be reason to believe that the secret key has fallen into the hands of someone else). Since a public key certificate is an electronic object and can reside in several places at the same time, it is neither practical nor possible to recall, delete, or erase all the copies of the subscriber certificate in a distributed environment. Thus, to invalidate a public key certificate by severing the binding between the subscriber and the subscriber public key, the CA creates a list of invalid certificates. This list is called a *certificate revocation list* (CRL). The relying parties must check that a certificate is not on the CRL prior to using the public key in the certificate. If the certificate is on the CRL, the relying party must not use it. The CA signs the CRL to allow the relying parties to verify the CRL's integrity and authenticity. The key information in the X.509 version 2 CRL is:

- Version number of CRL standard
- Algorithm and associated parameters used by CA to sign the certificate
- CA name
- This CRL issuance time
- Next CRL issuance time (optional)
- List of revoked certificates (listing these items for each certificate):
 - Certificate serial number
 - Time CA was notified of revocation
 - Extensions related to the revoked certificate (optional)
- Extensions related to CRL (optional)
- A's digital signature

It is important that a certificate contain only those elements specifically required for operation and that these be used properly otherwise confusion and improper/unexpected use may result.

Probably the most misused element is the *Key Usage* extension. (It is used so badly that often a second extension, the *Extended Key Usage*, is used to clarify the first.) The most common error is to have the *Non-Repudiation Bit* set on a key intended for

ENTERPRISE PUBLIC KEY INFRASTRUCTURE 37 · 7

identity only. Nonrepudiation should be asserted only on a key intended as a legal signature, not for identity alone.

A dangerous extension (and one that really does not belong in a user certificate) is the *SMIME Capabilities* extension. This only has meaning in the context of specific email clients but can cripple the use of encryption. Although it may have had some use in early mail applications, today it only establishes a maximum allowed symmetric encryption strength, which may be less than desired. If this extension is not present, RFC 5280 requires a default to triple-DES.

In the case of certificates, less is more. The minimum number of fields and extensions that is required for the assertions the certificate is intended to make is best. PKI systems in the past have failed when too many, and sometimes mission-conflicting extensions, were added.

Particular care needs to be taken that the CA server vendor chosen does not add anything not specified by the certificate template in the Certificate Policy.

37.4 ENTERPRISE PUBLIC KEY INFRASTRUCTURE. The *use* of a certificate is relatively simple, but establishing trust that the certificate is valid and appropriate for use requires a complex set of back-office elements, illustrated in Exhibit 37.4. Each group of users covered by a CA is called a *domain*. Subscribers in a domain receive public key certificates from the appropriate CA. The CA is responsible for generation of subscriber certificates and for CRL generation. The CA posts these signed objects to the repository where the relying parties can obtain them. The CA also archives the certificates and CRLs in case they are required in the future to resolve disputes among the subscribers and the relying parties.

The registration authority (RA) is the trusted representative of the CA and is responsible for authenticating the subscriber identity. The RA typically performs these functions:

- Authenticates (proves) the subscriber's claimed identity. For example, the RA could require the subscriber to provide a valid photo ID, such as a driver's license or a passport for minimum assurance. Both I-9 authentication, as required by the

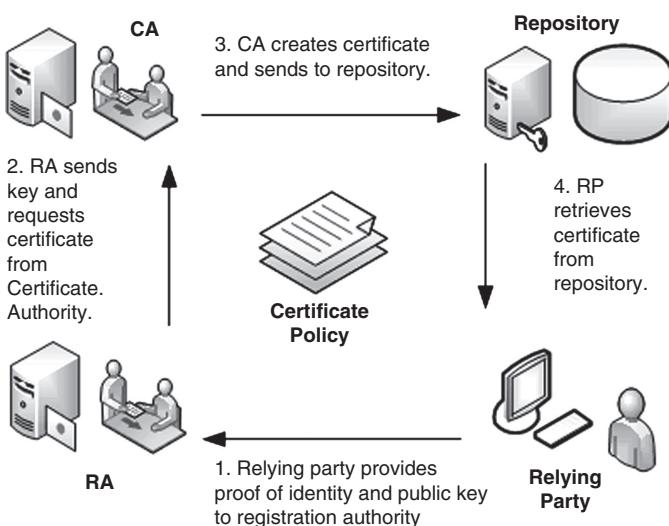


EXHIBIT 37.4 Certificate Issuance Cycle

37 · 8 PKI AND CERTIFICATE AUTHORITIES

Immigration Reform and Control Act (IRCA), and a Local Agency Check or a National Agency Check (LAC/NAC) plus “Need to Know” may be required for a higher level.⁴

- Obtains the subscriber public key from the subscriber.
- Provides the CA public key to the subscriber. A trust anchor is a CA’s public key that the relying party trusts. This trust generally is established by obtaining the public key from a trusted source using trusted means, such as physical hand-off or via Secure Sockets Layer (SSL) from a trusted or known Website. The CA public key becomes a subscriber trust anchor.
- Sends the certificate creation request to the CA. Typically, the RA creates an electronic mail message containing the subscriber name and the subscriber public key, digitally signs the message, and sends the message to the CA. Other transport means, such as manual or on the Web, also are appropriate as long as there is assurance that the subscriber identity and the public key are not changed. X.509 standard does not specify a protocol for certificate generation requests. The Public Key Infrastructure for X.509 Certificate (PKIX) working group of the Internet Engineering Task Force (IETF) has developed Internet standards in this area.⁵

37.5 CERTIFICATE POLICY. To ensure the security of the PKI, the PKI components need to operate with a high degree of security. To ensure this:

- Private keys must be kept confidential.
- Private keys must be used only by the owners of the keys.
- Trust anchors’ public key integrity must be ensured.
- Initial authentication of the subscriber (private key holder and the subject of the public key certificate) must be strong so that identity theft does not occur at the point of certificate creation.
- CA and RA computer systems and applications must be protected from tampering.
- Requirements for level of trust must be clearly defined

The Certificate Policy (CP) must specifically enumerate the certificate contents, both fields and extensions. Anything absent from the CP should not be found in the certificate.

In addition to the security requirements and in order to facilitate electronic commerce, the PKI must address obligations of all parties and their liabilities in case of dispute. These issues of security, liability, level of trust, and obligations are articulated in a CP.

According to the X.509 standard, a CP is “a named set of rules that indicates the applicability of a certificate to a particular community and/or class of application with common security requirements.”⁶ A certificate user may use a CP to decide whether a certificate, and the binding implied between the certificate and its owner, is sufficiently trustworthy for a particular application. The CP addresses security and obligations of all PKI components, not just the CA; this includes the CA, RA, repository, subscriber, and relying party.

A more detailed description of the practices followed by a CA in issuing and managing certificates is contained in a certification practice statement (CPS) published by or referenced by the CA. According to the American Bar Association’s *Digital*

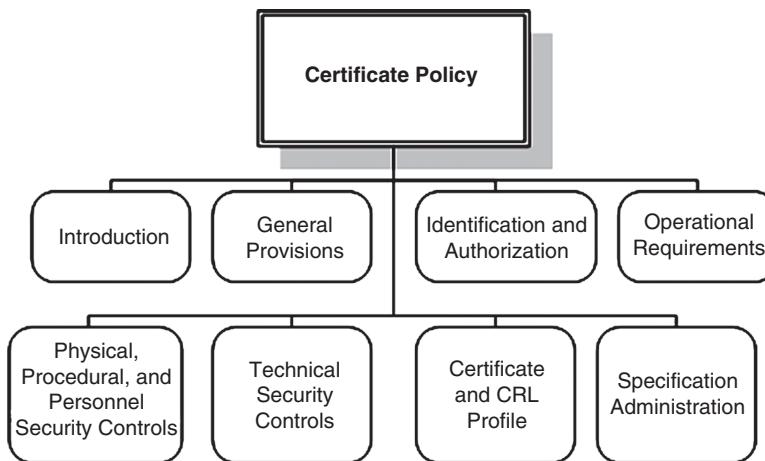
GLOBAL PUBLIC KEY INFRASTRUCTURE 37 · 9

EXHIBIT 37.5 Elements of a Comprehensive Certification Practice Statement

Signature Guidelines (hereinafter referred to as ABA *Guidelines*), “a CPS is a statement of the practices which a certification authority employs in issuing certificates.”⁷

Although a CP and a CPS both address the same topics, the CP defines the security requirements and obligations for an enterprise PKI, and the CPS describes how these requirements are satisfied by the enterprise PKI.

The CP and CPS also are used differently. The CP forms the basis for cross-certification across enterprise boundaries to facilitate secure, inter-enterprise electronic commerce. An object identifier (OID) pointing to the CP (which may be a Universal Resource Locator [URL]) is used to create a certificate that can be put into the “Certificate Policies” extension of X.509 certificates. The OID thus enables relying parties to learn the care taken during the generation of certificates, recommended usage, and obligations of the various parties.

Since certificates can be created with different levels of trust and can be either based in software or hardware (more secure), often the “Certificate Policies” extension is the only indicator of the level of trust that should be placed in a certificate.

The CPS enables PKI personnel to use and administer the PKI components. The CPS also forms the basis for compliance audits in order to ensure that the PKI components are operating in accordance with the stipulations of the CPS. Exhibit 37.5 illustrates the components of a comprehensive CPS. Components are divided further into subcomponents, which in turn are divided into elements. Components may be appropriate for various PKI entities but may be applied in the same way or differently. For example, technical security controls may apply to CA, RA, subscribers, and relying parties. These controls may be different for each of these entities, being most stringent for the CA, then the RA, and then the subscribers and relying parties.⁸

A sample policy may be found at www.verisign.com/repository/vtnCp.html

37.6 GLOBAL PUBLIC KEY INFRASTRUCTURE. The principles of an enterprise PKI with a single CA can be extended to support global, secure, electronic commerce by relying on multiple CAs and/or CAs to certify other CAs and each other. How the CAs cross-certify each other is also called *trust model*, *trust graph*, or *PKI architecture*. For one person to communicate securely with another, there must be a

37 · 10 PKI AND CERTIFICATE AUTHORITIES

EXHIBIT 37.6 Trust Level Determination

Category	Required Trust Level			
	1	2	3	4
Inconvenience or distress	Low	Med	High	High
Financial loss	Low	Med	Med	High
Harm to agency programs or public interests	N/A	Low	Med	High
Personal safety	N/A	N/A	Low	Med/high
Civil or criminal violations	N/A	Low	Med	High
Information classification				
Confidentiality	Low	Med	High	High
Integrity	Low	Med	High	High

Note: the trust **level** is sometimes confused with trust **factors** (something you have, something you know, something you know). A Smart Card may be a Level 3 (Medium, hardware) but part of a two factor (username, smartcard) authentication.

trust path from the trust anchor(s) of the relying party to the subscriber whose signature needs to be verified or to whom an encrypted message is to be sent.

37.6.1 Levels of Trust. As referenced in Office of Manpower and Budget Memorandum OMB M04-04, Section 2.1, there are four basic levels of trust⁹:

- Level 1.** Little or no confidence in the asserted identity's validity
- Level 2.** Some confidence in the asserted identity's validity
- Level 3.** High confidence in the asserted identity's validity
- Level 4.** Very high confidence in the asserted identity's validity

Each level requires a different initial authentication of identity, and higher levels (or those used for classified information) require background investigations. Level 3 is also known as Medium Assurance and is divided into two forms: (1) software (certificates and keys can be exported and moved between devices) and (2) hardware (certificates may be exported but keys and cryptographic functions are performed on a specific hardware device [e.g., a smart card, USB token, or PC-Card device]). Most commercial PKI uses the equivalent of the Medium level of assurance.

Exhibit 37.6 shows one way to determine the required trust level as a function of threat of misuse.¹⁰

37.6.2 Proofing. Each assurance level has a proofing (also known as *vetting*) requirement that increases with the trust level involved, as shown in Exhibit 37.7. Essentially none is required for level 1 (Basic) while extensive in-person requirements exist for level 4 (High).

37.6.3 Trusted Paths. The trust model can be viewed as a chain, with its tail a certificate-issuing CA and its head the subscriber (i.e., the subject of the certificate). The subscriber can be another CA or an end entity. To ascertain the trustworthiness of a certificate, it is necessary to start with the relying party trust anchor and to follow in the direction of the chain until the subscriber (of interest to the relying party) is reached.

GLOBAL PUBLIC KEY INFRASTRUCTURE 37 · 11**EXHIBIT 37.7** Trust Levels and Proofing

Level	Title	Proofing	Authentication
1	Default	Anonymous allowed.	None
2	Basic	Simple assertion—may be online.	Password
3	Medium (software)	I-9 employment eligibility verification and authorization. Must be in person.	Software certificate
3	Medium (hardware)	I-9 employment eligibility verification and authorization. Must be in person. Biometrics may be captured.	Hardware certificate
4	High	National agency check or local agency check, background investigation, and authorization required. Final proofing must be in person.	Hardware certificate

Global secure communications require that there be a trust path from every subscriber to every other subscriber.

The relying party can start with its trust anchor and verify the certificates issued by the trust anchor. Once that happens, the public keys can be trusted and used to verify the certificates issued by these CAs. This can be done recursively by the relying party until the public key certificate of the subscriber of interest is verified. Then the subscriber public key can be used to verify digital signatures and to perform encryption. Exhibit 37.8 illustrates this pathway. The arrows represent certificates. This is called a *certification path*.

37.6.4 Trust Models. Examples of trust models within PKI that relate to the trust in a certificate and are different from proofing (trust of an identity) include:

- Strict hierarchy
- Hierarchy
- Bridge

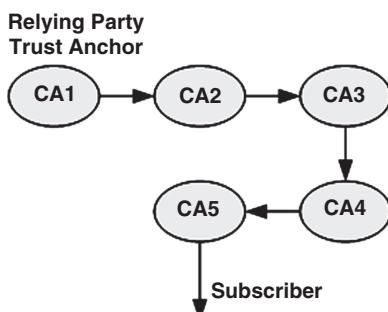


EXHIBIT 37.8 Trust Path through Multiple Trusting Certificate Authorities

37 · 12 PKI AND CERTIFICATE AUTHORITIES

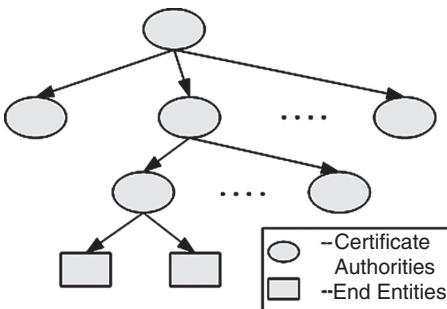


EXHIBIT 37.9 Strict Hierarchical Trust Chain

- Multiple trust anchors
- Mesh (aka anarchy or web)
- Combination

37.6.4.1 Strict Hierarchy. Exhibit 37.9 illustrates a strict hierarchy. It is a tree structure with a single root. In a strict hierarchy, for two parties to communicate with each other securely, they require the public key of their common ancestor as the trust anchor. Verifiable certificate chains require that the parties have a common ancestor. For all parties to communicate securely with each other, they require the single root as the trust anchor, since it is the only common trust anchor.

37.6.4.2 Hierarchy. In a (nonstrict) hierarchy, the subordinate CAs certify their parents. Since the directed graph is bidirectional, any CA can be the trust anchor for the relying parties. But from practical, operational, and performance (i.e., certificate path length) viewpoints, the local CA should be the trust anchor. The local CA is the CA that issued a certificate to the relying party.

37.6.4.3 Bridge. Another trust model is the bridge. Under this model, one CA cross-certifies with each CA from various domains. The domains can be organizations or vertical segments, such as banking or healthcare. Exhibit 37.10 illustrates the bridge CA. The bridge CA is not the trust anchor for any relying party. A CA in the domain of the relying party is the trust anchor.

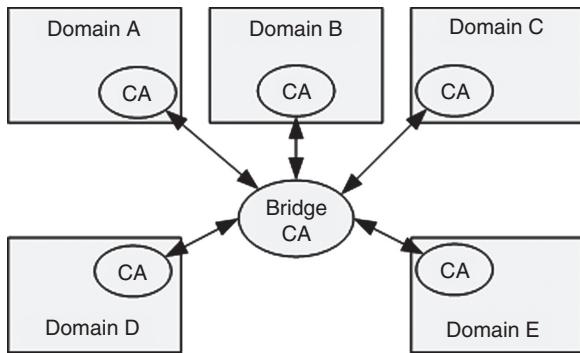


EXHIBIT 37.10 Bridge CA

GLOBAL PUBLIC KEY INFRASTRUCTURE 37 · 13

Within a domain, there are no constraints on the trust model. The domain PKI itself could be organized as any of the trusted models, including bridge, leading to possible layers of bridge CAs.

37.6.4.4 Multiple Trust Anchors. Another alternative is for the relying party to obtain the public keys of the various CAs in a trusted manner and then use those public keys as trust anchors. This approach is attractive when the CAs cannot, or are not willing to, cross-certify, and the relying party needs to communicate securely with the subscribers in the domains of the originating CAs. This approach is called *multiple trust anchors*. Each trust anchor representing a domain could be a single CA or a PKI with a collection of CAs in a trust model.

37.6.4.5 Mesh. The final example of trust model is a mesh (aka *web* or *anarchy*). The term “mesh” describes any depiction representing trust among CAs or certificates without any particular rules or patterns. This model sometimes is known as a *web of trust* and is particularly associated with the original design of Pretty Good Privacy, one of the first popular systems implementing public key certificates. Within this structure, each recipient must explicitly trust each other participant. The major problem is that it does not scale very well beyond a few hundred users.

37.6.5 Choosing a Public Key Infrastructure Architecture. Whether a domain (enterprise) chooses a single CA or multiple CAs for its intradomain operation should be determined from a variety of factors, including:

- Management culture
- Organization politics
- Certification path size
- Subscriber population size
- Subscriber population distribution
- Revocation information

In many situations, the politics or management structure may dictate that there be multiple CAs within the domain. In other words, organizations at business-unit level, regional office level, corporate level, or national level may want to create a CA in order to provide them with a certain degree of control, independence, autonomy, and prestige. How these CAs are organized (bilateral cross-certification, hierarchy, etc.) also will depend on the management and political landscape of the domain. The trust model should be such that it keeps the certification path size manageable; otherwise, end users will see unacceptable performance degradation in obtaining certificates and CRLs and in verifying digital signatures on the certificates and the CRLs.

Similarly, large subscriber populations may require more than a single CA in order to ensure that the CA can manage the subscribers and to keep the CRL size small. If CA products that issue partitioned CRLs are selected, the CRL sizes can be kept manageable even for a very large subscriber population. For further discussion of the CRL issue, see Section 37.7 on revocation alternatives.

When considering interdomain cross-certification, similar issues should be considered.

37 · 14 PKI AND CERTIFICATE AUTHORITIES

37.6.6 Cross-Certification. In the simplest form, cross-certification consists of two CAs that certify each other by issuing each other a certificate. The certificates can be stored in specific attributes of the directory entry in a certificate; examples include the cross-certificate attribute pair or the CA certificate.

There are two practical problems with cross-certification. One deals with the commercial products. If the two domains use different products, their CAs may not be able to exchange information to cross-certify, and their directories may not be able to chain to permit the relying parties to retrieve certificates.

The other problem is operational. Before certifying another CA, the certificate-issuing CA needs to make sure that the subject CA is operating in accordance with the acceptable controls, articulated in a CP. The issuing CA asserts the appropriate CP in the “certificate policies” extension of the X.509 version 3 certificate of the subject CA.

In practice, the two CAs cross-certify each other after reviewing each other’s CP and after ensuring that the CPs can be claimed to be equivalent. This does not mean that all the security controls and obligations are identical, but they need to offer roughly similar amounts of trust and of obligations and similar liability and financial relief.

When two CAs cross-certify each other, the trust generally is for a limited set of policies through assertions in “certificate policies” extensions, and trust is only bilateral. In other words, trust will not commute; it will remain between the two CAs. The CAs ensure this by inhibiting policy mapping through the “policy constraints” extension. Policy constraint extensions permit differing policy-mapping inhibitions down the certificate chain. In most direct cross-certifications, policy mapping should be inhibited immediately. In the case of cross-certification using the bridge CA model, in order to take advantage of the policy-mapping services of the bridge CA, the policy-mapping inhibition should be different for one certificate (namely the bridge CA certificate).

In addition, the two CAs should use the “name constraints” extension in the X.509 version 3 certificates to ensure that they trust the other domain for the names over which the other has control. The use of this extension also minimizes the chances of name collision.

Exhibits 37.11 and 37.12 illustrate cross-certification examples. These examples are for illustrative purposes only and do not represent real-world entities.

In the case of bilateral cross-certification, policy mapping should be inhibited immediately by using a value of “0” in the “inhibit policy mapping” field of the policy constraints extension in X.509 certificates. When bridge CA is used for interdomain interoperability, a value of “1” should be used in this field. This will permit the issuing CA domain to map its policies to the bridge CA policies and then permit the bridge CA to map its policies to the subject CA domain, in effect mapping from the issuing CA domain to the subject CA domain.

As long as the issuing CA uses its control on inhibit policy mapping, the bridge CA need not use inhibit policy mapping to control the mapping inhibition.

37.6.7 Public Key Infrastructure Interoperability. The complexity of the technology, standards, and products of PKI technology from one domain to another and from one product to another sometimes creates interoperability problems. Yet without interdomain interoperability there can be no global trust, only individual trust.

These factors play a critical role in ensuring PKI interoperability:

- Trust path
- Cryptographic algorithms

GLOBAL PUBLIC KEY INFRASTRUCTURE 37 · 15

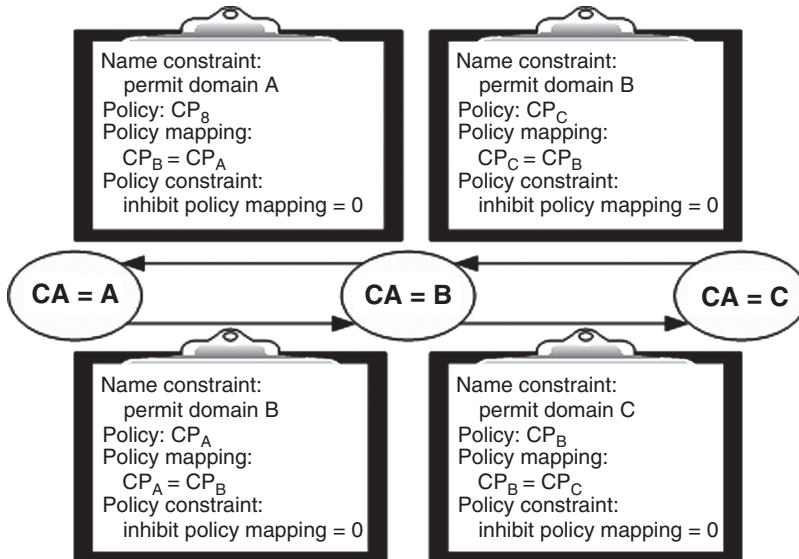


EXHIBIT 37.11 Trust Chain Mapping with Dissimilar Policies

- Certificate and CRL formats
- Certificate and CRL dissemination
- Certificate policies
- Names

37.6.7.1 Trust Path. The communicating parties must be able to form trust paths from their trust anchors to their subscribers. This can be achieved through multiple trust anchors, cross-certification, and other trust models described earlier.

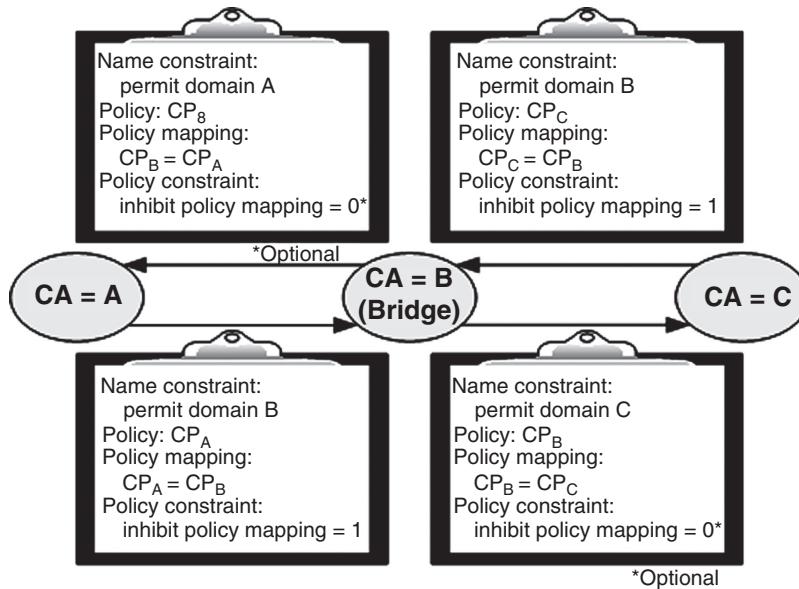


EXHIBIT 37.12 Policy Mapping with a Bridge

37 · 16 PKI AND CERTIFICATE AUTHORITIES

37.6.7.2 Cryptographic Algorithms. The communicating parties must implement the cryptographic algorithms, such as hashing, digital signatures, key encryption, and data encryption, used by each other.

In addition, the parties should be able to communicate to each other the algorithms they use. In X.509 certificates and CRL, this information may be contained in the objects themselves, as in the *algorithm* field. In X.509 certificates, for the information being communicated, algorithms such as the digital signature and key encryption algorithm may be carried in the end-entity certificate. The hashing algorithm and the data encryption algorithm can be part of the implicit agreement between the parties or can be carried with the information being communicated. The information also can be obtained from the *supported algorithms* attribute of the X.500 directory entry of the user, although this option is not widely used.

Although the expected public key algorithms used by the CA to create the certificate must be discernible by the recipient (in order to understand the certificate contents), it is important that the certificate not make any assertions as to the symmetric algorithms that may be used. The certificate should be application-agnostic and not impose any rules on applications when not necessary.

In all these situations, the algorithm is identified using the object identifiers. Different organizations may register the same algorithm under their OID arc. Thus, it is important that either the two domains use the same OID for the algorithms, or that their software interpret the multiple OIDs as the same algorithm. For this reason, OID proliferation for algorithms is not recommended.

Variants of the same base algorithm further exacerbate the problems of algorithm interoperability. For example, subtle padding and other differences exist between definitions of the RSA algorithm in the Public Key Cryptography Standards (PKCS) and in the American National Standards Institute (ANSI) X9 committee. Similarly, the Diffie-Hellman algorithm has various modes and various ways to reduce the calculated secret to symmetric key size (i.e., ways to make session keys smaller). Any of these differences in algorithms must be documented through different OIDs so that the OID invokes the appropriate implementation.

It is important that extensions be chosen carefully. Those pertaining to algorithms often have meaning within the context of a particular application. As mentioned earlier, a good example of an extension that should not be placed in a certificate is the *SMIME Capabilities* extension, since it does not set minimum expectations and instead may limit the symmetric key length that can be used by applications.

Such application-specific extensions should be avoided in user/subscriber certificates. The algorithms and key lengths used by the public key (asymmetric) elements are specified in the *Field* values such as *Subject Public Key Information*.

37.6.7.3 Certificate and Certificate Revocation List Format. The communicating parties must share, or must be able to understand, each other's certificate and CRL formats. The most common way to achieve this is to use a common standard, such as X.509. Many times this has not been sufficient due to the ambiguity in the standard and associated encoding schemes, although, over time, those bugs have been worked out. The primary reason today why certificates and CRLs issued by one product may not be understood by another is that either one or both are not compliant with the standard, or one product does not implement all the features of the standard used by the other product.

GLOBAL PUBLIC KEY INFRASTRUCTURE 37 · 17**37.6.7.4 Certificate and Certificate Revocation List Dissemination.**

The communicating parties must obtain the certificates and CRLs issued by the various CAs in each other's domain. These certificates and CRLs may be obtained from a repository such as an X.500 and Lightweight Directory Access Protocol (LDAP) server. Alternatively, the certificates and CRLs can be carried as part of the communication protocol between the parties, for example, as defined in S/MIME (Secure/Multipurpose Internet Mail Extension) version 3.

The X.500 and LDAP repositories are based on hierarchical databases. Each node in the hierarchical tree structure belongs to an object class. The node's object class determines the attributes that are stored for that node. Examples of attributes are job title, phone number, fax number, and the like. Certificates and CRLs are also attributes.¹¹

X.500 and LDAP have defined a standard schema for PKI certificates and CRLs. For certificates, these attributes are *userCertificate*, *cACertificate*, and *crossCertificatePair*. The end-entity certificates should be stored in *userCertificate* attribute. All CA certificates should be stored in the forward element of the *crossCertificatePair* attribute of the subject CA. In addition, all certificates issued to the CAs in the same domain should be stored in the *cACertificate* attribute of the subject CA.

Various revocation lists should be stored in the *cRL*, *aRL*, and *deltaCRL* attributes of the issuing CA as applicable.

If the certificates and CRLs are not stored in these standardized attributes, the relying-party software may not be able to obtain these objects. Furthermore, X.500 directory products still may not always interoperate due to additional complexity of the X.500 standard and to product differences. When implementing X.500 directories and connecting X.500 directory products from different vendors, implementers should allow time to make the products and directories interoperate.

37.6.7.5 Certificate Policies. In order to trust and cross-certify each other, the CAs in two domains need to operate under similar policies. Users in the two domains should be able to accept or reject certificates of each other's domains based on the security requirements of the application, and on the policy under which the certificates were issued.

In order to determine the similarity, or equivalence, of the policies of the two domains, the CP should be written using the IETF standard RFC-2527 framework. The CP is represented using an OID in the certificate. To ensure that the user software accepts and rejects certificates based on the application requirements and on the CP, PKI products should be selected and configured so that the CA asserts the certificate policies, policy mapping, and policy constraints extensions appropriately. The user's PKI-enabling software must process these extensions appropriately and fully in compliance with the requirements of X.509 certificate-path validation rules.

37.6.7.6 Names. The communicating domains must not assign the same name to two different entities. X.500 distinguished names (DNs) are steps in that direction but not sufficient to achieve this.

To illustrate the point, consider, for example, CygnaCom, a company incorporated in the Commonwealth of Virginia. While it is highly unlikely that there is another CygnaCom in Virginia, there is no assurance that there is no CygnaCom incorporated in other U.S. states. Thus, it would be possible that c=US, O=CygnaCom could be asserted by the CAs for several different domains.

37 · 18 PKI AND CERTIFICATE AUTHORITIES

In order to avoid this name collision and ambiguity, the *name constraints* extension in X.509 should be used. The CA for one domain can prevent any other entity from using a name registered in that domain. The issuing CA (CA “Y” in this example) uses the *name constraints* extension to assert priority and control over the specified identifier. For example, the first CA that certifies a company called *CygnaCom* in its domain should set the *name constraint* attribute in its certificate for its CygnaCom stating that only its CygnaCom is allowed to issue certificates under the name space c=US, O=CygnaCom. If another CygnaCom were to come along, CA “Y” would ask the second CygnaCom to choose another name in order to avoid name collision. Although this example focuses on the DN, the *name constraint* can be used for any hierarchical name forms, including DN, RFC 822-compliant names, and others.

PKI products should be selected and configured so that the CA asserts the *name constraints* extension appropriately. The user’s PKI-enabling software must process this extension appropriately and fully in compliance with the requirements of X.509 certificate-path validation rules.

37.7 FORMS OF REVOCATION. As discussed earlier, a PKI includes mechanisms for key revocation. It is necessary to make provisions for the revocation of compromised keys in order to maintain the trust relationships of any PKI employed in a real-world environment.

The first form of revocation designed was the CRL. It seems the most appropriate form of revocation, given the distributed authentication framework of PKI. The CRL mechanism allows the CA to generate the objects and the relying parties to process them securely without worrying about the security of the servers or system that supply the CRL, and without concerns about the network(s) over which the CRL has traveled.

37.7.1 Types of Revocation-Notification Mechanisms. However, there have been several concerns about the CRL, and these concerns have led to other forms of revocation-notification mechanisms. Many of these mechanisms are variations on the CRL in the sense that these are revocation lists, but they are not complete. The second category of revocation mechanisms defers the processing of revocation information to a server for example, through the Online Certificate Status Protocol (OCSP); see RFC 2560. A third category of mechanisms lets the users check the status of a single certificate from the directory and allows the CA to update the status of that certificate in the directory. A final category lets a CA or another trusted server organize the revocation information in a B-tree.

Which mechanism(s) to choose depends on a variety of factors, such as:

- The communication model (i.e., which class of users is communicating with which other class). For example, if a user communicates with several users who are subscribers to the same CA, a single CRL from that CA will provide relevant information about all those targeted users. If a user is communicating with users who belong to different CAs, each CRL provides information about only one user.
- The directory architecture: Where they are located and what portions of the directory information is replicated or shadowed?
- The communication bandwidth available.
- The bind time (i.e., the time to set a connection with the repository in order to perform retrievals and updates) to access the repository.
- The size of the revocation response from the repository (e.g., the CRL size).

FORMS OF REVOCATION 37 · 19

- The processing load on the repository, especially for digital signature *generation* on the revocation information.
- The processing load on the user workstation, especially for digital signature *verification* on the revocation information.

37.7.2 Certificate Revocation Lists and Their Variants. The first set of mechanisms, CRL and its various forms, is the most versatile, effective, and recommended approach for revocation notification. Like X.509 certificates, CRLs are expressed in ASN.1 format. There are several basic types of CRL, and they should be carefully considered, based on the user communication model and anticipated revocation rate:

- Full and complete CRL
- Authority revocation list (ARL)
- Distribution-point CRL
- Delta CRL

37.7.2.1 Full and Complete CRL. The full and complete CRL is a CRL that contains the revocation information for all certificates issued by a CA. This type of CRL is rarely seen; instead a normal CRL includes only information about revoked certificates and not currently valid ones. Expired certificates are not included, and revoked certificates will drop off the CRL when they expire.

37.7.2.2 Authority Revocation List. The ARL is a CRL that contains the revocation information for all the CA certificates issued by a CA; that is, the ARL is a subset of CRL for certificates issued to the CAs only. The ARL is a very desirable mechanism for these reasons:

- It is likely to be short. A CA is likely to certify fewer CAs than other types of subscribers. Also, given that CAs are expected to operate with a great deal of vigilance, and given that CAs are not going to be revoked for reasons such as name change or organizational affiliation change, CAs will be revoked far less often than the end entities. These factors will contribute to making the ARL very small.
- For all of the certificates except one, only the ARL needs to be checked, since in a certificate path all but the last certificate is issued to a CA.

Due to a security flaw in X.509 version 1, a CA should never issue ARLs defined using that version. In X.509 version 1, there is no difference between the CRL format and the ARL format. Since both CRLs and ARLs are signed by the same CA, if an adversary (directory or network adversary) were to supply an ARL to the relying party in lieu of a full CRL, the relying party would have no way of knowing that it had received an ARL instead of the requested CRL. The ARL would not have end-entity revocation information and therefore could mislead the relying party into using the revoked certificate of an end entity.

The X.509 version 2 ARL fixes this security flaw using an *issuing distribution point* extension. An ARL must use this extension and assert a field that states that the list contains only CA certificates. The presence of this field in the signed ARL tells the

37 · 20 PKI AND CERTIFICATE AUTHORITIES

relying party that it is not a full CRL. Now, if an adversary were to supply an ARL in lieu of a CRL, the relying party would detect this substitution by using the *issuing distribution point* field.

This is one of the several security reasons that PKI-enabling software must be able to process the various extensions properly in accordance with the requirements stated in X.509 standard.

37.7.2.3 Distribution-Point CRL. Distribution-point CRL is a mechanism that has several useful functions:

- To replicate a CRL
- To consolidate revocation information from the various CAs so that the relying parties need to obtain only one CRL
- To partition the revocation information for the subscribers of a CA into multiple smaller pieces

This latter function, partition, is achieved by asserting the *CRL Distribution Point* extension in the certificate that points to the name entry under which revocation information for the certificate will appear. The partitioned CRL will assert the same name in the *Distribution Point* field of the *issuing distribution point* extension in the CRL.

Since all the partitioned (distribution point) CRLs are signed by the same CA, it is not sufficient for the relying party simply to validate the CA's signature on the Distribution Point CRL. The relying party must match the Distribution Point name in the *issuing distribution point* extension of the CRL with the Distribution Point name in the *distribution point* extension in the certificate.

37.7.2.4 Delta Certificate Revocation List. Yet another way to reduce the size of the CRL is to publish changes to the revocation information since the last CRL. The CRL that contains changes only is called the *delta CRL*, and the CRL to which changes are published is called the *base CRL*. The delta CRL can be applied to any of these CRLs: CRL, ARL, and Distribution Point CRL. In order to construct current revocation information, the latest delta CRL and its base must be used. There is an algorithm that could be used to allow a subset of changes to be applied to an earlier CRL that would still match the digital signature of the new CRL.

37.7.3 Server-Based Revocation Protocols. Server-based revocation uses protocols, such as On-Line Certificate Status Protocol (OCSP) and Simple Certificate Validation Protocol (SCVP). In general, these protocols suffer from several flaws, including these:

- Since the revocation information is produced at the server, the communication channel between the relying party and the server must be secured, most likely by using digital signatures.
- Signed operations will limit server scalability since digital signature generation is computationally intensive.
- Since the revocation information is produced at the server, the scheme requires a trusted server as opposed to an untrusted repository.

FORMS OF REVOCATION 37 · 21

- Revocation of a server public key requires a method for checking the server public key status. This method is likely to use the server public key as an additional trust anchor or to rely on a CRL mechanism.
- There needs to be a nonsuppressible mechanism for the CA to provide revocation information to the trusted server; that is, the CA should know whether the revocation information has or has not reached the trusted server. Although a CA itself can act as a trusted server, this is not recommended for security reasons; in addition, we do not want to impose the high-performance requirement on the CA architecture. The trusted server must be a high-performance system.
- There are no standards in the area of CA to provide nonsuppressible mechanisms for transmitting the revocation information to the trusted server.

These mechanisms may be desirable under one of four situations:

1. Need to have thinnest possible PKI client
2. Need to generate revenue for CA services
3. Need to check changing credentials, such as available credit
4. Need to update dynamic credentials, such as the remaining credit line

The last two situations permit the trusted server to provide the revocation information and to check or change the credentials of the subscriber.

Delta CRLs and server-based revocation/authentication protocols such as OCSP (RFC 2560) are standards-compliant and can provide the same information as in the CRL for a single certificate in a significantly smaller bandwidth. They do require some form of acceptable authentication since the original CA will not be available for signing.

37.7.4 Summary of Recommendations for Revocation Notification.

The most scalable and versatile revocation-notification mechanism can be achieved by using a combination of:

- CRLs.
- Replication of the CA directory entry, at locations determined by the enterprise network topology, for fast access to CRL.
- Use of ARLs.
- Consolidation of ARLs for all CAs in a domain through the use of distribution points. Consolidation is achieved by placing the name of a CA that can revoke a certificate in the certificate's CRL *Distribution Point* extension.
- Consolidation of all the reason-codes of key compromise for all certificates in a domain through the use of the *Distribution Point* extension. This CRL can be issued very frequently to meet the freshness requirements of the domain. This mechanism makes the CRL mechanism as current as the OCSP.
- Partitioning routine revocation information using Distribution Point CRLs if CRLs become too large.

37 · 22 PKI AND CERTIFICATE AUTHORITIES

Several other techniques can help improve CRL retrieval efficiency:

- Repositories may store both enciphered CRLs to send to the relying parties and also deciphered (plaintext) CRLs to perform fast searches. Storing both forms reduces the overhead that would result from using encryption or decryption at the time of each request.
- If the repository does not store any private information, bind operations for retrieval can be configured to require no authentication, thus eliminating another potential performance bottleneck.
- CRL size can be reduced by having a short validity period for the certificates, by using a coarse domain name so that reorganization does not invalidate a name, and by allowing some changes (e.g., name change or transfer) without forcing revocation.

37.8 REKEY. The public key certificates for the subscribers have a defined validity period. Once the validity period expires, subscribers require new public key certificates. There are two primary reasons why public key certificates have a limited life. One relates to the life of a private key based on the potential cryptanalysis threat. Another reason is to help control CRL size since no certificate gets off a CRL until it expires.

No public key should be used longer than the estimated time for brute-force cryptanalysis using current technology (its *cryptanalysis threat period*). At that point, the certificate should be assigned a new public key (i.e., it should be *rekeyed*). However, before the cryptanalysis threat period expires, the same key can be renewed or *recertified*. Certificates can be renewed easily by having subscribers send a digitally signed request to the CA or by having the CA perform automatic renewal. During renewal, any information (other than the subscriber public key) may be changed.

For the foreseeable future, the life of a 1024-bit RSA key can be expected to be 25 years. By reusing the same key with new certificates for as long as possible, the number of past keys required to be escrowed in order to retrieve or validate older files is reduced. The level of trust should also be considered when rekeying, but keys stored in tamper-resistant hardware can be expected to achieve their full lifetime.

Note: NIST SP800-78-3 recommends that after 2013 only 2048 bit keys be used.

Elements using soft keys are inherently at risk to a determined attack and should be used only for lower levels of trust.

Certificates also can be rekeyed easily by having the subscriber send a digitally signed *rekey request message* that also contains the new public key. The message is signed using the current private key so that it can be verified using the current public key. If the subscriber being rekeyed is a CA, these requirements also come into play:

- The relying parties should be able to verify certificate chains after the CA is rekeyed.
- The relying parties should be able to verify CRLs issued by the CA.
- The rekey should not have a ripple effect on the PKI. Just because one CA rekeys, other CAs or end entities should not have to rekey.
- The length of certificate paths should be minimized.
- The operational impact on the PKI entities should be minimal.

KEY RECOVERY 37 · 23

A good way to meet these requirements is for the CA to:

- Issue all current valid certificates when it rekeys, without changing the validity periods in the subscriber certificates.
- Continue to sign CRLs with all current valid private keys. This will result in multiple CRLs, all with the same information. A CA private key is considered valid until all certificates signed using that key have expired.

If the CA is a trust anchor, it can use one of two approaches to rekey itself in-band, over the untrusted network:

1. The CA can send out a rekey message that contains its new public key and is signed using the current key. The CA needs to ensure that all the subscribers receive and process the rekey message prior to expiration of the current key.
2. The CA can provide the hash of the next public key and parameters (if the cryptographic algorithm has parameters; RSA does not have parameters, but Digital Signature Standard [DSS] does) with the current key. When it is time to publish the new public key, the CA can publish a new self-signed public key certificate that contains the new public key and parameters as well as the hash of the next public key and parameters.

37.9 KEY RECOVERY. Subscriber public keys can be used to encrypt data-encryption keys (for symmetric key encryption). Such data-encryption keys are used to encrypt data quickly with the lower overhead of symmetric key encryption. Subscribers require their private keys to decrypt the data-encryption keys and thus allow decryption of the data.

It is critically important to distinguish between *signing* keys and *data-encryption* keys. The former (any key that asserts **nonrepudiation**) may *never* be subjected to key recovery; the latter *may* be protected using key-recovery techniques.

High levels of trust keys may require a separate identification key that may not be used to sign a document; this login key should have only the signing bit and not the nonrepudiation bit set in the Key Usage extension.

Sometimes a subscriber's private key token (e.g., a diskette, hard drive, smart card, etc.) may be corrupted or the subscriber may forget the password associated with the token. Similarly, sometimes a subscriber may not be available, yet the subscriber's employer may need to decrypt corporate information encrypted by the missing subscriber. Key-recovery techniques are designed to meet these emergency needs for access to encrypted information. Inherently, they provide a form of back door to the keys, but they also impose additional overhead costs. Thus, the need to provide key recovery should be balanced carefully against potential costs and complexity.

The two most popular forms of key-recovery mechanisms are:

1. **Key escrow.** Under this form, the subscriber's long-term private decryption key is provided to a trusted third party called a *key recovery agent* (KRA)
2. **Key encapsulation.** Under this form, the subscriber encrypts the data-encrypting key using the public key of the KRA so that the KRA can decrypt the data.

Of these two schemes, key escrow is becoming more widely available in the PKI products because it is simpler to implement at the infrastructure level. It is also

37 · 24 PKI AND CERTIFICATE AUTHORITIES

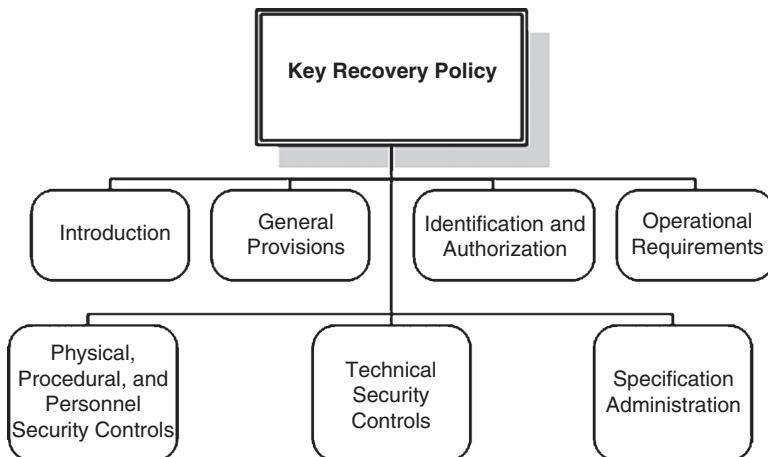


EXHIBIT 37.13 Elements of a Key Recovery Policy

independent of organization boundaries between the sender and receiver of encrypted communications. If a party's private data-encryption key is escrowed, than communications to the party can be decrypted.

Subscribers may always recover their own data-encryption key from the key-recovery system. Authorized third parties, such as a subscriber's employer, also may request keys. Such an authorized party is called a *key recovery requester* (KRR). All of the components are governed by a *key recovery policy* (KRP) and associated *key recovery practices statement* (KRPS). The KRP and KRPS are akin to the *Certificate Policy* and *Certification Practice Statement* but have some differences. One of the main differences is in the technical security-controls sections. There are several requirements to check the communication protocols among the components to ensure confidentiality, integrity, and authorization. Exhibit 37.13 illustrates the components of a key-recovery system, as would be expressed in the KPRS.

A general criticism of key recovery is that it provides secrets to a single party, namely a KRA. One way to mitigate that concern is to share the secret among multiple recipients in a way that requires cooperation (if it is authorized) or collusion (if it is not) among two or more holders of the escrowed secret. For example, superencryption (encryption of ciphertext) can make unauthorized discovery of a key more difficult. The secret key S is encrypted with one recipient's public key (say, K_1), producing a ciphertext, represented as $E(S, K_1)$, and then that ciphertext is super-encrypted using a second recipient's public key, K_2 , to produce the ciphertext $E(E(S, K_1), K_2)$. Unlike encryption of the same message for two recipients, where each recipient can decrypt the ciphertext independently, superencryption requires decryption by each recipient in the reverse order of priority. Thus, if a user encrypts a secret key using A's public key and then superencrypts the ciphertext using B's public key, key recovery requires decryption by B using the corresponding private key and then decryption of the resulting ciphertext by A with that user's private key.

To use superencryption so that fewer than all recipients may decrypt the key, the sender encrypts to a first group of recipients and then superencrypts to a second group of recipients. Thus, any one of the members of the first group and any member of the second group of recipients can cooperate to decrypt the secret key.

This technique allows key escrow even in the absence of a formal PKI, such as in informal webs of trust using PGP.

TRUSTED ARCHIVAL SERVICES AND TRUSTED TIME STAMPS 37 · 25

Another solution to making collusion more difficult in key escrow is to split the key using Shamir's n out of m rule.¹² In a properly implemented key-splitting scheme, parts of a secret key are distributed to m destinations and at least n recipients are required to reconstitute the secret key. Thus, $n-1$ or fewer persons colluding together cannot determine even a single bit of the escrowed key. Successful collusion requires at least n individuals. The split-key approach can be applied to key escrow, in which case the private key can be split and different splits can be provided to different KRAs. Alternatively, the split-key approach may be applied to encapsulation, where the session key can be split, and different splits can be encrypted using public keys of different KRAs.

This strategy also theoretically allows full reconstitution of the secret key to be performed by the authorized recipient of the partial keys, instead of by any of the escrow agents.

37.10 PRIVILEGE MANAGEMENT. The primary purpose of PKI is to provide entity authentication in a global, distributed environment. In most systems and applications, authentication becomes the basis for access control. There are three fundamental ways to implement access control:

1. The systems and applications can perform access control on their own. The PKI continues to provide the authentication framework.
2. The privileges, attributes, roles, rights, and authorizations can be carried in a public key certificate in the *subject directory attribute* extension.
3. The privileges, attributes, roles, rights, and authorizations can be carried in an *attribute certificate*. The X.509 standard is being revised to include the concept of attribute certificates. The attribute certificates carry privileges and authorizations instead of the public key, and thus provide a distributed authorization framework.

User login, signing, and encryption certificates should never be used to carry authorizations or application information. Where needed, separate certificates can be issued for those tasks. Such authorization certificates can be issued daily, or even for shorter periods, eliminating the need for a CRL since the certificate would expire before one would be issued. These authorization certificates could be used in place of Kerberos (the Greek mythological three-headed dog that guards the entrance to Hades) Tickets. Exhibit 37.14 summarizes the pros and cons of each of the three approaches.

These factors should be considered when architecting an access infrastructure. Although the attribute certificate seems to be the latest fad in the PKI world, users should carefully study putting privileges in public key certificates to save the cost of implementing a Privilege Management Infrastructure (PMI) over and above the cost of PKI.

37.11 TRUSTED ARCHIVAL SERVICES AND TRUSTED TIME STAMPS. PKI technology supports global electronic commerce through the use of digital signature technology. Digital signature technology is a detection or passive mechanism. In other words, the technology does not *prevent* someone from modifying data in storage or in transit, nor from impersonating someone else. The technology merely detects that an attempt had been made to modify the data or that someone had tried to impersonate someone else.

37 · 26 PKI AND CERTIFICATE AUTHORITIES

EXHIBIT 37.14 Privilege Management

Alternative	Pros	Cons
Application-based access control	Easy to implement. Does not require additional infrastructure, so saves cost.	Need to manage privileges on an application-by-application basis. Synchronization of privileges may be hard as applications increase and as they are distributed. Security may be compromised if privileges are not removed from all applications. Higher operational costs.
Public key Certificate	Easy to add to PKI. Privileges can be managed easily by revoking certificate.	Changes in privileges require revocation of identity certificate. Sometimes this is a small price to pay for savings that result from not having to deploy and operate a separate privilege management infrastructure (PMI). Parties issuing identity certificate may not have authority to bestow privileges.
Attribute Certificate	Privileges can be managed easily by revoking attribute certificates. Change in privilege does not require revocation of public key certificate.	Cost of privilege management infrastructure (PMI).

In a court of law, digital signatures may be disputed long after they were applied, if, for example, the cryptanalysis threat period for the keys has expired. In those circumstances, producing a document with a verified digital signature may not prevent repudiation. A party could claim that the cryptanalysis threat period has passed and the private key might have been discovered or broken by an adversary.

To mitigate both of these threats—data corruption and expiry of the cryptanalysis threat period—trusted archival services are required for transactions with a potential for this kind of dispute. Such an archival service also would be able to safeguard associated certificates and CRLs. Trusted archival services should depend on controls such as physical security, stable media (e.g., write-once read-many [WORM] devices), and appropriate techniques to maintain readability despite changing technologies. Such services should be capable of error-free transcription of data on older media and of translating outdated encoding into more modern media, and into current encoding schemes. For example, Extended Binary Coded Decimal Interchange Code (EBCDIC)—encoded data on nine-track magnetic tapes could be copied onto optical disks using ASCII encoding.

A related technology is the trusted time stamp, where a trusted third party attaches the current valid time to a document and signs it to prove the existence of the document

FURTHER READING 37 · 27

at a given time. If the document owner does not want to reveal the contents of the document to the time stamp server, then a hash of the document may be stamped and signed instead.

In most applications, a trusted archival service may obviate the need for trusted time stamp service because, for the long term, the trusted archival service can attest to the time of the transaction. For the short term, both parties can date and digitally sign the transaction when it is consummated, rendering it valid as a contract. If the date is not acceptable to any party—it is either too far in the future or in the past—that party can either reject the transaction or invoke some dispute-resolution procedure immediately.

37.12 COST OF PUBLIC KEY INFRASTRUCTURE. One of the misconceptions about Public Key Infrastructure technology is that it is too costly, but these costs should be compared to the alternatives.

Short of taking a major risk, there is no practical technology other than cryptography to protect data in transit over untrusted network. The only choice, then, is between symmetric key cryptography and public key cryptography. Aside from the difficulties of distributing and supporting symmetric, secret keys, such cryptosystems require approximately n^2 keys for a group of n individuals who must communicate with each other.

A secret key cryptosystem requires n^2 keys to be kept confidential, with their integrity maintained; PKI requires managing only n keys. Clearly, it should be cheaper to maintain the integrity of n keys than to manage n^2 keys and their confidentiality and integrity. PKI is needed to manage the public keys.

PKI costs seem large when they must provide global trust and interoperability, something not asked of most systems or infrastructures. Currently, no other technology works as well or is as cost effective as PKI to achieve global, secure, trusted communication and electronic commerce. The alternative is to assume the risk without PKI, or to continue with the last century's approach of paper-based trust.

37.13 FURTHER READING

- Adams, C., S. Lloyd, and S. Kent. *Understanding the Public-Key Infrastructure*, 2nd ed. Upper Saddle River, NJ: Addison-Wesley Professional, 2002.
- Atzeni, A. S., and A. Liou, eds. *Public Key Infrastructure, Third European PKI Workshop: Theory and Practice, EuroPKI 2006, Turin, Italy, June 19–20, 2006. Proceedings*. New York: Springer, 2006.
- Austin, T. *PKI: A Wiley Technical Brief*. New York: John Wiley & Sons, 2000.
- Australian Government Public Key Infrastructure; see: www.govonline.gov.au/projects/publickey/index.asp
- Ballad, B. *Access Control, Authentication, and Public Key Infrastructure*. Jones & Bartlett Learning, 2010
- Barker, E., W. C. Barker, and A. Lee. *Guideline for Implementing Cryptography in the Federal Government*, 2nd ed., SP 800-21. NIST, 2005. <http://csrc.nist.gov/publications/nistpubs/800-21-1/sp800-21-1-Dec2005.pdf> (URL inactive)
- Barker, E., and A. Roginsky. *Recommendation for Cryptographic Key Generation*. SP 800-113. NIST, 2012. <http://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-133.pdf>
- CertiPath. CertiPath X.509 Certificate Policy, version 3.22. CertiPath, 2013. <https://www.certipath.com/images/stories/data/policy-docs/2013-05-22%20CertiPath%20CPv.3.22.pdf>

37 · 28 PKI AND CERTIFICATE AUTHORITIES

- Chadwick, D., and G. Zhao, eds. *Public Key Infrastructure. Second European PKI Workshop: Research and Applications, EuroPKI 2005, Canterbury, UK, June 30–July 1, 2005. Revised Selected Papers*. New York: Springer, 2005.
- Desmedt, Y. G., ed. *Secure Public Key Infrastructure: Standards, PGP and Beyond*, 2nd ed. New York: Springer, 2013.
- ECA Certificate Policies: <http://iase.disa.mil/pki/eca/documents.html>
- Ford, W., and M. S. Baum. *Secure Electronic Commerce: Building the Infrastructure for Digital Signatures and Encryption*. Upper Saddle River, NJ: Prentice-Hall, 1997.
- Housley, R., and T. Polk. *Planning for PKI: Best Practices Guide for Deploying Public Key Infrastructure*. New York: John Wiley & Sons, 2001.
- Howes, T., M. C. Smith, and G. S. Good. *Understanding and Deploying LDAP Directory Services*, 2nd ed. New York: Macmillan, 2003.
- IETF. Public-Key Infrastructure (X.509) (pkix) Working Group Charter, version 4.50. IETF, 2013. <http://datatracker.ietf.org/wg/pkix/charter/>
- Karamanian, A., S. Tenneti, and F. Dessart. *PKI Uncovered: Certificate-Based Security Solutions for Next-Generation Networks*. Cisco Press, 2011.
- Kuhn, D. R., V. C. Hu, W. T. Polk, S.-J. Chang. *Introduction to Public Key Technology and the Federal PKI Infrastructure*. NIST SP800-32, 2001. <http://csrc.nist.gov/publications/nistpubs/800-32/sp800-32.pdf>
- Lopez, J., P. Samarati, and J. L. Ferrer, eds. *Public Key Infrastructure: Fourth European PKI Workshop: Theory and Practice, EuroPKI 2007, Palma de Mallorca, Spain, June 28–30, 2007. Proceedings*. New York: Springer, 2007.
- Menezes, A., P. van Oorschot, and S. Vanstone. *Handbook of Applied Cryptography*. Boca Raton, FL: CRC Press, 1996.
- NIST. *DRAFT A Profile for U. S. Federal Cryptographic Key Management Systems (CKMS)*. SP 800-152. NIST, 2012. <http://csrc.nist.gov/publications/drafts/800-152/draft-sp-800-152.pdf>
- Sabo, J. T., and Y. A. Dzambasow. “PKI Policy in the Business Environment,” *SC Magazine Asia* (June 2002), <http://scmagazine.com/asia/news/article/419807/pkipolicy-business-environment>
- US Department of Defense External Certification Authority X.509 Certificate Policy, version 4.3. DoD, 2012. http://iase.disa.mil/pki/eca/downloads/pdf/eca_cp_v4-3_final_signed.pdf (URL inactive)
- Vacca, J. R. *Public Key Infrastructure: Building Trusted Applications and Web Services*. Boca Raton, FL: Auerbach, 2004.

37.14 NOTES

1. B. Schneier, *Applied Cryptography: Protocols, Algorithms, and Source Code in C*, 2nd ed. (New York: John Wiley & Sons, 1995).
2. Additional details on X.509 certificates can be obtained from X.509 standard ISO/IEC 9594-8.
3. For more information on ASN.1, see <http://asn1.elibel.tm.fr>
4. The I-9 “Employment Eligibility Verification” form, properly known as OMB No. 1615-0047, can be found at: www.uscis.gov/files/form/i-9.pdf
5. These standards can be obtained from www.ietf.org
6. RFC 2527, Internet X.509 Public Key Infrastructure Certificate Policy and Certification Practices Framework: www.ietf.org/rfc/rfc2527.txt

NOTES 37 · 29

7. American Bar Association, "Digital Signatures Guidelines: Legal Infrastructure for Certification Authorities and Electronic Commerce." Draft, 1995. www.abanet.org/scitech/ec/isc/dsgfree.html (URL inactive)
8. Further details on CP and CPS format and contents can be found in S. Chokhani and W. Ford, Certificate Policy and Certification Practices Framework, RFC 2527 (April 1998), www.ietf.org/rfc/rfc2527.txt
9. www.whitehouse.gov/omb/memoranda/fy04/m04-04.pdf
10. New York Office for Technology, "Potential Impacts of Authentication Errors," www.oft.state.ny.us/policy/G07-001/table1.htm (URL inactive)
11. For further discussion on repositories, see T. Howes et al., *Understanding and Deploying LDAP Directory Services* (New York: Macmillan, 1998).
12. A. Shamir, "How to Share a Secret," *Communications of ACM* 22, No. 11 (1979): 612–613, and W. T. Polk, D. F. Dodson, W. E. Burr, H. Ferriolo, and D. Cooper, "Cryptographic Algorithms and Key Sizes for Personal Identity Verification," NIST Special Publication 800-78-3, December, 2010. <http://csrc.nist.gov/publications/nistpubs/800-78-3/sp800-78-3.pdf>

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 38

WRITING SECURE CODE

**Lester E. Nichols, M. E. Kabay, and
Timothy Braithwaite**

38.1 INTRODUCTION	38·1	38.4 TYPES OF SOFTWARE ERRORS	38·10
38.2 POLICY AND MANAGEMENT ISSUES	38·1	38.4.1 Internal Design or Implementation Errors	38·10
38.2.1 Software Total Quality Management	38·2		
38.2.2 Due Diligence	38·4	38.5 ASSURANCE TOOLS AND TECHNIQUES	38·15
38.2.3 Regulatory and Compliance Considerations	38·5	38.5.1 Education Resources	38·15
		38.5.2 Code Examination and Application Testing Techniques	38·16
38.3 TECHNICAL AND PROCEDURAL ISSUES	38·5	38.5.3 Standards and Best Practices	38·20
38.3.1 Requirements Analysis	38·5		
38.3.2 Design	38·6	38.6 CONCLUDING REMARKS	38·20
38.3.3 Operating System	38·6		
38.3.4 Best Practices and Guidelines	38·7	38.7 FURTHER READING	38·21
38.3.5 Languages	38·9	38.8 NOTES	38·22

38.1 INTRODUCTION. The topic of secure coding cannot be adequately addressed in a single chapter. Unfortunately, programs are inherently difficult to secure because of the large number of ways that execution can traverse the code as a result of different input sequences and data values. The additional pressures to get the applications to market can also take a toll on the secure coding process.

This chapter provides a starting point and additional resources for security professionals, system architects, and developers to build a successful and secure development methodology. Writing secure code takes coordination and cooperation of various functional areas within an organization, and may require fundamental changes in the way software development currently is designed, written, tested, and implemented, as well as a change in corporate culture as it pertains to software development.

38.2 POLICY AND MANAGEMENT ISSUES. There are countless security hurdles facing those writing code and developing software. Today dependence on the

38 · 2 WRITING SECURE CODE

reliability and security of the automated system is nearly total. For an increasing number of organizations, distributed information processes, implemented via networked environments, have become the critical operating element of their business. Not only must the processing system work when needed, but also the information processed must retain its integrity so that it can be trusted in use. Because of a general lack of basic IT organizational discipline, and adherence to fundamental software and systems development and maintenance principles, most software development efforts have been significantly lacking at best. If the same inadequate practices are left unchanged in the face of increasing cyber threats, they will continue contributing to the insecure systems of tomorrow, just as they have contributed to the insecure systems of today.

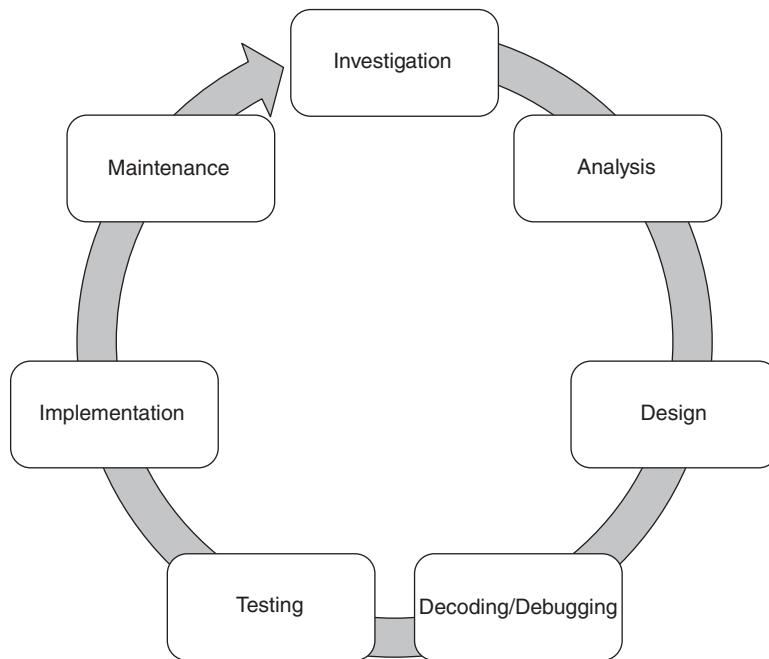
The fundamental problem underlying this kind of assertion is the difficulty of justifying what are often perceived as non-revenue-generating activities and the need to get the product to market on or before delivery timelines. The difficulty also stems from the infeasibility of proving a negative. For example, one cannot prove to the uninvolved or unimpressed observer that the time, money, and human resources spent on prevention are well spent simply because a disaster or incident was averted. Such an observer will take the fact that it did not occur as evidence that the possibility never really existed, or that the problem was exaggerated. In the same way, a skeptical observer may take the absence of security breaches as evidence that no such threat existed, or that it was never as bad as claimed. The problem is exacerbated after spending time and effort on establishing security policies, security controls, and enforcement mechanisms.

In both cases, the money, time, and effort expended to prevent adverse computing consequences from harming an organization are viewed with suspicion and not qualifying at a level to justify the change in delivery schedule or funding. However, such suspicions are not new in the world of information processing. Throughout the history of IT, this attitude has existed regarding problem-prevention activities such as quality assurance, quality control testing, documentation, configuration management, change management, and other system-management controls designed to improve quality, reliability, and system maintainability.

If you consider the analogy of application development as being similar to the way a house is built, it is much harder to implement and integrate something after the fact. Security is not and should not be considered a superficial exercise similar to hanging the drapes or changing a light bulb. It is like the plumbing of a house, and it needs to be implemented with that in mind.

Problems with security are management problems. Management dictates how much time and money can be spent on security activities such as risk management and software testing. How much time and energy are invested in designing, building, testing, and fixing software is in the end a management call. Thus, the primary concern of those striving to improve software security must be the correction of these same basic and fundamental IT management defects; otherwise effective security can never be attained or maintained.

38.2.1 Software Total Quality Management. The nature of the computer security challenge is that of dynamically changing vulnerabilities and threats, based on rapidly changing technologies, used in increasingly complex business applications. Organizations need to build software around the idea of a continuous improvement model, as popularized by the total quality management (TQM) articulated in the ISO 9000 family of standards. For purposes of quality improvement, this model was viewed as a continuing cycle of plan-do-check-and-act. For software security, it can be thought of as a cycle of plan-fix-monitor-and-assess. A security management process to

POLICY AND MANAGEMENT ISSUES 38 · 3**EXHIBIT 38.1** Common Lifecycle Process

execute this model must be established and integrated into the day-to-day operations of the business. Integration means that the security initiative does not stand alone as a watchdog or act merely as an audit function, and it must move away from a policing function. Security must be considered an essential part of all other systems development and management activities, including requirements definition, systems design, programming, package integration, testing, systems documentation, training, configuration management, operations, and maintenance. Security must be viewed as good architecture. In other words, security must be a part of the entire development lifecycle and management gates of the project. Exhibit 38.1 provides a common lifecycle process as described in Chapter 39.3.1.

A security management process must be designed in such a way that the activities of planning, fixing, monitoring, and assessing are accomplished in an iterative fashion for systems being defined, designed, and implemented and for systems that are currently executing the day-to-day business operations.

For many organizations, it has not been uncommon for software to be built with little or no attention to security concerns. Consequently, many business-critical applications have made their way into full production with inadequate or nonexistent security controls in place. For example, a quickly built application to resolve an initial need for tracking information within a database through a Web interface may not get all the access control configurations designed into the application that are needed for proper security. It serves a quick purpose but becomes a long-term tool because of its success. Because of the nature of business (moving from one task to the next), nothing is really done to develop it into the function it now fills. In addition, the longer the application stays in use, the greater the likelihood that the original developers are to transition, and any historical knowledge or documentation will be lost. This makes the ability to update or correct the security of the application less desirable. As a result, users or

38 · 4 WRITING SECURE CODE

hackers may be unintentionally or intentionally capable of accessing and modifying critical information within the database or within the system itself. Fortunately, this trend has started to change, but it is still in need of significant attention. Because such systems are already operating, the only way realistically to identify and assess their security risks outside of replacing the application outright with an application that includes the security requirements is to begin aggressively monitoring for intrusions and then to design and implement corrective security policies and controls based on what is discovered. The effective sequence of security activities for systems already in operational status will be to monitor current operations, assess suspicious activity, plan a corrective policy and technology control, implement the control on the system, and then monitor for continuing effectiveness.

38.2.2 Due Diligence. To address the software and code security effectively, it must become an issue for management, just as other security concerns have begun to demand management attention. Business executives and IT managers are now all too familiar with the concept of due diligence as applied to the uses of information technology. Because of the potential for legal fallout, it has become necessary to view IT-related actions through the definitions of due diligence and reasonable care.

The significance of the concepts of due diligence and reasonable care is that they allow for an evolving metric against which an organization's application security deliberations, decisions, and actions can be compared. The comparison is usually against a similar organization, in like circumstances of vulnerability and threat, and with similar predictable adverse impacts on customers, partners, shareholders, employees, and the public.

For example, if one company employs a security control and does not experience any security breaches that the technique was supposed to prevent, that fact could be used to establish a baseline against which other similar organizations could be compared. If enough similar organizations employ the same technique, the security control may become categorized as a best practice for that industry.

If, however, another company in the same industry did not employ the security control and did experience security breaches of the type the technique was suppose to prevent, it might clearly indicate a lack of due diligence or reasonable care. Software security decisions and implementations are not a one-time event, but need to be under a continuous process of risk evaluation, management, and improvement. The decisions start with the selection of the development language and continue from there. The choice of an interpreted versus compiled language could have an impact not only to development skills required and timelines, but how the use of the languages libraries are used and the impact future updates can be accomplished by the customer. A good example is the recent¹ security issues that have plagued Java and Microsoft products.

It is therefore imperative, in order to demonstrate an ability to exercise continual due diligence to management, to establish a documented computer security risk management program, regardless of internal or external software development, and to integrate it into the overall management processes of the software development process. Nothing less will demonstrate that a company, and its board, is capable of assessing computer security threats and of acting in a reasonable manner.

For software developers working on commercial software for sale to a specific client under contract, or even to a wide range of customers governed by end-user license agreements, the exercise of due diligence and improved security may reduce

TECHNICAL AND PROCEDURAL ISSUES 38 · 5

the risk of lawsuits claiming damages for negligence. Whether successful or not, such lawsuits are never positive publicity for makers of software products.

38.2.3 Regulatory and Compliance Considerations. Regulatory and compliance considerations have become a significant aspect to most businesses. The additional aspect of software development increases the need to be aware of the implication state, federal, or international regulation may have on a development project. Compliance with nongovernmental industries may also need to be considered.

If an organization is subject to specific regulations (e.g., Sarbanes-Oxley legislation), then documentation of the errors encountered, and the resulting internal reporting and remediation efforts, is critical. The error should clearly indicate how it was identified, who identified it, how it was reported to management, what the remediation will be, and when it is anticipated to be completed. Without these details, management may encounter significant difficulties in confirming that adequate internal control mechanisms exist, together with the appropriate and adequate involvement of management. In addition, an error could result in a control weakness being identified by the external auditors or, even worse, as a material weakness (a material weakness would be the use of a programming language that has a known inherent flaw rather than another more appropriate and without the flaw, similar to using aluminum instead of steel for a bridge), depending on its nature and severity.

Additionally, when errors are noted that affect multilocation systems or applications, the significance and materiality—that is to say the extent to which the errors impact or may impact the system or any other system to which the application interfaces—must be considered from the aspect of both the subsidiary and the headquarters locales.

Since laws differ among states and countries, what might be considered legally acceptable standards of privacy or accounting in one locale might not be acceptable elsewhere. As a result, careful consideration of the requirements and intended uses or customers must be established as part of the project initiation and due diligence processes.

38.3 TECHNICAL AND PROCEDURAL ISSUES. Security should be integrated into every stage of the application life cycle. This includes the requirements analysis, the design stage for software, and the operating system security kernel, in corporate policy development, and in human awareness, training, and education programs.

Writing secure code can be a daunting and technically challenging undertaking. That is not to say that it cannot be done, but it requires the developer to work in conjunction with the rest of the project team and other key stakeholders to meet the challenges that must be defined, reconciled, and overcome prior to the release of the software or system. The additional push to market that can drive the rapid development process can also have a negative impact on the security of the application. The problem must be dealt with in two ways, technical and procedural. Technically, the programmers must be aware of the pitfalls associated with application development and avoid them. Procedurally, the development team and organization need to adhere to a consistent methodology of development. This consistency also needs to include the identification of security risks at every stage of the process, including the requirements analysis.

38.3.1 Requirements Analysis. Security must be a part of the requirements analysis within any development project. Most application developers know that adding

38 · 6 WRITING SECURE CODE

security after the fact increases cost and time, and becomes more complicated. Requirements can be derived from different sources:

- Functional needs of the system (this includes the operating system the application will run on)
- National, international, or organizational standards or guidelines
- Regulatory restrictions
- Sensitivity level of data
- Existing security policies
- Cost/benefit analysis
- Risk rating

The purpose of analysis is to determine what information and processes are needed to support the desired objectives and functions of the software and the business model.

38.3.2 Design. The informational data gathered during the analysis gathering goes into the software design as requirements. What comes out of the analysis are the data, logic, and procedural design.

Developers take the data and the informational model data and transform them into the data structures that will be required to implement the software. The logic design defines the relationships between the major structures and components of the application. The procedural design transforms structural components into descriptive procedures. Access control mechanisms are chosen, rights and permissions are defined, any other security specifications are appraised, and solutions are determined. A work breakdown structure includes the development and implementation stages. The structure includes a timeline and detailed activities for testing, development, staging, integration testing, and product delivery.

The decisions made during the design phase are pivotal to application development. The design is the only way the requirements are translated into software components. It is in this way that software design is the foundation of the development process, and it greatly affects software quality and maintenance. If good product design is not put in place in the beginning, the rest of the development process will be that much more challenging.

38.3.3 Operating System. The operating system security kernel is responsible for enforcing the security policy within the operating system and the application. As such, the architecture of the kernel operating system is typically layered, with the kernel at the most privileged layer. This is a small portion of the operating system, but all references to information and changes to authorization pass through the kernel.

To be secure the kernel must meet three basic criteria:

1. **Completeness.** All access to information must go through the kernel.
2. **Isolation.** The kernel must be protected from unauthorized access.
3. **Verifiability.** The kernel must be proven to meet design specifications.

As mentioned, the kernel runs in the most privileged layer. Most operating systems have two processor access modes, user and kernel. General application code runs in user mode, while the operating system runs in kernel mode. Kernel mode allows

TECHNICAL AND PROCEDURAL ISSUES 38 · 7

the processor full access to all system memory and CPU instructions. When applications are not written securely, this separation of modes can become compromised, enabling exploitation of vulnerabilities through arbitrary code, buffer overflows, and other techniques.

38.3.4 Best Practices and Guidelines. “Best practices” is a term that can cause vigorous debate, especially with regard to security. Best practices in general, and particularly with regard to security, often fall prey to a dogmatic, and sometimes blind, devotion to nonsensical practices that have little to do with security and more to do with faith or tradition. Regardless, best practices help provide a set of guidelines that help provide structure. In addition, best practices can be adapted to help meet the needs of a particular situation. An example would be the National Institute of Standards and Technology (NIST) Special Publication Series 800. This series provides best practices and recommendations that can be adapted or integrated into other practices to assist in improving and eliminating the tradition to a particular practice that may no longer have a legitimate use. Considering this, most general security textbooks contain recommendations on security-related aspects of programming—see, for example, Stallings—that do in fact have a very real benefit in creating more secure software.

In addition to designing security into a system from the start, there are also some obvious guidelines that will help in developing secure software:

- Impose strong identification and authentication (I&A) for critical and sensitive systems in addition to the I&A available from the operating system; ideally, use token-based or biometric authentication as part of the initialization phase of your application.
- Document your code thoroughly, including using data dictionaries for full definition of allowable input and output to functions and allowable range and type of values for all variables.
- Use local variables, not global variables, when storing sensitive data that should be used only within a specific routine (i.e., use the architecture of the process stack to limit inadvertent or unauthorized access to data in the stack).
- Reinitialize temporary storage immediately after the last legitimate use for the variable, thus making scavenging harder for malefactors.
- Limit functionality in a specific module to what is required for a specific job (e.g., do not use the same module for supervisory functions and for routine functions carried out by clerical staff).
- Define views of data in databases that conform to functional requirements and limit access to sensitive data (e.g., the view of data from a medical records database should exclude patient identifiers when a worker in the finance department is using the database for statistical aggregation).
- Use strong encryption (*not* homegrown encryption) that has industry-standard routines to safeguard sensitive and critical data on disk. Locally developed, home-grown encryption is generally not as safe.
- Disallow access by programmers to production databases.
- Randomize or otherwise mask sensitive data when generating test subsets from production data.

38 · 8 WRITING SECURE CODE

- Use test-coverage monitors to verify that all sections of source code are in fact exercised during quality assurance tests; investigate the functions of code that never is executed.
- Integrate logging capability into all applications for debugging work, for data recovery after crashes in the middle of a transaction, and for security purposes such as forensic analysis.
- Create log-file records that include a cryptographically sound message authentication code (MAC) that itself includes the MAC of the preceding record as input for the algorithm; this technique ensures that forging a log file or modifying it will be more difficult for a malefactor.
- Log all process initiations for a program and log process termination; include full details of who loaded the program or module.
- Log all modifications to records and optionally provide logging for read access as well.
- Use record-level locking to prevent inadvertent overwriting of data on records that are accessed concurrently. Be sure to unlock a sequence of locks in the inverse order of the lock sequence to prevent deadlocks. (Thus, if you lock resource A, B, and C in that order, unlock C, then B, then A.)
- Sign your source code using digital signatures.
- Use checksums in production executables to make unauthorized modifications more difficult to conceal.

Mike Gerdes, former manager at AtomicTangerine, contributed these suggestions in the course of a discussion:

- Recommend the readers adopt a practice of designing code in a more holistic fashion. A common practice is to write and test routines in a way that verifies the code processes the data in the way intended. To avoid the effects of malicious code and data input attacks, the programmer must also write code that deals with what is *not* supposed to be processed. A more complete design methodology would also include testing of all inbound information to ensure exclusion of any data which did not fit the requirements for acceptable data. This method should be applied to high-risk applications and those with an extremely arduous test cycle and will eliminate many of the common attack methods used today.
- Establish the criteria for determining the sensitivity level of information contained in, or processed by, the application and subroutines.
- If they are not already present, consider implementing formal control procedures in the software programming methodology to ensure that all data is reviewed during QA processes and to be sure it is classified and handled appropriately for the level assigned.
- Identify and include any mandatory operating system, and network security characteristics for the production system in the specifications of the software. In addition to providing the development and QA teams some definition of the environment the software is designed to run in, giving the administrator and end users an idea of what your expectations were when you created the code can be extremely useful in determining where software can, or cannot, be used.

TECHNICAL AND PROCEDURAL ISSUES 38 · 9

- Where appropriate, verify the digital signatures of routines that process sensitive data when the code is being loaded for execution.
- If you include checksums on executables for production code, include routines that verify the checksums at every system restart.

The Carnegie-Mellon Software Engineering Institute Computer Emergency Response Team (CERT) provides a good “Top 10 Secure Coding Practices”² list. The list below provides some additional items to consider in addition to the items previously listed:

- Validate input from data sources. Proper input validation can eliminate many software vulnerabilities.
- Compile code using the highest warning level available for the compiler and eliminate warnings by modifying the code.
- Create a software architecture and design the software to implement and enforce security. Carnegie-Mellon has pre-established programming standard by language that can be found at www.cert.org/secure-coding/scstandards.html
- Remember the K-I-S-S (Keep It Simple, Stupid) principle. Complex designs increase the likelihood of errors or vulnerabilities and the need for security mechanisms or resulting mitigations to become more complex.
- Base access decisions on permission rather than exclusion to adhere to the principle of least privilege.
- Sanitize data sent to other systems, such as command shells, relational databases, and commercial off-the-shelf (COTS) components.
- Practice defense in depth within the application, so that if one layer of defense turns out to be inadequate, another layer of defense can prevent or limit a security flaw from becoming an exploitable vulnerability.
- Establish good quality assurance techniques to be more effective in identifying and eliminating vulnerabilities.
- Develop and/or apply a secure coding standard for your target development language and platform.
- Define security requirements at the start of the project.
- Think like a hacker and use threat modeling to anticipate the threats to which the software may be subjected. If you consider how it can be attacked, you will be able to prevent the attack before it can occur.

38.3.5 Languages. To date, no common computer languages have security specific features built-in. Java does include provisions for limiting access to resources outside the “sandbox” reserved for a process, as described in the books by Felten and McGraw. Nevertheless, Java is still one of the languages targeted as an attack vector by many malware implementations. PASCAL uses strong typing and requires full definition of data structures, thus making it harder to access data and code outside the virtual machine defined for a given process. In contrast, C and C++ allow programmers to access any region of memory at any time the operating system permits it.

Computer languages allow developers to write code as well as they can. Strongly typed languages may offer better constraints on programmers, but the essential

38 · 10 WRITING SECURE CODE

requirement is that the programmers continue to think about security as they design and build code.

There are several sets of security utilities and resources available for programmers; for example, RSA has a number of cryptographic toolkits. Some textbooks (e.g., Schneier's *Applied Cryptography*) include CD-ROMs with sample code. In addition, the Computer Emergency Response Team Coordination Center (CERT/CC) was started in December 1988 by the Defense Advanced Research Projects Agency (DARPA), which was part of the U.S. Department of Defense. CERT/CC is located at the Software Engineering Institute, a federally funded research center operated by Carnegie Mellon University.

CERT/CC studies Internet security vulnerabilities, provides services to Websites that have been attacked, and publishes security alerts. CERT/CC's research activities include the area of WAN computing and developing improved Internet security.

As mentioned in Section 38.3.4, Carnegie-Mellon has established secure coding standards based on language that can help establish a uniform set of principles and guidelines across your development group. Regardless of the language selected for software development, it is important to make the decision based on the strengths and weaknesses of the language. This means that a risk analysis needs to be a part of the design of any development project, including the selection of the development language. By understanding those strengths and weaknesses, better decisions can be made to prevent the inherent vulnerabilities within the language from becoming exposed inadvertently during development.

38.4 TYPES OF SOFTWARE ERRORS. New programmers should review the range of root causes for software errors. Such review is particularly useful for students who have completed training that did not include discussions of systematic quality assurance methodology. The next sections can serve as a basis for creating effective sets of test data and test procedures for unit tests of new or modified routines.

38.4.1 Internal Design or Implementation Errors. A general definition of a software error is a mismatch between a program and its specifications; a more specific definition is the failure of a program to do what the end user reasonably expects. There are many types of software errors. Some of the most important include:

- Initialization
- Logic flow
- Calculation
- Boundary condition violations
- Parameter passing
- Race condition
- Load condition
- Resource exhaustion
- Resource, address, or program conflict with the operating system or application(s)
- Regulatory compliance considerations
- Other errors

TYPES OF SOFTWARE ERRORS 38 · 11

38.4.1.1 Initialization. Initialization errors are insidious and difficult to find. The most insidious programs save initialization information to disk and fail only the first time used—that is, before they create the initialization file. The second time a given user activates the program, there are no further initialization errors. Thus, the bugs appear only to employees and customers when they activate a fresh copy of the defective program. Other programs with initialization errors may show odd calculations or other flaws the first time they are used or initialized; because they do not store their initialization values, these initialization errors will continue to reappear each time the program is used.

38.4.1.2 Logic Flow. Modules pass control to each other or to other programs. If execution passes to the wrong module, a logic-flow error has occurred. Examples include calling the wrong function, or branching to a subroutine that lacks a RETURN instruction, so that execution falls through the logical end of a module and begin executing some other code module.

38.4.1.3 Calculation. When a program misinterprets complicated formulas and loses precision as it calculates, it is likely that a calculation error has occurred; for example, an intermediate value may be stored in an array with 16 bits of precision when it needs 32 bits. This category of errors also includes computational errors due to incorrect algorithms.

38.4.1.4 Boundary Condition Violations. Boundaries refer to the largest and smallest values with which a program can cope; for example, an array may be dimensioned with 365 values to account for days of the year, and then fail in a leap year when the program increments the day-counter to 366 and thereby attempts to store a value in an illegal address. Programs that set variable ranges and memory allocation may work correctly within the boundaries but, if incorrectly designed, may crash at or outside the boundaries. The first use of a program also can be considered a boundary condition.

One of the most important types of boundary violations is the buffer overflow. In this error, data placed into storage exceed the defined maximum size and overflow into a section of memory identified as belonging to one or more different variables. The consequences can include data corruption (e.g., if the overflow overwrites data that are interpreted as numerical or literal values) or changes in the flow of execution (e.g., if the altered data include logical flags that are tested in branch instructions).

Buffer overflows have been exploited by writers of malicious code who insert data that overflows into memory areas of interpreted programs and thus become executed as code.

All programs should verify that the data being stored in an array or in any memory location do not exceed the expected size of the input. Data exceeding the expected size should be rejected or, at least, truncated to prevent buffer overflow.

38.4.1.5 Parameter Passing. Sometimes there are errors in passing data back and forth among modules. For instance, a call to a function accidentally might pass the wrong variable name so that the function acts on the wrong values. When these parameter-passing errors occur, data may be corrupted, and the execution path may be affected because of incorrect results of calculations or comparisons. As a result, the

38 · 12 WRITING SECURE CODE

latest changes to the data might be lost, or execution might fall into error-handling routines even though the intended data were correct.

38.4.1.6 Race Condition. When a race occurs between event A and event B, a specific sequence of events is required for correct operation, but the program does not ensure this sequence. For example, if process A locks resource 1 and waits for resource 2 to be unlocked while process B locks resource 2 and waits for resource 1 to be unlocked, there may be a deadly embrace that freezes the operations if the processes overlap in execution. If they do not happen to overlap, there is no problem at that time.

Race conditions can be expected in multiprocessing systems and interactive systems, but they can be difficult to replicate; for example, the deadly embrace just described might happen only once in 1,000 transactions if the average transaction time is very short. Consequently, race conditions are among the most difficult to detect during quality assurance testing and are best identified in code reviews. Programmers should establish and comply with standards on sequential operations that require exclusive access to more than one resource. For example, if all processes exclusive-lock resources in a given sequence and unlock them in the reverse order, there can be no deadly embrace.

38.4.1.7 Load Condition. All programs and systems have limits to storage capacity, numbers of users, transactions, and throughput. Load errors are caused by exceeding the volume limitations of storage, transactions, users, and networks can occur due to high volume, which includes a great deal of work over a long period, or high stress, which includes the maximum load all at one time. For example, if the total theoretical number of transactions causes a demand on the disk I/O system that exceeds the throughput of the disk controller, processes will necessarily begin piling up in a queue waiting for completion of disk I/Os. Although theoretical calculations can help to identify where possible bottlenecks can occur in CPU, memory, disk, and network resources, a useful adjunct is automated testing that permits simulation of maximum loads defined by service-level agreements.

38.4.1.8 Resource Exhaustion. The program running out of high-speed memory (RAM), mass storage (disk), central processing unit (CPU) cycles, operating system table entries, semaphores, network bandwidth, or other resources can cause failure of the program. For example, inadequate main memory may cause excessive swapping of data to disk (thrashing), typically causing drastic reductions in throughput, because disk I/O is typically 1,000 times slower than memory access.

38.4.1.9 Interapplication Conflicts. With operating systems (OS) as complex as they are, OS manufacturers routinely distribute the code requirements and certain parameters to the application software manufacturers, so that the likelihood of program conflicts or unexpected stoppages are minimized. While this certainly helps reduce the number of problems and improves the forward and backward compatibility with previous OS versions, even the OS vendors on occasion experience or cause difficulties when they do not conform to the parameters established for their own programs.

38.4.1.10 Other Sources of Error. It is not unusual for errors to occur where programs send bad data to devices, ignore error codes coming back, and even try to use devices that are busy or missing. The hardware might well be broken, but the

TYPES OF SOFTWARE ERRORS 38 · 13

software also is considered to be in error when it does not recover from such hardware conditions.

Additional errors can occur through improper builds of the executable; for example, if an old version of a module is linked to the latest version of the rest of the program, the wrong sign-on screens may pop up, the wrong copyright messages may be displayed, the wrong version numbers may appear, and various other inaccuracies may occur.

38.4.1.11 User Interface. Generally speaking, the term “user interface” denotes all aspects of a system that are relevant to a user. It can be broadly described as the user virtual machine (UVM). This would include all screens, the mouse and keyboard, printed outputs, and all other elements with which the user interacts. A major problem arises when system designers cannot put themselves in the user’s place and cannot foresee the problems that technologically challenged users will have with an interface designed by a technically knowledgeable person.

Documentation is a crucial part of every system. Each phase of development—requirements, analysis, development, coding, testing, errors, error solutions and modifications, implementation, and maintenance—needs to be documented. All documents and their various versions need to be retained for both future reference and auditing purposes. Additionally, it is important to document the correct use of the system and to provide adequate instructional and reference materials to the user. Security policies and related enforcement and penalties also need to be documented. Ideally, the documentation should enable any technically qualified person to repair or modify any element, as long as the system remains operational.

38.4.1.12 Functionality. A program has a functionality error if performance that can reasonably be expected is confusing, awkward, difficult, or impossible. Functionality errors often involve key features or functions that have never been implemented. Additional functionality errors exist when:

- Features are not documented.
- Required information is missing.
- A program fails to acknowledge legitimate input.
- There are factual errors or conflicting names for features.
- There is information overload.
- The material is written to an inappropriate reading level.
- The cursor disappears, or is in the wrong place.
- Screen displays are wrong.
- Instructions are obscured.
- Identical functions require different operations in different screens.
- Improperly formatted input screens exist.
- Passwords or other confidential information are not obscured or protected adequately.
- Tracing the user data entry or changes is unavailable or incomplete.
- Segregation of duties is not enforced. (This can be particularly critical for organizations subject to legal and regulatory requirements.)

38 · 14 WRITING SECURE CODE

38.4.1.13 Control (Command) Structure. Control structure errors can cause serious problems because they can result in:

- Users getting lost in a program.
- Users wasting time because they must deal with confusing commands.
- Loss of data or the unwanted exposure of data.
- Work delay.
- Financial cost.
- Unanticipated exposure to data leakage or compromise; this can result in significant liability if consumers' personal identifying information (PII) is compromised.
- Data not being encrypted as intended or being visible to unauthorized users.

Some common errors include:

- Inability to move between menus.
- Confusing and repetitive menus.
- Failure to allow adequate command-line entries.
- Requiring command-line entries that are neither intuitive nor clearly defined on screen.
- Failure of the application program to follow the operating system's conventions.
- Failure to distinguish between source and parameter files, resulting in the wrong values being made available to the user through the interface, or failure to identify the source of the error
- Inappropriate use of the keyboard, when new programs do not meet the standard of a keyboard that has labeled function keys tied to standard meanings.
- Missing commands from the code and screens resulting in the user being unable to access information, to utilize programs, or to provide for the system to be backed up and recoverable. There are a host of other commands that can leave the system in a state of less-than-optimum operability.
- Inadequate privacy or security that can result in confidential information being divulged, in the complete change or loss of data without recoverability, in poor reporting, and even in undesired access by outside parties.

38.4.1.14 Performance. Speed is important in interactive software. If a user feels that the program is working slowly, that can be an immediate problem.

Performance problems include slow response, unannounced case sensitivity, uncontrollable and excessively frequent automatic saves, inability to save, and limited scrolling speed.

There are five fundamental potential bottlenecks to consider in analyzing any computer and network performance problem relating to access to and speed of:

- The central processing unit(s)
- Main memory
- Secondary memory (e.g., disk drives)

ASSURANCE TOOLS AND TECHNIQUES 38 · 15

- Application design
- Network bandwidth

Slow operation can depend on (but is not limited to) the OS, the other applications running, memory saturation and thrashing (excessive swapping of memory contents to virtual memory on disk or on other, relatively slow, memory resources such as flash drives), memory leakage (the failure to de-allocate memory that is no longer needed), disk I/O inefficiencies (e.g., reading single records from very large blocks), and program conflicts (e.g., locking errors). For more information on program performance.

At another level, performance suffers when program designs make it difficult to change their functionality in response to changing requirements. In a database design, defining a primary index field that determines the sequence in which records are stored on disk can greatly speed access to records during sequential reads on key values for that index—but it can be counterproductive if the predominant method for accessing the records is sequential reads on a completely different index.

38.4.1.15 Output Format. Output format errors can be frustrating and time consuming. An error is considered to have occurred when the user cannot change fonts, underlining, boldface, and spacing that influence the final look of the output; alternatively, delays or errors when printing or saving document may occur. Errors occur when the user cannot control the content, scaling, and look of tables, figures, and graphs. Additionally, there are output errors that involve expression of the data to an inappropriate level of precision.

38.5 ASSURANCE TOOLS AND TECHNIQUES. The security community has known about buffer overflow vulnerabilities, as an example, for the past 25 or 30 years, yet the vulnerabilities are still increasing. For example, the National Vulnerability Database reported 418 new buffer overflow errors in 2013 compared with 392 in 2011. Clearly, greater efforts are necessary to teach and to enforce awareness, implementation, and testing of source code.

38.5.1 Education Resources. Build Security In (BSI) is a project of the Software Assurance program of the Strategic Initiatives Branch of the National Cyber Security Division (NCSD) of the U.S. Department of Homeland Security. BSI (<http://buildsecurityin.us-cert.gov>) contains and links to best practices, tools, guidelines, rules, principles, and other resources that software developers, architects, and security practitioners can use to build security into software in every phase of its development. BSI content is based on the principle that software security is fundamentally a software engineering problem and must be addressed in a systematic way throughout the software development life cycle.

In addition to BSI, and out of the multitudes of sites on the Web that address some facet of secure coding, some sites are listed next to assist in secure application development:

- Carnegie-Mellon Software Engineering Institute: www.cert.org/secure-coding
- FreeBSD Security Information: www.freebsd.org/security/security.html
- International Systems Security Engineering Association (ISSEA): www.issea.org
- Open Web Application Security Project: www.owasp.org

38 · 16 WRITING SECURE CODE

- Secure, Efficient, and Easy C Programming: <http://irccrew.org/~cras/security/c-guide.html>
- Secure Programming for Linux and UNIX HOWTO: <http://www.dwheeler.com/secure-programs/>
- Windows Security: www.windowsecurity.com
- World Wide Web Security FAQ: <http://www.w3.org/Security/Faq/>

38.5.2 Code Examination and Application Testing Techniques. Scanning once is not enough; ongoing application assessment is essential to implementing effective secure development practices and, in turn, a secure application. Like many other aspects of information technology and computer security, a layered approach to testing is also needed. As a result, there are numerous techniques and approaches, both manual and automated, to software security testing. This section will provide a cursory overview of the major techniques and approaches.

38.5.2.1 Static Analysis. Static analysis or source code review is the process of manually checking source code for security weaknesses. Sometimes looking at the code is the only way to resolve issues. Understanding the intent of the inner workings of the application can identify serious security issues that cannot be identified any other way.

Some issues that are particularly conducive for discovery through static analysis include concurrency problems, flawed business logic, access control problems, backdoors, Trojans, Easter eggs, time/logic bombs, and other forms of malicious code.

38.5.2.2 Binary Code Analysis. Binary code analysis utilizes tools that support analysis of binary executables such as decompilers, disassemblers, and binary code scanners. This is to say that binary analysis is the reverse engineering of the executable. The least intrusive technique is binary scanning. Binary scanners analyze machine code to model a language-neutral representation of the program's behaviors, control/data flows, call trees, and external function calls. An automated vulnerability scanner then traverses the model in order to locate vulnerabilities caused by common coding errors and simple back doors.

The next least intrusive technique is disassembly of the binary. The binary code is reverse engineered to an intermediate assembly language, but the disadvantage of disassembly is that it requires an expert who both thoroughly understands that particular assembler language and who is skilled in detecting security-relevant constructs within assembler code.

The most intrusive reverse engineering technique is decompilation of the code. The binary code is reverse engineered completely to source code, which can then be subjected to the same security code review techniques and other white box tests as original source code. Note, however, that decompilation is technically problematic in that the quality of the source code generated through decompilation is often very poor. As a result, code is rarely as navigable or comprehensible as the original source code. The analysis of decompiled source code will always be more difficult and time consuming than that of the original source code. For this reason, the use of decompilation for security analysis should only be used for the most critical systems.

Another point to consider regarding reverse engineering is that it may also be legally prohibited. A majority of software vendors' license agreements prohibit reverse engineering. In addition, many software vendors repeatedly cite the Digital Millennium

ASSURANCE TOOLS AND TECHNIQUES 38 · 17

Copyright Act (DMCA) of 1999 as justification, even though the DMCA explicitly exempts reverse engineering and encryption research from the prohibitions against copy protection circumvention.

38.5.2.3 White Box Approach. White box testing strategy deals with the internal logic and structure of the code. White box testing is also known as glass, structural, open box, or clear box testing. In order to implement white box testing, the tester has to deal with the code, and therefore needs to possess knowledge of coding and logic (i.e., internal working of the code). White box tests also need the tester to look into the code and to find out which unit/statement/chunk of the code is malfunctioning.

Advantages of white box testing are:

- It becomes very easy to find out which type of input/data can help in testing the application effectively.
- It helps in optimizing the code.
- It helps in removing the extra lines of code, which can bring in hidden defects.

Disadvantages of white box testing are:

- A skilled tester is needed to carry out this type of testing, which increases the cost.
- It is nearly impossible to look into every bit of code to find hidden errors that may create problems, resulting in failure of the application.

38.5.2.4 Black Box Approach. Black box testing is testing without knowledge of the internal workings of the item being tested. Black box testing is also known as behavioral, functional, opaque box, and closed box. For this reason, the tester and the programmer can be independent of one another, avoiding programmer bias toward his own work. Test groups are often used. Due to the nature of black box testing, the test planning can begin as soon as the specifications are written.

Advantages of black box testing are:

- It is more effective on larger units of code than glass box testing.
- The tester needs no knowledge of implementation, including specific programming languages.
- Tester and programmer are independent of each other.
- Tests are done from a user's point of view.
- It will help to expose any ambiguities or inconsistencies in the specifications.
- Test cases can be designed as soon as the specifications are complete.

Disadvantages of black box testing are:

- Only a small number of possible inputs can actually be tested; to test every possible input stream on a complex system would take nearly forever.
- Without clear and concise specifications, test cases are hard to design.
- There may be unnecessary repetition of test inputs if the tester is not informed of test cases the programmer has already tried. It is always critical to validate any

38 · 18 WRITING SECURE CODE

testing claims through reports or result output to help avoid duplication. Likewise, it is always important to rerun tests that may have ambiguous results.

- It may leave many program paths untested.
- It cannot be directed toward specific segments of code that may be very complex.

38.5.2.5 Gray Box—Blended Approach. Gray box testing is a software testing technique that uses a combination of black box testing and white box testing. Gray box testing is not black box testing, because the tester does know some of the internal workings of the software under test. In gray box testing, the tester applies a limited number of test cases to the internal workings of the software under test. In the remaining part of the gray box testing, the tester takes a black box approach in applying inputs to the software under test and observing the outputs.

38.5.2.6 Fault Injection. There are two primary categories of fault injection: source code or binary fault injection. Fault injection is used to induce stress in the software, create interoperability problems among components, simulate faults in the execution environment, and thereby revealing safety-threatening faults that are not made apparent by traditional testing techniques. Security extends standard fault injection by adding error injection, enabling testers to analyze the security impacts of the behaviors, and state changes resulting in the software when exposed to changes in the environment data. These data changes are meant to simulate the types of faults that would occur during unintentional user errors or intentional attacks on the software through its environment, including attacks on the environment itself. A data change is simply an alteration of the data the execution environment passes to the software, or between software components. Fault injection reveals both the effects of security faults on individual component behaviors, and the behavior of the system as a whole.

Binary fault injection is useful when performed as part of penetration testing to enable the tester to obtain a more holistic picture of how the software responds to attacks. However, faults should not be limited to simulating real world attacks. Fault injection scenarios should be designed to give a complete understanding of the security impacts that the application's behaviors, states, and properties exhibit under all possible operating conditions.

The main challenges in fault injection testing are determining the most meaningful combination and volume of faults to inject and, as with all tests, interpreting the results. Furthermore, just because fault injection does not cause the software to behave in a nonsecure manner or fail in a nonsecure state, this cannot be interpreted to mean that the software will perform equally when exposed to the more complex inputs it typically receives during execution.

38.5.2.7 Fuzzing. Fuzz testing or fuzzing is a technique that is often automated or semi-automated, that involves providing invalid, unexpected, or random data to the inputs of a computer program. The program is then monitored for exceptions, failing built-in code assertions, or for finding potential memory leaks. There are two forms of fuzzing, mutation-based and generation-based, which can be employed as part of white, black, or grey box testing.

The idea behind fuzz testing is to look for interesting program behavior that results from noise injection and may indicate the presence of vulnerabilities or faults. Fuzzing

ASSURANCE TOOLS AND TECHNIQUES 38 · 19

applications, fuzzers, are generally specific to a particular type of input and written to test a specific program; they are not easily reusable, but file formats and network protocols are the most common targets of testing. Any type of program input can be fuzzed, including items not normally considered “input” (e.g., databases, shared memory, etc.). The value is the specificity, that is to say, fuzzing can often reveal security vulnerabilities that generic testing tools such as vulnerability scanners and fault injectors cannot. On the other hand, completely random fuzzing is a typically ineffective way to uncover problems in an application, and as a result fuzzing technology has evolved to include more intelligent techniques.

38.5.2.8 Vulnerability Scanning. Automated vulnerability scanning is an automated scanning technique for application level software, as well as for Web servers, database management systems, and some operating systems. The tools scan the application software for input and output of known patterns associated with known vulnerabilities. The vulnerability patterns, or signatures, are similar to the signatures used by virus scanners, or the coding constructs searched for by automated source code scanners, making the vulnerability scanner, essentially, an automated pattern-matching tool. Some Web application vulnerability scanners additionally attempt to perform automated stateful assessment using simulated reconnaissance attack patterns and fuzzing techniques to probe for known or common vulnerabilities.

Like any signature-based scanner, vulnerability scanners can report false positives. As a result, the tester must have enough software and security expertise to meaningfully interpret the scanner’s results to limit the false positives or negatives, to prevent the creation of an incident where it does not exist. This is why it is important to use a layered approach and combine different test techniques to examine the software for vulnerabilities.

Vulnerabilities in software are often masked by environmental protections such as network and application-level firewalls. Furthermore, environment conditions may create unique vulnerabilities that cannot be found by a signature-based tool, but using a combination of other black box tests, most especially penetration testing of the deployed software in the actual target environment. Vulnerability scanners are most effective as part of:

- Security assessments of binary components;
- Or before penetration testing (to eliminate the need to run penetration test scenarios for common vulnerabilities).

38.5.2.9 Software Penetration Testing. Penetration testing has been a common technique to test network security for many years and is also known as black box testing and ethical hacking. As described in the black box approach, penetration testing is essentially the testing of a running application (without knowing the inner workings of the application) to identify vulnerabilities. In penetration testing, the whole software system is a “live” environment. Typically, the penetration team would have access to the application like typical users. The testers then act like an attacker, attempting to find and exploit potential vulnerabilities.

Penetration testing should focus on aspects of system behavior, interaction, and vulnerability that cannot be observed through other testing techniques performed outside of a live environment. That means that the penetration should subject the system to

38 · 20 WRITING SECURE CODE

sophisticated multi-pattern attacks designed to trigger complex behaviors across system components. In addition, the testing should attempt to find security problems likely to originate in the architecture and design, rather than coding flaws, since this type of problem tends to be overlooked by other testing techniques. This also means test plan should include worst-case scenarios that reproduce threat vectors considered highly damaging, such as an insider threat scenario.

Many people use application penetration testing as the primary testing technique and while it has a place in a testing program; it should not be considered the primary or only testing technique. Penetration testing can potentially lead to a false sense of security. Just because an application passes a penetration test, it does not mean it is vulnerability free. Likewise, if an application fails, it is a strong indication that there are serious problems with it. However, focused penetration testing, which is testing that attempts to exploit vulnerabilities detected in previous tests, can be useful in detecting if those vulnerabilities are actually fixed within the code.

38.5.3 Standards and Best Practices. Testing is a complete software engineering discipline in its own right. Volumes have been written about the various techniques that software engineers use to test and validate their applications. Some basic best practices can be employed regardless of the testing methodology selected. These are:

- Perform automated testing.
- Make a test plan.
- Follow a specific methodology.
- Test at every stage.
- Test all system components.

In addition to these best practices, use of these standards can provide additional resources:

- ISO 17799, Information Technology: Code of Practice for Information Security Management
- ISO/IEC 15408, Evaluation Criteria for IT Security (the Common Criteria)
- SSE-CMM, System Security Engineering Capability Maturity Model
- ISO/IEC WD 15443, Information Technology: Security Techniques

38.6 CONCLUDING REMARKS. Computer security practitioners and top-level managers must understand that IT management deficiencies, left unchanged, will sabotage efforts to establish and sustain effective computer security programs. Furthermore, it is incumbent on those coding to help steer the project and educate management toward a model that is inclusive of security throughout the development life cycle.

The costs of neglecting security at the start will continue to be a serious issue until software security is fully integrated into every aspect of the development life cycle. Security must be built in, and building security into systems starts with the software development teams, practices, and techniques used to build those systems.

FURTHER READING 38 · 21**38.7 FURTHER READING**

- Campanella, J., ed. *Principles of Quality Costs: Implementation and Use*, 3rd ed. Milwaukee, Wisconsin: ASQ Quality Press, 1999.
- Eilam, E. *Reversing: Secrets of Reverse Engineering*. Indianapolis, IN: Wiley Publishing, 2005.
- Felten, E., & G. McGraw. *Securing Java: Getting Down to Business with Mobile Code*. New York: John Wiley & Sons, 1999. Also free and unlimited Web access from www.securingjava.com
- Fox, C., and P. Zonneveld. *IT Control Objectives for Sarbanes-Oxley: The Role of IT in the Design and Implementation of Internal Control Over Financial Reporting*, 2nd ed. IT Governance Institute, 2006.
- Hoglund, G., and G. McGraw. *Exploiting Software: How to Break Code*. Boston: Addison-Wesley, 2004.
- Institute of Electrical and Electronics Engineers. *IEEE Standard for Software Test Documentation*. ANSI/IEEE Std. 829-1983. 1983.
- Institute of Electrical and Electronics Engineers. *IEEE Standard for Software Quality Assurance Plans*. ANSI/IEEE Std. 730-1981. 1981.
- Long, F., D. Mohindra, R. C. Seacord, D. F. Sutherland, and D. Svoboda. *Java Coding Guidelines: 75 Recommendations for Reliable and Secure Programs*. Addison-Wesley Professional, 2013
- Martin, R. et al. *Clean Code: A Handbook of Agile Software Craftsmanship*. New York: Prentice-Hall, 2008.
- McConnell, S. *Code Complete*, 2nd ed. Redmond, WA: Microsoft Press, 2004.
- McGraw, G., and E. W. Felten. *Java Security: Hostile Applets, Holes and Antidotes—What Every Netscape and Internet Explorer User Needs to Know*. New York: John Wiley & Sons, 1997.
- McGraw, G., and E. W. Felten. “Understanding the Keys to Java Security—the Sandbox and Authentication,” *JavaWorld*. May 1, 1997. www.javaworld.com/javaworld/jw-05-1997/jw-05-security.html.
- Markantonakis, K., and K. Mayes, eds. *Secure Smart Embedded Devices, Platforms and Applications*. Springer, 2013.
- NASA Software Assurance Technology Center. <http://satc.gsfc.nasa.gov> (URL inactive)
- Ould, M.A. *Strategies for Software Engineering: The Management of Risk and Quality*. New York: John Wiley & Sons, 1990.
- Open Web Application Security Project (OWASP), www.owasp.org
- RSA Data Security, www.rsasecurity.com/products
- Sebesta, R. *Concepts of Programming Languages*. 9th ed. Boston: Pearson Education, 2010.
- Schneier, B. *Applied Cryptography: Protocols, Algorithms, and Source Code in C*, 2nd ed. New York: John Wiley & Sons, 1995.
- Secord, R. C. *Secure Coding in C and C++*, 2nd ed. Addison-Wesley Professional, 2013
- Shore, J., and J. Warden. *The Art of Agile Development*. Sebastopol: O'Reilly Media, 2008.
- Stallings, W. *Network and Internetwork Security: Principles and Practice*. Englewood Cliffs, NJ: Prentice Hall, 1995.
- Tipton, H. F., and K. Henry. *Official (ISC) Guide to the CISSP CBK*. Boca Raton, FL: Auerbach Publications, 2006.

38 · 22 WRITING SECURE CODE

Wilhelm, T. et al. *Ninja Hacking: Unconventional Penetration Testing Tactics and Techniques*. New York: Syngress, 2011.

38.8 NOTES

1. Consider the following articles: “Ms. Smith,” “Microsoft Mega-Patch Closes Critical IE Flaws, Fixes 57 Vulnerabilities,” NetworkWorld, February 12, 2013, www.networkworld.com/community/blog/microsoft-mega-patch-closes-critical-ie-hole-fixes-57-vulnerabilities?source=NWWNL_E_nlt_security_2013-02-13 or Danielle Walker, “Oracle Speaks, Promises To Get Java ‘Fixed Up.’” SC Magazine, January 28, 2013, www.scmagazine.com/oracle-speaks-promises-to-get-java-fixed-up/article/277898/?DCMP=EMC-SCUS_Newswire
2. CERT, “Top 10 Secure Coding Practices,” last edited by Robert Seacord on March 1, 2011, www.securecoding.cert.org/confluence/display/seccode/Top+10+Secure+Coding+Practices

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 39

SOFTWARE DEVELOPMENT AND QUALITY ASSURANCE

**Diane E. Levine, John Mason, and
Jennifer Hadley**

39.1	INTRODUCTION	39·2	39.4	TYPES OF SOFTWARE ERRORS	39·8
39.2	GOALS OF SOFTWARE QUALITY ASSURANCE	39·2		39.4.1 Internal Design or Implementation Errors	39·8
39.2.1	Uncover All of a Program's Problems	39·2		39.4.2 User Interface	39·10
39.2.2	Reduce the Likelihood that Defective Programs Will Enter Production	39·3	39.5	DESIGNING SOFTWARE TEST CASES	39·12
39.2.3	Safeguard the Interests of Users	39·3		39.5.1 Good Tests	39·12
39.2.4	Safeguard the Interests of Software Producers	39·3		39.5.2 Emphasize Boundary Conditions	39·13
39.3	SOFTWARE DEVELOPMENT LIFE CYCLE	39·3		39.5.3 Check All State Transitions	39·14
39.3.1	Phases of the Traditional Software Development Life Cycle	39·4	39.6	39.5.4 Use Test-Coverage Monitors	39·14
39.3.2	Classic Waterfall Model	39·6		39.5.5 Seeding	39·15
39.3.3	Rapid Application Development and Joint Application Design	39·7		39.5.6 Building Test Data Sets	39·15
39.3.4	Extreme Programming	39·8	39.7	39.5.7 Fuzzing	39·16
39.3.5	Importance of Integrating Security at Every Phase	39·8		BEFORE GOING INTO PRODUCTION	39·16
				39.6.1 Regression Testing	39·16
				39.6.2 Automated Testing	39·16
				39.6.3 Tracking Bugs from Discovery to Removal	39·17
				MANAGING CHANGE	39·17
				39.7.1 Change Request	39·17
				39.7.2 Tracking System	39·18
				39.7.3 Regression Testing	39·18
				39.7.4 Documentation	39·18

39 · 2 SOFTWARE DEVELOPMENT AND QUALITY ASSURANCE

39.8 SOURCES OF BUGS AND PROBLEMS	39·19	39.8.4 Insufficient or Substandard Programming Quality	39·19
39.8.1 Design Flaws	39·19	39.8.5 Data Corruption	39·19
39.8.2 Implementation Flaws	39·19	39.8.6 Hacking	39·20
39.8.3 Unauthorized Changes to Production Code	39·19	CONCLUSION	39·21
		39.10 FURTHER READING	39·21

39.1 INTRODUCTION. Software development can affect all of the six fundamental principles of information security as described in Chapter 3 in this *Handbook*, but the most frequent problems caused by poor software involve integrity, availability, and utility.

Despite the ready availability of packaged software on the open market, such software frequently has to be customized to meet its users' particular needs. Where this is not possible, programs must be developed from scratch. Unfortunately, during any software development project, despite careful planning, unforeseen problems inevitably arise. Custom software and customized packages frequently are delivered late, are faulty, and do not meet specifications. Generally, software project managers tend to underestimate the impact of technical as well as nontechnical difficulties. Because of this experience, the field of software engineering has developed; its purpose is to find reasonable answers to questions that occur during software development projects.

39.2 GOALS OF SOFTWARE QUALITY ASSURANCE. In the IEEE *Glossary of Software Engineering Terminology*, quality is defined as "the degree to which a system, component, or process meets customer or user needs or expectations." In accordance with this definition, software should be measured primarily by the degree to which user needs are met. Because software frequently needs to be adapted to changing requirements, it should be adaptable at reasonable cost. Therefore, in addition to being concerned with correctness, reliability, and usability, the customer also is concerned with testability, maintainability, portability, and compliance with the established standards and procedures for the software and development process. Software quality assurance (SQA) is an element of software engineering that tries to ensure that software meets acceptable standards of completeness and quality. SQA acts as a watchdog overseeing all quality-related activities involved in software development.

The principal goals of SQA are:

- Uncover all of a program's problems.
- Reduce the likelihood that defective programs will enter production.
- Safeguard the interests of users.
- Safeguard the interests of the software producer.

39.2.1 Uncover All of a Program's Problems. Quality must be built into a product. Testing can only reveal the presence of defects in the product. To ensure quality, SQA monitors both the development process and the behavior of software. The goal is not to pass or certify the software; the goal is to identify all the inadequacies and problems in the software so that they can be corrected.

SOFTWARE DEVELOPMENT LIFE CYCLE 39 · 3

39.2.2 Reduce the Likelihood that Defective Programs Will Enter Production. All software development must comply with the standards and policies established within the organization to identify and eliminate errors before defective programs enter production.

39.2.3 Safeguard the Interests of Users. The ultimate goal is to achieve software that meets the users' requirements and provides them with needed functionality. To meet these goals, SQA must review and audit the software development process and provide the results of those reviews and audits to management. SQA, as a functioning department, can be successful only if it has the support of management and if it reports to management at the same level as software development. Likewise, each SQA project will possess specific attributes, and the SQA program should be tailored to accommodate project needs. Characteristics to be considered include: mission criticality, schedule and budget, size and complexity of project, and size and competence of project staff organization.

39.2.4 Safeguard the Interests of Software Producers. By ensuring that software meets requirements, SQA can help prevent legal conflicts that may arise if purchased software fails to meet contractual obligations. When software is developed in-house, SQA can prevent the finger-pointing that otherwise would damage relations between software developers and corporate users.

39.3 SOFTWARE DEVELOPMENT LIFE CYCLE. Good software, whether developed in-house or bought from an external vendor, needs to be constructed using sound principles. Because software development projects often are large and have many people working on them for long periods of time, the development process needs to be monitored and controlled.

Although progress of such projects is difficult to measure, using a phased approach can control the projects. In a phased approach, a number of clearly identifiable milestones are established between the start and the finish of the project. A common analogy that is used is that of constructing a house, where the foundation is laid initially and each phase of construction is achieved in an orderly and controlled manner. Frequently, with both house construction and software development, payments for phases are dependent on reaching the designated milestones.

When developing systems, we refer to the phased process as the system development life cycle (SDLC). The SDLC is the process of developing information systems through investigation, analysis, design, coding and debugging, testing, implementation, and maintenance. These seven phases are common, although different models and techniques may contain more or fewer phases.

Generally, milestones identified in the SDLC correspond to the points in time when specified documents become available. Frequently, the documents explain and reinforce the actions taken during the just-completed phase of the SDLC. We therefore say that traditional models for the phased development are document driven to a large extent.

There are problems and drawbacks inherent in the document-driven process. The method of viewing the development process via the SDLC is not totally realistic to the actual projects. In reality, errors found in earlier phases are noted, and fixes are developed, prototyping is introduced, and solutions are implemented. This, in effect, is more than what is assumed will be necessary in the debugging and maintenance phases. Also, often the problems are solved before those later phases are reached. Because of

39 · 4 SOFTWARE DEVELOPMENT AND QUALITY ASSURANCE

the recognition that much of what in traditional models is referred to as maintenance is really evolution, other models of the SDLC, known as evolutionary models, have been developed.

In traditional models, the initial development of a system is kept strictly separate from the maintenance phase. The major goal of the SDLC is to deliver a first production version of the software system to the user. It is not unusual for this approach to result in excessive maintenance costs to make the functioning or final system fit the needs of the real user. There are other models available, and it is necessary, for each project, to choose a specific SDLC model. To do this, it is necessary to identify the need for a new system and to follow this by identifying the individual steps and phases, possible interaction of the phases, the necessary deliverables, and all related materials.

There are five major system development techniques:

1. Traditional systems development life cycle (SDLC)
2. Waterfall model (a variant of the traditional SDLC)
3. Rapid application development (RAD)
4. Joint application development (JAD)
5. Extreme programming

Although these five development techniques frequently are seen as mutually exclusive, in truth they represent solutions that place different emphasis on common elements of systems design. Defined at different times, each methodology's strengths demonstrate the technology, economics, and organizational issues that were current at the time the methodology was first defined.

Software generally is constructed in phases, and tests should be conducted at the end of each phase before development continues in the next phase. Four major phases that are always present in software construction are:

1. The analysis phase, where software requirements are defined
2. The design phase, based on the previously described requirements
3. The construction, programming, or coding phase
4. The implementation phase, where software is actually installed on production hardware and finally tested before release to production

The programming phase always includes unit testing of individual modules and system testing of overall functions. Sometimes, in addition to testing during programming and implementation, a fifth phase, solely devoted to functional testing, is established. Review and modification generally follow all phases, where weaknesses and inadequacies are corrected.

39.3.1 Phases of the Traditional Software Development Life Cycle.

The seven phases discussed here comprise the most common traditional SDLC:

1. Investigation
2. Analysis
3. Design
4. Decoding and debugging

SOFTWARE DEVELOPMENT LIFE CYCLE 39 · 5

5. Testing
6. Implementation
7. Maintenance

From this traditional model, other models with fewer or more phases have been developed.

39.3.1.1 *Investigation.* This phase involves the determination of the need for the system. It involves determining whether a business problem or opportunity exists and conducting a feasibility study to determine the cost effectiveness of the proposed solution.

39.3.1.2 *Analysis.* Requirements analysis is the process of analyzing the end users' information needs, the environment of the organization and the current system, and developing functional requirements for a system to meet users' needs. This phase includes recording all of the requirements. The documentation must be referred to continually during the remainder of the system development process.

39.3.1.3 *Design.* The architectural design phase lists and describes all of the necessary specifications for the hardware, software, people and data resources, and information products that will satisfy the proposed system's functional requirements. The design can best be described as a blueprint for the system. It is a crucial tool in detecting and eliminating problems and errors before they are built into the system.

39.3.1.4 *Coding and Debugging.* This phase involves the actual creation of the system. It is done by information technology professionals, sometimes on staff within a company and sometimes from an external company that specializes in this type of work. The system is coded, and attempts are made to catch and eliminate all coding errors before the system is implemented.

39.3.1.5 *Testing.* Once the system is created, testing is essential. Testing proves the functionality and reliability of the system and acts as another milestone for finding problems and errors before implementation. This phase is instrumental in determining whether the system will meet users' needs. The testing should be documented in a project plan format as to the

- Objective of the testing
- Scope
- Expected results
- Test plan
- Test cases; ideally, the end users and/or business process owner(s) should be involved so that the test cases represent real-life situations as much as possible
- Test results and analysis
- Conclusion and determination of next steps
- End user/business process owner [signed] acceptance of the test results

39 · 6 SOFTWARE DEVELOPMENT AND QUALITY ASSURANCE

These steps help ensure that the tests align with the users' needs and that the results are representative of what can be expected.

39.3.1.6 Implementation. Once the previous five phases have been completed and accepted, with substantiating documentation, the system is implemented and users are permitted to utilize it.

39.3.1.7 Maintenance. This phase is ongoing and takes place after the system is implemented. Because systems are subject to variances, flaws, and breakdowns, as well as difficulties when integrating with other systems and hardware, maintenance continues for as long as the system is in use.

39.3.2 Classic Waterfall Model. Although phased approaches to software development appeared in the 1960s, the waterfall model, attributed to Roy Royce, appeared in the 1970s. The waterfall model demands a sequential approach to software development and contains only five phases:

1. Requirements analysis phase
2. Design
3. Implementation
4. Testing
5. Maintenance

39.3.2.1 Analysis or Requirements Analysis. The waterfall model emphasizes analysis as part of the requirements analysis phase. Since software is always part of a larger system, requirements are first established for all system elements, and then some subset of these requirements is allocated to software. Identified requirements, for both the system and the software, are then documented in the requirements specification. A series of validation tests is required to ensure that, as the system is developed, it continues to meet these specifications.

The waterfall model includes validation and verification in each of its five phases. This means that in each phase of the software development process, it is necessary to compare the obtained results against the required results. Testing is done within every phase to answer this question and does not occur strictly in the testing phase that follows the implementation phase.

39.3.2.2 Design. The design phase is actually a multistep process focusing on data structure, software architecture, procedural detail, and interface characterization. During the design phase, requirements are translated into a representation of the software that then can be assessed for quality before actual coding begins. Once again, documentation plays an important role because the documented design becomes a part of the software configuration. Coding is incorporated into the design phase instead of being a separate phase. During coding, the design is translated into machine-readable form.

39.3.2.3 Implementation. During the implementation phase, the system is given to the user. Many developers feel that a large problem with this model is that the system is implemented before it actually is ready to be given to the user. However, waterfall model advocates claim that the consistent testing throughout the development

SOFTWARE DEVELOPMENT LIFE CYCLE 39 · 7

process permits the system to be ready by the time this phase is reached. Note that in the waterfall model, the implementation phase precedes the testing phase.

39.3.2.4 Testing. Although testing has been going on during the entire development process, the testing phase begins after code has been generated and the implementation phase has occurred. This phase focuses on both the logical internals of the software and the functional externals. Ideally, the end users and the process/business owners should be enlisted to assist in this phase, since they usually:

- Are familiar with the business requirements
- Know how the current software and hardware operates
- Know what types of activity are anticipated and planned
- Provide the ultimate acceptance approval of the software and hardware

Although not directly required by many of the regulatory reviews, involving the end users and process owners provides the IT area with additional resources, not only from a human resource perspective but also from a business knowledge and regulatory compliance perspective.

39.3.2.5 Maintenance. The maintenance phase reapplies each of the preceding life cycle steps to existing programs. Maintenance is necessary because of errors that are detected, necessary adaptation to the external environment, and functional and performance enhancements required and requested by the customer.

The waterfall model is considered to be unreliable, partly due to failure to obey the strict sequence of phases advocated by the traditional model. In addition, the waterfall model often fails because it delivers products whose usefulness is limited when requirements have changed during the development process.

39.3.3 Rapid Application Development and Joint Application Design.

Rapid application development (RAD) supports the iteration and flexibility necessary for building robust business process support. RAD emphasizes user involvement and small development teams, prototyping of software, software reuse, and automated tools. Activities must be carried out within a specified time frame known as a time box. This approach differs from other development models where the requirements are fixed first and the time frame is decided later. RAD, in an effort to keep within the time box and its immovable deadline, may sacrifice some functionality.

RAD has four phases within its cycle:

1. Requirements planning
2. User design
3. Construction
4. Cutover

The main techniques used in RAD are joint requirements planning (JRP) and joint application design (JAD). The word “joint” refers to developers and users working together through the heavy use of workshops. JAD enables the identification, definition, and implementation of information infrastructures. The JAD technique is discussed with RAD because it enhances RAD.

39 · 8 SOFTWARE DEVELOPMENT AND QUALITY ASSURANCE

39.3.4 Extreme Programming. A form of *agile software development* that has gained popularity in the last two decades is *extreme programming* (XP). Programmers work closely with each other and the users to produce usable code quickly. Unit tests and integrated integration tests are constructed before the coding and constantly applied throughout the development process. A variant of the method is known as *Scrum*.

39.3.5 Importance of Integrating Security at Every Phase. Security should never be something added to software at the end of a project. Security must be considered continuously during an entire project in order to safeguard both the software and the entire system within which it functions. Regardless of which software development model is used, it is essential to integrate security within each phase of development and to include security in the testing at each phase. Rough estimates indicate that the cost of correcting an error rises tenfold with every additional phase. For example, catching an error at the analysis phase might require only several minutes—time for a user to correct the analyst—and therefore may cost only a few dollars. Catching the same error once it has been incorporated into the specifications document might cost 10 times more; after implementation, as much as 1,000 times more (this excludes any compliance-related penalties, losses, or loss responsibility).

39.4 TYPES OF SOFTWARE ERRORS

39.4.1 Internal Design or Implementation Errors. A general definition of a software error is a mismatch between a program and its specifications; a more specific definition is “the failure of a program to do what the end user reasonably expects.” There are many types of software errors. Some of the most important include:

- Initialization
- Logic flow
- Calculation
- Boundary condition violations
- Parameter passing
- Race condition
- Load condition
- Resource exhaustion
- Resource, address, or program conflict with the operating system or application(s)
- Regulatory compliance considerations
- Other errors

39.4.1.1 Initialization. Initialization errors are insidious and difficult to find. The most insidious programs save initialization information to disk and fail only the first time they are used—that is, before they create the initialization file. The second time a given user activates the program, there are no further initialization errors. Thus, the bugs appear only to employees and customers when they activate a fresh copy or installation of the defective program. Other programs with initialization errors may show odd calculations or other flaws the first time they are used or initialized; because

TYPES OF SOFTWARE ERRORS 39 · 9

they do not store their initialization values, these initialization errors will continue to reappear each time the program is used.

39.4.1.2 Logic Flow. Modules pass control to each other or to other programs. If execution passes to the wrong module, a logic-flow error has occurred. Examples include calling the wrong function, or branching to a subroutine that lacks a RETURN instruction, so that execution falls through the logical end of a module and begins executing some other code module.

39.4.1.3 Calculation. When a program misinterprets complicated formulas and loses precision as it calculates, it is likely that a calculation error has occurred; for example, an intermediate value may be stored in an array with 16 bits of precision when it needs 32 bits. This category of errors also includes computational errors due to incorrect algorithms.

39.4.1.4 Boundary Condition Violations. The term “boundaries” refers to the largest and smallest values with which a program can cope; for example, an array may be dimensioned with 365 values to account for days of the year and then fail in a leap year when the program increments the day-counter to 366 and thereby attempts to store a value in an illegal address. Programs that set variable ranges and memory allocation may work within the boundaries but, if incorrectly designed, may crash at or outside the boundaries. The first use of a program also can be considered a boundary condition.

39.4.1.5 Parameter Passing. Sometimes there are errors in passing data back and forth among modules. For instance, a call to a function accidentally might pass the wrong variable name so that the function acts on the wrong values. When these parameter-passing errors occur, data may be corrupted and the execution path may be affected because of incorrect results of calculations or comparisons. As a result, the latest changes to the data might be lost or execution might fall into error-handling routines even though the intended data were correct.

39.4.1.6 Race Condition. When a race occurs between event A and event B, a specific sequence of events is required for correct operation, but this sequence is not ensured by the program. For example, if process A locks resource 1 and waits for resource 2 to be unlocked while process B locks resource 2 and waits for resource 1 to be unlocked, there will be a deadly embrace that freezes the operations.

Race conditions can be expected in multiprocessing systems and interactive systems, but they can be difficult to replicate; for example, the deadly embrace just described might happen only once in 1,000 transactions if the average transaction time is short. Consequently, race conditions are among the least tested.

39.4.1.7 Load Condition. All programs and systems have limits to storage capacity, numbers of users, transactions, and throughput. Load errors can occur due to high volume, which includes a great deal of work over a long period of time, or high stress, which includes the maximum load all at one time.

39.4.1.8 Resource Exhaustion. The program’s running out of high-speed memory (random access memory, or RAM), mass storage (disk), central processing unit (CPU) cycles, operating system table entries, semaphores, or other resources

39 · 10 SOFTWARE DEVELOPMENT AND QUALITY ASSURANCE

can cause failure of the program. For example, inadequate main memory may cause swapping of data to disk, typically causing drastic reductions in throughput.

39.4.1.9 Interapplication Conflicts. With operating systems (OS) as complex as they are, OS manufacturers routinely distribute the code requirements and certain parameters to the application software manufacturers, so that the likelihood of program conflicts or unexpected stoppages are minimized. Although this certainly helps reduce the number of problems and improves the forward and backward compatibility with previous OS versions, on occasion even the OS vendors experience or cause difficulties when they do not conform to the parameters established for their own programs.

39.4.1.10 Other Sources of Error. It is not unusual for errors to occur where programs send bad data to devices, ignore error codes coming back, and even try to use devices that are busy or missing. The hardware might well be broken, but the software also is considered to be in error when it does not recover from such hardware conditions.

Additional errors can occur through improper builds of the executable; for example, if an old version of a module is linked to the latest version of the rest of the program, the wrong sign-on screens may pop up, the wrong copyright messages may be displayed, the wrong version numbers may appear, and various other inaccuracies may occur.

39.4.1.11 Regulatory Compliance Considerations. If an organization is subject to Sarbanes-Oxley (SOX), then documentation of the errors encountered and the resulting internal reporting and remediation efforts is critical. The error should clearly indicate how it was identified, who identified it, how it was reported to management, what the remediation will be, and when it is anticipated to be completed. Without these details, management may encounter significant difficulties in confirming that adequate internal control mechanisms exist and that it is informed and involved appropriately and adequately. Also, an error could result in a control weakness being identified by the external auditors or, even worse, as a material weakness, depending on its nature and severity.

Additionally, if and when errors are noted that affect multilocation systems or applications, the significance and materiality must be considered from the aspect of both the subsidiary and the headquarters locales. Since laws differ among states and countries, what might be considered legally acceptable standards of privacy or accounting in one locale might not be acceptable elsewhere.

39.4.2 User Interface. Generally speaking, the term “user interface” denotes all aspects of a system that are relevant to a user. It can be broadly described as the user virtual machine (UVM). This would include all screens, the mouse and keyboard, printed outputs, and all other elements with which the user interacts. A major problem arises when system designers cannot put themselves in the user’s place and cannot foresee the problems that a technologically challenged user will have with an interface designed by a technologically knowledgeable person.

Documentation is a crucial part of every system. Each phase of development—requirements, analysis, development, coding, testing, errors, error solutions and modifications, implementation, and maintenance—needs to be documented. All documents and their various versions need to be retained for both future reference and auditing purposes. Additionally, it is important to document the correct use of the system and

TYPES OF SOFTWARE ERRORS 39 · 11

provide adequate instructional and reference materials to the user. Security policies and related enforcement and penalties also need to be documented. Ideally, the documentation should enable any technically qualified person to repair or modify any element, as long as the system remains operational.

39.4.2.1 Functionality. A program has a functionality error if performance that can reasonably be expected is confusing, awkward, difficult, or impossible. Functionality errors often involve key features or functions that have never been implemented. Additional functionality errors exist when:

- Features are not documented.
- Required information is missing.
- A program fails to acknowledge legitimate input.
- There are factual errors or conflicting names for features.
- There is information overload.
- The material is written to an inappropriate reading level.
- The cursor disappears or is in the wrong place.
- Screen displays are wrong.
- Instructions are obscured.
- Identical functions require different operations in different screens.
- Improperly formatted input screens exist.
- Passwords or other confidential information are not obscured or protected adequately.
- Tracing the user data entry or changes is unavailable or incomplete.
- Segregation of duties is not enforced. (This can be particularly critical for organizations subject to the Health Insurance Portability and Accountability Act [HIPAA], SOX, International Organization for Standardization [ISO] 17799, and/or the Gramm-Leach-Bliley Act [GLBA].)

39.4.2.2 Control (Command) Structure. Control structure errors can cause serious problems because they can result in:

- Users getting lost in a program
- Users wasting time because they must deal with confusing commands
- Loss of data or the unwanted exposure of data
- Work delay
- Financial cost
- Unanticipated exposure to data leakage or compromise; this can result in significant liability if consumers' personal identifying information (PII) is compromised
- Data not being encrypted as intended or being visible to unauthorized users

Some common errors include:

- Inability to move between menus
- Confusing and repetitive menus

39 · 12 SOFTWARE DEVELOPMENT AND QUALITY ASSURANCE

- Failure to allow adequate command-line entries
- Requiring command-line entries that are neither intuitive nor clearly defined on screen
- Failure of the application program to follow the operating system's conventions
- Failure to distinguish between source and parameter files, resulting in the wrong values being made available to the user through the interface, or failure to identify the source of the error
- Inappropriate use of the keyboard when new programs do not meet the standard of a keyboard that has labeled function keys tied to standard meanings
- Missing commands from the code and screens, resulting in the user being unable to access information, to utilize programs, or to provide for the system to be backed up and recoverable, as well as a host of other commands that leave the system in a state of less-than-optimum operability
- Inadequate privacy or security that can result in confidential information being divulged, the complete change or loss of data without recoverability, poor reporting, and even undesired access by outside parties

39.4.2.3 Performance. Speed is important in interactive software. If a user feels that the program is working slowly, that can be an immediate problem. Slow operation can depend on (but is not limited to) the OS, the other applications running, memory allocation, memory leakage, and program conflicts. At another level, performance suffers when program designs make it difficult to change their functionality in response to changing requirements. Performance errors include slow response, unannounced case sensitivity, uncontrollable and excessively frequent automatic saves, inability to save, and limited scrolling speed.

39.4.2.4 Output Format. Output format errors can be frustrating and time consuming. An error is considered to have occurred when the user cannot change fonts, underlining, boldface, and spacing that influence the final look of the output; alternatively, delays or errors when printing or saving document may occur. Errors occur when the user cannot control the content, scaling, and look of tables, figures, and graphs. Additionally, there are output errors that involve expression of the data to an inappropriate level of precision.

39.5 DESIGNING SOFTWARE TEST CASES

39.5.1 Good Tests. No software program can ever be tested completely, since it would not be cost effective or efficient to test the validity, parameters, syntax, and boundaries of every line of code. Even if all valid inputs are defined and tested, there is no way to test all invalid inputs and all the variations on input timing. It is also difficult, if not impossible, to test every path the program might take, find every design error, and prove programs to be logically correct. Nevertheless, a good test procedure will find most of the problems that would occur, allowing the designers and developers to correct those problems and ensure that the software works properly. Generally, OS and application manufacturers classify the severity level of errors encountered, and many (as evidenced by the frequency and content of program patches) correct only the most serious ones that either they or the users notice.

DESIGNING SOFTWARE TEST CASES 39 · 13

Over the past several years, there has been an increasing debate on whether software vulnerabilities should be publicized immediately by researchers, or if the researchers should automatically notify and give the manufacturer time to correct the flaw. Generally, from a casual review of the vulnerabilities discovered concerning Microsoft's Internet Explorer version 6 over the past few years, the trend appears to have been that the researcher notifies Microsoft and allows Microsoft a varying period of time to create and distribute the patch (e.g., three to six months) before the researcher publicly announces the weakness. This trend is based on published reports in the technical news media.

Whenever you expect the same results from two tests, you consider the tests equivalent. A group of tests forms an equivalence class if the tester believes that the tests all test the same thing, and if one test catches or does not catch a bug, the others probably will do the same. Classical boundary tests check a program's response to input and output data; equivalence tests, however, teach a way of thinking about analyzing programs that enhances and strengthens test planning.

Finding equivalence classes is a subjective process. Different people analyzing the same program will come up with different lists of equivalence classes because of what the programs appear to achieve. Test cases often are lumped into the same equivalence class when they involve the same input variables, result in similar operations, affect the same output variables, or handle errors in the same manner.

Equivalence classes can be groups of tests dealing with ranges or multiple ranges of numbers, members of a group, and even time determined. Equivalence classes are generally the most extreme values, such as the biggest, smallest, fastest, and slowest. When testing, it is important to test each edge of an equivalence class, on all sides of each edge. Testers should use only one or two test cases from each equivalence class because a program that passes the tests generally will pass any test drawn from that class. Invalid input equivalence classes often allow program bugs to be overlooked during debugging.

39.5.2 Emphasize Boundary Conditions. Boundaries are crucial for checking each program's response to input and output data. Equivalence class boundaries are the most extreme values of the class; for example, boundary conditions may consist of the biggest and smallest, soonest and latest, shortest and longest, or slowest and fastest members of an equivalence class.

Tests should include values below, at, and above boundary values. Additionally, they should be tested at various user levels (e.g., administrator, root or super user, data entry, and read-only access). This testing must ensure that there are no conditions at any level that permit unintended access or unauthorized escalation of privileges.

When programs fail with nonboundary values, they generally fail at the boundaries too. Programs passing these tests probably also will pass any other test drawn from that class. Tests should also be able to generate the largest and smallest legitimate output values, remembering that input-boundary values may not generate output-boundary values.

Today, the most prevalent security breaches involve buffer overflows in active code (ActiveX and Java), in which inadequate bounds checking on input strings allows overflowing text to be interpreted as code that then can carry out improper operations. Restrictions on input length and type would prevent these exploits. Such overflows can result in the exposure of consumers' PII and incurring unwanted liability for the exposure; currently, the liability exists principally at the state level and may include criminal sanctions, depending on the state in which the consumers reside.

39 · 14 SOFTWARE DEVELOPMENT AND QUALITY ASSURANCE

39.5.3 Check All State Transitions. All interactive programs move from one state to another. A program's state is changed whenever something causes the program to alter output (e.g., to display something different on the screen) or to change the range of available choices (e.g., displaying a new menu of user options).

To test state transitions, the test designer should lay out a transition probability matrix to show all the paths people are likely to follow. State transitions can be complex and often depend not on the simple choices provided but on the numbers the user enters. Testers often find it useful to construct menu maps that show exactly where to go from each choice available.

Menu maps also can show when and where users go when menu or keyboard commands take them to different states or dialogs. Maps are particularly handy when working with spaghetti code (code that has been poorly designed or badly maintained so that logical relationships are difficult to see); maps allow the designer or user to reach a dialog box in several ways and then proceed from the dialog box to several places. For spaghetti code, the menu maps afford a simpler method of spotting relationships between states than trying to work exclusively from the program itself because the map shows transition between states on paper or on screen. After mapping the relationships, the designer or user can check the program against the map for correctness.

Full testing theoretically requires that all possible paths are tested. However, in practice, complete testing may be unattainable for complex programs. Therefore, it makes sense to test the most frequently used paths first.

39.5.3.1 Test Every Limit. It is necessary to test every limit on a program's behavior that is specified by any of the program's documents. Limits include the size of files, the number of terminals, the memory size the program can manage, the maximum size it requires, and the number of printers the program can drive. It is important to check on the ability of the program to handle large numbers of open files on an immediate basis and on a long-term, continuing basis as well. Load testing is actually boundary condition testing and should include running tests the program ought to pass and tests the program should not pass.

For Web applications, tools exist to simplify this task, as some will check the boundaries and limits of fields at varying depths of data entry to test the protections against buffer overflow, invalid data, and incorrect data syntax; such tools include but are not limited to WebInspect by Spi Dynamics.

39.5.3.2 Test for Race Conditions. After testing the system under "normal" load, it should be tested for race conditions. A "race condition" is usually defined as anomalous behavior due to unexpected critical dependence on the relative timing of events. Race conditions generally involve one or more processes accessing a shared resource such as a file or variable, where this multiple access has not been properly controlled. Systems that are vulnerable to races, especially multiuser systems with concurrent access to resources, should undergo a full cycle of testing under load. Race conditions sometimes can be identified through testing under heavy load, light load, fast speed, slow speed, multiprocessors running concurrent programs, enhanced and more numerous input/output devices, frequent interrupts, less memory, slower memory, and related variables.

39.5.4 Use Test-Coverage Monitors. For complex programs, path testing usually cannot or would not test every possible path throughout a program for practical

DESIGNING SOFTWARE TEST CASES 39 · 15

as well as cost-effectiveness reasons. A more practical glass box approach (i.e., with knowledge of the internal design of a program) is to use the source code listing to force the program down every branch visible in the code. When programmers add special debugging code during development, a unique message will print out or be added to a log file whenever the program reaches a specified point. Typically, such messages can be generated by calling a print routine with a unique parameter for each block of code. When source code is compiled, many languages permit a switch to be set to allow conditional compilation, so that a test version of the code contains active debugging statements whereas a production version does not. For interpreted code, such as most fourth-generation languages, similar switches allow for inclusion or activation of debugging instructions.

Since distinct messages are planted, it is possible to know exactly what point in the program the test has reached. Programmers specifically insert these messages at significant parts of the program. Once these messages have been added, anyone can run the program and conclude whether the different parts actually have run and been tested.

Special devices or tools also can be used to add these messages to the code automatically. Once source code is fed to such a coverage monitor, it analyzes the control structures in the code and adds probes for each branch of the program. Adding these probes or lines of code is called instrumenting the program. It is possible to tell that the program has been pushed down every branch when all the probes have been printed. Besides instrumenting code, a good coverage monitor can capture probe outputs, perform analyses, and summarize them. Some coverage monitors also log the time used for each routine, thus supporting performance analysis and optimization.

A coverage monitor is designed for glass box testing, but knowledge of the internals of the program under test is not needed in order to use the monitor. The coverage monitor counts the number of probe messages, reports on the number of probes triggered, and even reports on the thoroughness of the testing. Because it is possible for the coverage monitor to report on untriggered branches, the monitor can find and identify code that does not belong, such as routines included by default but never used or deliberately inserted undocumented code (a Trojan horse). Some commercially available coverage monitors are, unfortunately, themselves full of bugs. It is important to obtain a full list of the known bugs and patches from the developer. In addition, such tools should themselves be tested with programs that contain known errors, in the process known as seeding, to verify their correctness.

39.5.5 Seeding. Quality assurance procedures themselves benefit from testing. One productive method, known as seeding, is to add known bugs to a program and measure how many of them are discovered through normal testing procedures. The success rate in identifying such known bugs can help estimate the proportion of unknown bugs left in a program. Such seeding is particularly important when establishing automated testing procedures and test-coverage monitors; also, it is useful for SOX or other compliance-related testing.

39.5.6 Building Test Data Sets. One of the most serious errors in quality assurance is to allow programmers to use production data sets for testing. Production data may include confidential information that should not be accessible to programming staff, so access should be forbidden not only to production data but even to copies of production data, or even to copies of subsets or samples from production data. However, *anonymizing* processes applied to copies of production data (e.g., scrambling names

39 · 16 SOFTWARE DEVELOPMENT AND QUALITY ASSURANCE

and addresses of patients in a medical records database) may produce test data sets suitable for use in quality assurance.

Alternatively, if the organization has a test system that is completely (logically and physically) segregated from the production system, using historical data may be appropriate. In doing so, particular care and appropriate documentation must be maintained to ensure that there is no chance of the test system interfacing with the production system. A side benefit of such testing is that it may help the organization's SOX or HIPAA periodic testing and documentation regarding the availability, utility, and integrity of backup media and systems.

39.5.7 Fuzzing. A useful addition to the techniques described above is the automated injection of random data into the input stream. *Fuzzing* can support the planned challenges to a system with unexpected inputs that programmers may not have thought of.

39.6 BEFORE GOING INTO PRODUCTION

39.6.1 Regression Testing. Regression testing is fundamental work done by glass box and black box testers. The term is used in two different ways. The first definition involves those tests where an error is found and fixed and the test that exposed the problem is performed again. The second definition involves finding and fixing an error and then performing a standard series of tests to make certain the changes or fix made did not disturb anything else.

The first type of regression testing serves to test that a fix does what it is intended to do. The second type tests the fix and also tests the overall integrity of the program. It is not unusual for people who mention regression testing to be referring to both definitions, since both involve fixing and then retesting. It is recommended that both types of testing be done whenever errors are fixed. If an organization is subject to SOX or HIPAA, appropriate documentation should be maintained that details the test methodology, the sample selection, the frequency, the results, and the remediation (if any).

39.6.2 Automated Testing. In some enterprises, every time a bug is fixed, every time a program is modified, and every time a new version is produced, a tester runs a regression test. All of this testing takes a considerable amount of time and consumes both personnel and machine resources. In addition, repetitive work can become mind-numbing, so testers may accidentally omit test cases or overlook erroneous results.

To cut down on the time consumption of personnel and the repetitive nature of the task, it is possible to program the computer to run acceptance and regression tests. This type of automation results in execution of the tests, collection of the results, comparison of the results with known good results, and a report of the results to the tester.

Early test automation consisted essentially of keyboard macros or scripts, with limited capacity for responding to errors; typically, an error would produce invalid results that would be fed into the next step of processing and lead to long chains of meaningless results. Slightly more advanced test harnesses would halt the test process at each error and require human intervention to resume. Both methods were of only modest help to the test team. However, today's test-automation software can include test databases that allow orderly configuration of tests, specific instructions for restarting test sequences after errors, and full documentation of inputs and results. For realistic load testing, some systems can be configured to simulate users on workstations, with

MANAGING CHANGE 39 · 17

scripts that define a range of randomly generated values, specific parameter limits, and even variable response times. Large-scale load testing can connect computers together through networks to simulate thousands of concurrent users and thus speed identification of resource exhaustion or race conditions.

Test automation can be well worth the expenditures required. Automated testing is usually more precise, more complete, faster, and less expensive than the tests done by human personnel. Typically, automated tests can accomplish tenfold or hundredfold increases in the number of tests achievable through manual methods. In addition to freeing personnel to do more rewarding work, such methods greatly reduce overall maintenance costs due to detection of errors before systems reach production. Additionally, for regulatory purposes, the automated testing requires significantly less sampling, documentation, and testing. Appropriate management review and reporting of results still are required, as is documentation of remediation follow-up.

39.6.3 Tracking Bugs from Discovery to Removal. Finding a bug is not enough. Even eliminating a bug is insufficient. A system also must allow the causes of each bug to be identified, documented, and rectified. For these reasons, it is important to track and document all problems from the time of their discovery to their removal, to be certain that the problems have been resolved and will not affect the system. Too, this provides important historical records in diagnosing and remediating incidents and provides appropriate documentation for audit and regulatory purposes.

Problem-tracking systems must be used to report bugs, track solutions, and write summary reports about them. An organized system is essential to ensure accountability and communication regarding the bugs. Typical reports include: where bugs originate (e.g., which programmers and which teams are responsible for the greatest number of bugs); types of problems encountered (e.g., typographical errors, logic errors, boundary violations); and time to repair. However, problem-tracking systems can raise political issues, such as project accountability, personal monitoring, control issues, and issue remediation regarding the data in the database and who owns them.

Once a tracking system is established, a bug report is entered into the database system and a copy goes to the project manager, who either prioritizes it and passes it along or responds to it personally. Eventually, the programmers will be brought into the loop, will fix the problem, and will mark the problem as fixed. The fixed problem then is tested and a status report is issued. If a fix proves not to work, that becomes a new problem that needs to be addressed.

39.7 MANAGING CHANGE. Many people may be involved in creating and managing systems; whenever changes are made to the system, those changes need to be managed and monitored in an organized fashion to avoid chaos.

39.7.1 Change Request. The change request is an important document that requests either a fix or some other change to a program or system. Information contained on the change request form generally includes who is requesting the change, the date the request is being made, the program and area affected, the date by which the change is needed, and authorization to go ahead and make the change.

Most companies have paper forms that are retained for monitoring and auditing purposes. As the popularity of digital signatures grows and society continues to move toward a paperless office, it seems logical that some of these change requests eventually will become totally automated.

39 · 18 SOFTWARE DEVELOPMENT AND QUALITY ASSURANCE

39.7.2 Tracking System. The tracking system may be manual, automated, or a combination of methods; for regulatory and control purposes, a system that incorporates automated tracking is best, since it reduces the chance for error, and it can enhance follow-up and remediation. Regardless of how it is kept, the original change form generally is logged into a system, either manually filed or entered into an automated database by some identifying aspect, such as date, type of change, or requesting department.

The system is used to track what happens to the change. When an action is taken, information must be entered into the system to show who worked on the request, what was done, when action was taken, what the result of the action was, who was notified of the actions taken, whether the change was accepted, and related information. This information is crucial in ensuring that issues are managed appropriately, escalated as needed to senior management, and tracked until remediated.

39.7.3 Regression Testing. When a change or fix has been made, regression testing verifies that the change has indeed been made and that the program now works in the fashion desired, including all other functions. The successful completion of the testing should be documented in writing; this can be through the tracking system or in documentation appended to the tracking system entry.

39.7.4 Documentation. Documentation is crucial when considering, approving, and implementing changes. If information is retained strictly in an individual's head, what happens when the individual goes on vacation, is out sick, or leaves the enterprise permanently? What if the individual simply does not remember what changes were made, or when; how does an organization know who was involved, what actually was done, and who signed off that the changes were made and accepted?

Undocumented changes can result in:

- System crashes
- Inconsistent data entry
- Inappropriate segregation of duties
- Data theft or corruption
- Embezzlement
- Other serious crimes

Lack of documentation can mean that unauthorized changes were made and likely can result in regulatory and audit violations. Undocumented changes may violate segregation of duties (e.g., a person might be informally approving changes to his or her own account, approving changes carried out by a supervisor, or authorizing changes by a person who is not under the responsibility of the approver).

Lack of documentation is often a violation of corporate policy and is often cited in audits; it often will cause significant difficulties regarding SOX, GLBA, and HIPAA compliance efforts and reviews, since management's assertions might not be supported adequately. Because constructing adequate documentation after the fact is almost impossible (and frequently can be considered to be falsification of records in sectors such as financial services), documentation must proceed in step with every phase of a system's development and operation. Thereafter, documentation must be retained according to company policy and legal requirements.

SOURCES OF BUGS AND PROBLEMS 39 · 19

39.8 SOURCES OF BUGS AND PROBLEMS. Locating and fixing bugs is an important task, but it is also essential to try to determine where and why the bugs and related problems originated. By finding and studying the source, often it is possible to prevent a recurrence of the same or similar type of problem.

39.8.1 Design Flaws. Design flaws often occur because of poor communication between users and designers. Users are not always clear about what they want and need, while designers misunderstand, misinterpret, or simply ignore what the users relate to them. Within the design process, even when users' feelings and requirements are known and understood, design flaws can occur. However, they are easier to identify and remedy if the appropriate documentation is affected. Without the proper documentation, it may not be possible to identify the source of an error; this can result in patches or remedies that need to be patched. Often flaws result from attempts to comply with unrealistic delivery schedules. Managers should support their staff in resisting the pressure to rush through any of the design, development, and implementation stages.

39.8.2 Implementation Flaws. Whenever a program is developed and implemented, there are time limits and deliverable dates. Most problems during development cause delay, so developers and testers often are rushed to get the job done. Sometimes they sacrifice the documentation, the review portion of the project, or even cut short testing, leaving unrecognized design and implementation flaws in place. Managers should emphasize the value of thorough review and testing, and allocate enough time to avoid such blunders.

39.8.3 Unauthorized Changes to Production Code. If problems are traced to unauthorized changes, project managers should examine their policies. If the policies are clear, perhaps employee training and awareness programs need improvement. Managers also should examine their own behavior to ensure that they are not putting such pressure on their staff that cutting corners is perceived to be acceptable.

39.8.4 Insufficient or Substandard Programming Quality. A programmer can make or break a program and a project. Programmers play essential roles in software development. Therefore, it is important that all programmers on a project be capable and reliable. Project programmers should be carefully screened before being placed on a software development project to ensure that their skills meet project requirements. Sometimes holes in programmers' skills can be filled with appropriate training and coaching. However, if problems appear to be consistent, management may want to check a programmer's background to verify that he or she did not falsify information when applying for the job. Truly incompetent programmers need to be removed from the project. Additionally, if management suspects that the programming quality is inadequate, then it should consider having an independent programmer review the code, performing particularly rigorous quality assurance (QA) testing.

39.8.5 Data Corruption. Data corruption can occur because of poor programming, invalid data entry, inadequate locking during concurrent data access and modification, illegal access by one process to another process data stack, and hardware failures. Data corruption can occur even when a program is automatically tested or run without human intervention. In any event, when searching for the sources of bugs and other problems, it is important to do a careful review of the data after each round of testing, in order to identify deviations from the correct end state. The review should

39 · 20 SOFTWARE DEVELOPMENT AND QUALITY ASSURANCE

be documented appropriately and escalated to management; it then will assist the regulatory compliance efforts.

39.8.6 Hacking. When bugs and problems do occur, hacking—both internal and external—should be considered a possibility and searched for by reviewing the logs of who worked on the software and when that work was done. Managers should be able to spot the use of legitimate IDs when the individuals involved were on vacation, out sick, or not available to log on for other reasons.

Archiving logs and retrieving them is an increasingly critical capability, aside from the regulatory requirements of SOX, HIPAA, GLBA, and so on. Logs are often used for:

- Incident response activities, particularly in determining:
 - If an incident occurred
 - When an incident occurred
 - What happened prior to and subsequent to when an incident is suspected to have occurred
- Research; this can be to review performance capabilities or issues of hardware or software or networking

To ensure that the logs are available when needed, the organization needs a sound strategy that preserves and archives logs periodically, so that an attacker or a system anomaly does not wipe them from the system's hard drive. Although there are many means to accomplish this, a frequently used inexpensive method involves establishing a log server. For example, an organization might identify a discarded desktop computer, outfit it with a new, high-speed, high-capacity optical read-only memory (ROM) drive, and increase memory capacity (if needed). By routing the logs to the optical drive frequently and recording them throughout the day (perhaps multiple times during an hour) and changing the optical disk at least daily, the organization will create a forensically sound archive that can be duplicated easily for research, law enforcement, or legal use. Often the daily logs volume will not fill a DVD or a CD. However, having the regular routine of changing the media daily helps ensure that the media does not reach its capacity at an inopportune time—indeed, if a worm such as Blaster or Nimda infects the organization, then the logs likely will swell rapidly. Too, if an attacker within the system stops the logging, the organization should encounter less difficulty determining the date, time (or time period), and extent of the intrusion because ROM cannot be overwritten, in contrast with read-write (RW) media.

Unauthorized changes to code and data sometimes can be identified even if the perpetrator stops the logging and deletes or modifies log files. One method is to create checksums for production or other official versions of software and to protect the checksums against unauthorized access and modification using encryption and digital signatures. Similarly, unauthorized data modifications sometimes can be made more difficult by creating checksums for records and linking the checksums to time and data stamps and to the checksums for authorized programs. Under those conditions, unauthorized personnel and intruders find it difficult to create valid checksums for modified records. Other products, such as intrusion detection systems (IDSs) or intrusion prevention systems (IPSs), may be useful in an organization's protection and risk-management strategy.

FURTHER READING 39 · 21

39.9 CONCLUSION. This chapter has presented a comprehensive overview of the software development and quality assurance processes that must be utilized when developing, implementing, and modifying software. Software development involves more than simply selecting and utilizing an approach such as the traditional, waterfall, or RAD methodology. It means working as a team to develop, review, refine, and implement a viable working product. Many good techniques and products can be applied to both the SDLC and the quality assurance portions of producing software; some of the key elements are good documentation, allowing sufficient time for testing in the development and maintenance processes, building good tests, establishing test data, automating testing, and keeping track of change requests.

39.10 FURTHER READING

- Basandra, S. *Software Architecture, Data Structures, Algorithms, Programming and Testing Questions and Answers*. Basandra Books, 2013.
- Black, D. B. *Software Quality Assurance (Building Better Software Better)*. Amazon Digital Publishing, 2012.
- Chemuturi, M. *Mastering Software Quality Assurance: Best Practices, Tools and Techniques for Software Developers*. J. Ross Publishing, 2010.
- Fox, C., and P. Zonneveld. *IT Control Objectives for Sarbanes-Oxley: The Role of IT in the Design and Implementation of Internal Control Over Financial Reporting*, 2nd ed. IT Governance Institute, 2006.
- Godbole, N. S. *Software Quality Assurance: Principles and Practices*, 2nd ed. Alpha Science International Ltd., 2013.
- Hambling, B., and P van Goethem. *User Acceptance Testing—A Step-by-Step Guide*. BCS, The Chartered Institute for IT, 2013.
- Horch, J. W. *Software Quality Assurance: A Standards-based Guide*. IEEE Computer Society, 2013.
- ISO/IEC. “ISO/IEC 29119 Software Testing: The new international software testing standard.” ISO/IEC, 2013. www.softwaretestingstandard.org/index.php
- McCaffrey, J. D. *Software Testing: Fundamental Principles and Essential Knowledge*. BookSurge Publishing, 2009.
- OWASP. “Fuzzing.” Open Web Application Security Project, 2013. <https://www.owasp.org/index.php/Fuzzing>
- Rasmussen, J. *The Agile Samurai: How Agile Masters Deliver Great Software*. Pragmatic Bookshelf, 2013.
- Rubin, K. S. *Essential Scrum: A Practical Guide to the Most Popular Agile Process*. Addison-Wesley Professional, 2012.
- Shaffer, S. C. *A Brief Introduction to Software Development and Quality Assurance Management*. S. C. Shaffer Publications, 2013.
- Tassey, G. *The Economic Impacts of Inadequate Infrastructure for Software Testing: Final Report*. RTI Project Number 7007.011. NIST, 2002.
- Wells, D. *Extreme Programming: A Gentle Introduction*. XP, 2009, www.extremeprogramming.org
- Wood, D. C. *Principles of Quality Costs: Financial Measures for Strategic Implementation of Quality Management*, 4th ed. Quality Press, 2013.

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 40

MANAGING SOFTWARE PATCHES AND VULNERABILITIES

Karen Scarfone, Peter Mell, and Murugiah Souppaya

40.1 INTRODUCTION	40·1	40.4 ENTERPRISE PATCH MANAGEMENT TECHNOLOGIES	40·7
40.2 THE IMPORTANCE OF PATCH MANAGEMENT	40·2	40.4.1 Components and Architecture	40·7
40.3 THE CHALLENGES OF PATCH MANAGEMENT	40·2	40.4.2 Security Capabilities	40·8
40.3.1 Timing, Prioritization, and Testing	40·3	40.4.3 Management Capabilities	40·9
40.3.2 Patch Management Configuration	40·4		
40.3.3 Alternative Host Architectures	40·5		
40.3.4 Other Challenges	40·6		
		40.5 METRICS AND MEASURES	40·10
		40.6 FURTHER READING	40·11
		40.7 NOTES	40·11

40.1 INTRODUCTION. *Vulnerabilities* are flaws that can be exploited by a malicious entity to gain greater access or privileges than it is authorized to have on a computer system. *Patches* are additional pieces of code developed to address problems (commonly called “bugs”) in software. Patches correct security and functionality problems in software and firmware. *Patch management* is the process for identifying, acquiring, installing, and verifying patches for products and systems. Patch management is a security practice designed to proactively prevent the exploitation of IT vulnerabilities that exist within an organization. The expected result is to reduce the time and money spent dealing with vulnerabilities and the exploitation of those vulnerabilities. Proactively managing vulnerabilities of systems will reduce or eliminate the potential for exploitation and involve considerably less time and effort than responding after an exploitation has occurred.

This chapter is designed to assist organizations in implementing patch and vulnerability remediation programs. It first explains the importance of patch management and then discusses the challenges inherent in performing patch management. Next, the chapter provides an overview of enterprise patch management technologies and gives recommendations for their use. It also seeks to inform the reader about possible measures and metrics for enterprise patch management.¹

40 · 2 MANAGING SOFTWARE PATCHES AND VULNERABILITIES

For more general information about related topics, see selected chapters in this *Handbook*:

- Chapter 38: Writing Secure Code
- Chapter 39: Software Development and Quality Assurance
- Chapter 46: Vulnerability Assessment
- Chapter 47: Operations Security and Production Controls
- Chapter 52: Application Controls
- Chapter 53: Monitoring and Control Systems

40.2 THE IMPORTANCE OF PATCH MANAGEMENT. From a security perspective, patches are most often of interest because they are mitigating software flaw vulnerabilities; applying patches to eliminate these vulnerabilities significantly reduces the opportunities for exploitation. Also, patches are usually the most effective way to mitigate software flaw vulnerabilities, and are often the only fully effective solution. Sometimes there are alternatives to patches, such as temporary workarounds involving software or security control reconfiguration, but these workarounds often negatively impact functionality.

Patches serve other purposes than just fixing software flaws; they can also add new features to software and firmware, including security capabilities. New features can also be added through upgrades, which bring software or firmware to a newer version in a much broader change than just applying a patch. Upgrades may also fix security and functionality problems in previous versions of software and firmware. Also, vendors often stop supporting older versions of their products, which includes no longer releasing patches to address new vulnerabilities, thus making older versions less secure over time. Upgrades are necessary to get such products to a supported version that is patched.

As Section 40.3 of this chapter explains, there are several challenges that complicate patch management. If organizations do not overcome these challenges, they will be unable to patch systems effectively and efficiently, leading to easily preventable compromises. Organizations that can minimize the time they spend dealing with patching can use those resources for addressing other security concerns. Already many organizations have largely operationalized their patch management, making it more of a core IT function than a part of security. However, it is still important for all organizations to carefully consider patch management in the context of security because patch management is so important to achieving and maintaining sound security.

Patch management is required by various security compliance frameworks, mandates, and other policies. For example, NIST Special Publication (SP) 800-53² requires the SI-2, Flaw Remediation security control, which includes installing security-relevant software and firmware patches, testing patches before installing them, and incorporating patches into the organization's configuration management processes. Another example is the Payment Card Industry (PCI) Data Security Standard (DSS),³ which requires that the latest patches be installed and sets a maximum time frame for installing the most critical patches.

40.3 THE CHALLENGES OF PATCH MANAGEMENT. This section briefly examines the challenges inherent in performing patch management. These are the challenges that the patch management technologies discussed in Section 40.4 of this chapter are trying to solve.

THE CHALLENGES OF PATCH MANAGEMENT 40 · 3

40.3.1 Timing, Prioritization, and Testing. Timing, prioritization, and testing are intertwined issues for enterprise patch management. Ideally, an organization would deploy every new patch immediately to minimize the time that systems are vulnerable. However, in reality this is simply not possible because organizations have limited resources, which makes it necessary to prioritize which patches should be installed before other patches. Further complicating this is the significant risk of installing patches without first testing them, which could cause serious operational disruptions, potentially even more damaging than the corresponding security impact of not pushing the patches out. Unfortunately, testing patches consumes even more of the limited resources and makes prioritization even more important. For patch management, timing, prioritization, and testing are often in conflict.

Product vendors have responded to this conflict by bundling patches for their products. Instead of releasing dozens of patches one at a time over a period of three months, necessitating testing and patch deployment every few days, a vendor might release their patches in a single bundle once a quarter. This allows an organization to perform testing once and roll out patches once, which is far more efficient than testing and rolling out all the patches separately. It also reduces the need to prioritize patches—the organization just needs to prioritize the bundle instead of separately prioritizing each patch it contains. Vendors who bundle patches tend to release them monthly or quarterly, except for cases when an unpatched vulnerability is actively being exploited, in which case they usually issue the appropriate patch immediately instead of delaying it for the next bundle.

There is a downside to patch bundling; it lengthens the time from the discovery of a vulnerability to the time a patch for it becomes publicly available. If an attacker discovers the same vulnerability before the patch is released, the attacker may have a longer window of opportunity to exploit the vulnerability because of the intentional delay in releasing the patch. However, there are two mitigating factors here. One is that if exploitation is known to be occurring, the vendor is likely to release the patch immediately. The other factor is that patches may be installed more quickly if they are bundled than if they are all released separately. So that effectively helps to shrink the window of opportunity for vulnerabilities associated with bundled patches.

There are even more issues to consider with timing. The release of a patch may provide attackers with the information that they need to exploit the corresponding vulnerability (e.g., reverse engineer the vulnerability from the patch), meaning that a newly released patch might need to be applied immediately to avoid compromises. However, if a vulnerability is not being exploited yet, organizations should carefully weigh the security risks of not patching with the operational risks of patching without performing thorough testing first. In some operational environments, such as virtual hosts with snapshot capabilities enabled, it may be preferable to patch without testing as long as the organization is fully prepared to roll back the patches if there are usability or functionality problems caused by them.

Another fundamental issue with timing is forcing the implementation of changes to make a patch take effect; this can require restarting a patched application or service, rebooting the operating system, or making other changes to the state of the host. Ultimately what matters is not when the patch was installed, but when the patch actually takes effect. In some cases it may make more sense to mitigate a vulnerability through an alternative method, at least until patches are fully operational. An example is changing configuration settings for vulnerable software to temporarily block vulnerable application functionality. Each mitigation option has different implications for the security, functionality, and operations of the vulnerable host, so it is not a trivial

40 · 4 MANAGING SOFTWARE PATCHES AND VULNERABILITIES

matter to select one option over others. Also, if configuration settings are changed, this necessitates preserving the old setting values and restoring them at the appropriate time. Another problem with changing configuration settings is that they often require a state change to the host to take effect, such as restarting an application. Implementing configuration changes may be as disruptive to the operations of a host as installing a patch.

Forcing the implementation of changes that requires rebooting the operating system can be problematic when the host requires authentication before booting, such as the use of full disk encryption (FDE) software. Organizations using FDE software or other technologies that require authentication before booting should carefully consider the impact that these technologies may have on patch installation.

Prioritizing which patches to apply and when to apply them is closely related to timing, but there are other considerations as well. It can depend on the relative importance of the vulnerable systems (for example, servers versus clients) and the relative severity of each vulnerability (e.g., vulnerability severity metrics such as the Common Vulnerability Scoring System [CVSS]). Another consideration is dependencies that patches may have on each other; installing one patch may require installing other patches first, and in some cases restarting an application or rebooting a host multiple times to make the patches take effect sequentially.

In summary, organizations should carefully consider the relevant issues related to timing, prioritization, and testing when planning and executing their enterprise patch management processes.

40.3.2 Patch Management Configuration. Another major challenge in enterprise patch management is that there are usually multiple mechanisms for applying patches. For example:

- A piece of software may be able to automatically update itself.
- A centralized OS management tool may be able to initiate patching.
- Third-party patch management applications may be able to initiate patching.
- Network access control, health check technologies, and similar technologies may be able to initiate patching.
- A user may be able to manually direct software to update itself.
- A user may be able to manually install a patch or a new version of the software.

Having multiple ways of applying patches can cause conflicts. Multiple methods might each try to patch the same software, which is particularly problematic when the organization doesn't want certain patches applied because of issues with those patches, testing delays, etc. Multiple methods can also cause patches to be delayed or missed, because each tool or administrator may assume another one is already taking care of a particular patch. Organizations should identify all the ways in which patches could be applied and act to resolve any conflicts among patch application methods.

A related problem with patch management configuration is that users may override or circumvent patch management processes. If users are able to make changes to their hosts' software, such as altering settings (e.g., enabling direct updates, disabling patch

THE CHALLENGES OF PATCH MANAGEMENT 40 · 5

management software), installing old versions of software, and uninstalling patches, they can undermine patch management integrity. To address these problems, organizations should ensure that users cannot disable or otherwise negatively affect enterprise patch management technologies, and organizations should perform continuous monitoring of enterprise patch management technologies to identify any issues that occur.

40.3.3 Alternative Host Architectures. Enterprise patch management is relatively straightforward when all of the hosts are fully managed and running typical applications and operating systems on a regular platform. When alternative host architectures are employed, patch management can be considerably more challenging. Examples of these architectures include the following:

Unmanaged hosts. As discussed in Section 40.3.2 of this chapter, it can be much more difficult to control patching when hosts are not centrally managed (i.e., users manage their own hosts).

Out-of-office hosts (e.g., telework laptops). Hosts on other networks are not protected by the enterprise’s network security controls (firewalls, network intrusion detection systems, vulnerability scanners, etc.).

Nonstandard IT components (e.g., appliances). On such hosts, it’s often not possible to patch individual applications independently. Rather, the organization must wait for the component vendor to release updated software.

Mobile devices. Smartphones, tablets, and other mobile devices (excluding laptops) typically run mobile operating systems, and patching for these devices is fundamentally different. It is often necessary to connect the mobile device to a desktop or laptop and to acquire and download updates through that desktop or laptop. Some mobile devices can directly download updates, but this can be problematic because of bandwidth considerations (such as taking a long time to download large updates and paying data charges for the downloads). Another option for keeping mobile devices updated is the use of enterprise mobile device management software. Enterprise mobile device management software is used to manage mobile devices, even personally owned devices not controlled by the organization. It can install, update, and remove applications, and it can restrict enterprise access if the phone’s operating system and mobile device management software are not up to date. See Section 3 of SP 800-124 Revision 1, Guidelines for Managing and Securing Mobile Devices in the Enterprise, for more information.

Operating system (OS) virtualization. Patches need to be maintained for every OS image and snapshot used for full virtualization. Patching capabilities are often built into virtualized environments, such as the ability to patch offline images and quarantine dormant virtual machine instances. See NIST SP 800-125, Guide to Security for Full Virtualization Technologies, for additional information—specifically, Section 3.3 discusses virtual machine image and snapshot management.

Firmware. Firmware updates, such as updating the system BIOS, generally require special privileges and involve different procedures than other types of updates. See NIST SP 800-147, BIOS Protection Guidelines, for additional information on BIOS updates.

40 · 6 MANAGING SOFTWARE PATCHES AND VULNERABILITIES

Organizations should carefully consider all alternative host architectures in use for the enterprise when designing enterprise patch management policies and solutions.

40.3.4 Other Challenges. This section briefly discusses other challenges not covered earlier in this section.

40.3.4.1 Software Inventory Management. Enterprise patch management is dependent on having a current and complete inventory of the patchable software (applications and operating systems) installed on each host. This inventory should include not only which software is currently installed on each host, but also what version of each piece of software is installed. Without this information, the correct patches cannot be identified, acquired, and installed. This inventory information is also necessary for identifying older versions of installed software so that they can be brought up to date. A major benefit of updating older versions is that it reduces the number of software versions that need to be patched and have their patches tested.

40.3.4.2 Resource Overload. Enterprise patch management can cause resources to become overloaded. For example, many hosts might start downloading the same large patch (or bundle of patches) at the same time. This could consume excessive network bandwidth or, if the patches are coming from an organization patch server, overwhelm the resources of that server. Organizations should ensure that their enterprise patch management can avoid resource overload situations, such as by sizing the solution to meet expected volumes of requests, and staggering the delivery of patches so that the enterprise patch management system does not try to transfer patches to too many hosts at the same time.

40.3.4.3 Installation Side Effects. Installing a patch may cause side effects to occur. A common example is the installation inadvertently altering existing security configuration settings or adding new settings. This may create a new security problem in the process of fixing the original vulnerability via patching. Organizations should be capable of detecting side effects, such as changes to security configuration settings, caused by patch installation.

40.3.4.4 Patch Implementation Verification. As discussed in Section 40.3.1 of this chapter, an installed patch might not take effect until the affected software is restarted or other state changes are made. It can be surprisingly difficult to examine a host and determine whether or not a particular patch has taken effect. This is further complicated when there is no indication for a patch when it would take effect (reboot required/not required, etc.). One option is to attempt to exploit the vulnerability, but this is generally only feasible if an exploit already exists, and there are substantial risks with attempting exploitation, even under highly controlled conditions. Organizations should use other methods of confirming installation, such as a vulnerability scanner that is independent from the patch management system.

40.3.4.5 Application Whitelisting. Application whitelisting technologies can conflict with patch management technologies because the application whitelisting technologies function based on known characteristics of executables and other application components, which may be changed by patching. If the vendor is providing the whitelist information, the vendor will have to acquire the patch, record its files' characteristics, and send the corresponding information to customers. If the organization

ENTERPRISE PATCH MANAGEMENT TECHNOLOGIES 40 · 7

is building its own whitelist information, it will have to acquire each patch, record its files' characteristics, and update its whitelists with the new information. Either method may cause problematic delays for organizations that apply patches quickly, especially automatically; patched software may be seen as unknown software and prohibited from running.

To avoid these problems with updates, most application whitelisting technologies offer maintenance options. For example, many technologies allow the administrator to select certain services (e.g., patch management software) to be trusted updaters. This means that any files that they add to or modify on a host are automatically added to the whitelist. Similar options are available for designating trusted publishers (i.e., software vendors), users (such as system administrators), sources (such as trusted network paths), and other trusted entities that may update whitelists. Organizations using application whitelisting technologies should ensure that they are configured to avoid problems with updates.

40.4 ENTERPRISE PATCH MANAGEMENT TECHNOLOGIES. This section provides an overview of enterprise patch management technologies. It discusses their composition, focuses on the security and management capabilities that they provide, and gives recommendations for their use.

40.4.1 Components and Architecture. Enterprise patch management technologies are similar architecturally to other enterprise security solutions: one or more centralized servers that provide management and reporting, and one or more consoles. Enterprise patch management technologies can also be offered as a managed service. What distinguishes enterprise patch management technologies from each other architecturally are the techniques they use to identify missing patches. The three prevalent techniques are agent-based, agentless scanning, and passive network monitoring. Many products support only one of these techniques, while other products support more than one. All the techniques are explained in more detail in the following subsections. Organizations should carefully consider the advantages and disadvantages of each technique when selecting enterprise patch management technologies.

40.4.1.1 Agent-Based. An agent-based patch management technology requires an agent to be running on each host to be patched, with one or more servers that manage the patching process and coordinate with the agents. (Note that agent-based patch management technology is built into some operating systems.) Each agent is responsible for determining what vulnerable software is installed on the host, communicating with the patch management servers, determining what new patches are available for the host, installing those patches, and executing any state changes needed to make the patches take effect (e.g., application restart, OS reboot). Each agent runs with administrator privileges so it can perform these actions. The patch management server is responsible for providing the agents with information on vulnerable software and available patches, including where patches can be acquired and what state changes are needed.

Compared to agentless scanning and passive network monitoring, agent-based patch management technologies are strongly preferred for hosts that are not on the local network all the time, such as telecommuter laptops and smartphones.

There are a few limitations to agent-based patch management technologies. Hosts that don't permit direct administrator access to the operating system, such as many

40 · 8 MANAGING SOFTWARE PATCHES AND VULNERABILITIES

appliances, generally cannot run agents. Also, agents may not be available for all of the organization's platforms.

40.4.1.2 Agentless Scanning. An agentless scanning patch management technology has one or more servers that perform network scanning of each host to be patched and determine what patches each host needs. Generally, agentless scanning requires the servers to have administrative privileges on each host, so that they can return more accurate scanning results and they have the ability to install patches and implement state changes on the hosts (application restarts, OS reboots, etc.).

The main advantage of agentless scanning is that it doesn't require the installation and execution of an agent on each host.

One of the primary limitations of agentless scanning is that it omits hosts not on the local network, such as telecommuter laptops and mobile devices. Also, network security controls (e.g., host-based firewalls) and network technologies (e.g., network address translation) may inadvertently block scanning or otherwise negatively affect scanning results. Agentless scanning may also negatively impact operations by consuming excessive amounts of bandwidth. Finally, agentless scanning may not support all of the organization's platforms.

40.4.1.3 Passive Network Monitoring. Passive network monitoring technologies for patch management monitor local network traffic to identify applications (and in some cases, operating systems) that are in need of patching.

These technologies can be effective at identifying hosts that are not being maintained by other patch management solutions (agent-based, agentless scanning). They do not require any privileges on the hosts to be monitored, so they can be used to monitor the patch status of hosts that the organization does not control (unmanaged systems, visitor systems, contractor systems, etc.).

The primary disadvantage of passive network monitoring is that it only works with software where you can identify the version based on its network traffic (assumed to be unencrypted). Also, of course, it only works with hosts on the local network.

40.4.2 Security Capabilities. This section describes common security capabilities provided by patch management technologies, divided into three categories: inventory management, patch management, and other.

40.4.2.1 Inventory Management Capabilities. Patch management technologies typically have capabilities for identifying which software and versions of software are installed on each host, or alternately, just identifying vulnerable versions of software that are installed. In addition, some products have features for installing new versions of software, installing or uninstalling software features, and uninstalling software.

40.4.2.2 Patch Management Capabilities. Patch management technologies obviously provide a range of patch management capabilities. Common features include identifying which patches are needed, bundling and sequencing patches for distribution, allowing administrators to select which patches may or may not be deployed, and installing patches and verifying installation. Many patch management technologies also allow patches to be stored centrally (within the organization) or downloaded as needed from external sources.

ENTERPRISE PATCH MANAGEMENT TECHNOLOGIES 40 · 9

40.4.2.3 Other Capabilities. Many host-based products that have patch management capabilities also provide a variety of other security capabilities, such as antivirus software, configuration management, and vulnerability scanning. Further discussion of these capabilities is outside the scope of this chapter.

40.4.3 Management Capabilities. Once a patch management technology has been selected, its administrators should design a solution architecture, perform testing, deploy and secure the solution, and maintain its operations and security. This section highlights issues of particular interest with the management—the implementation, operation, and maintenance—of patch management technologies, and provides recommendations for performing them effectively and efficiently.

40.4.3.1 Technology Security. Deploying enterprise patch management tools within an enterprise can create additional security risks for an organization; however, a much greater risk is faced by organizations that do not effectively patch their systems. Such tools usually increase security far more than they decrease security, especially when the tools contain built-in security measures to protect against security risks and threats. The following are some risks with using these tools:

- A patch may have been altered (inadvertently or intentionally).
- Credentials may be misused.
- Vulnerabilities in the solution components (including agents) may be exploited.
- An entity could monitor tool communications to identify vulnerabilities (particularly when the host is on an external network).

Organizations should reduce these risks through the application of standard security techniques that should be used when deploying any enterprise-wide application. Examples of countermeasures include the following:

- Keeping the patching solution components tightly secured (including patching them)
- Encrypting network communications
- Verifying integrity of patches before installing them
- Testing patches before deployment (to identify corruption)

40.4.3.2 Phased Deployment. Organizations should deploy enterprise patch management tools using a phased approach. This allows process and user communication issues to be addressed with a small group before deploying the patch application universally. Most organizations deploy patch management tools first to standardized desktop systems and single-platform server farms of similarly configured servers. Once this has been accomplished, organizations should address the more difficult issue of integrating multiplatform environments, nonstandard desktop systems, legacy computers, and computers with unusual configurations. Manual methods may need to be used for operating systems and applications not supported by automated patching tools, as well as some computers with unusual configurations; examples include embedded systems, industrial control systems, medical devices, and experimental systems. For such computers, there should be a written and implemented procedure for the manual patching process.

40 · 10 MANAGING SOFTWARE PATCHES AND VULNERABILITIES

40.4.3.3 Usability and Availability. Organizations should balance their security needs with their needs for usability and availability. For example, installing a patch may “break” other applications; this can best be addressed by testing patches before deployment. Another example is that forcing application restarts, OS reboots, and other host state changes is disruptive and could cause loss of data or services. Again, organizations need to balance the need to get patches applied with the need to support operations. A final example, particularly important for mobile devices, is the acquisition of updates over low-bandwidth or metered connections; it may be technically or financially infeasible to download large patches over such connections. Organizations should make provisions for ensuring that their enterprise patching solution works for mobile hosts and other hosts used on low-bandwidth or metered networks.

40.5 METRICS AND MEASURES. As explained in Section 3.3 of NIST SP 800-55 Revision 1, Performance Measurement Guide for Information Security, there are three types of measures:

1. Implementation measures are used to demonstrate progress in implementing security programs, specific security controls, and associated policies and procedures...
2. Effectiveness/efficiency measures are used to monitor if program-level processes and system-level security controls are implemented correctly, operating as intended, and meeting the desired outcome...
3. Impact measures are used to articulate the impact of information security on an organization’s mission...⁴

Regarding these types of measures, “less mature information security programs need to develop their goals and objectives before being able to implement effective measurement. More mature programs use implementation measures to evaluate performance, while the most mature programs use effectiveness/efficiency and business impact measures to determine the effect of their information security processes and procedures.” Accordingly, organizations should implement and use appropriate measures for their enterprise patch management technologies and processes.

Examples of possible implementation measures include:

- What percentage of the organization’s desktops and laptops are being covered by the enterprise patch management technologies?
- What percentage of the organization’s servers have their applications automatically inventoried by the enterprise patch management technologies?

Examples of possible effectiveness/efficiency measures include:

- How often are hosts checked for missing updates?
- How often are asset inventories for host applications updated?
- What is the minimum/average/maximum time to apply patches to X% of hosts?
- What percentage of the organization’s desktops and laptops are patched within X days of patch release? Y days? Z days? (where X, Y, and Z are different values, such as 10, 20, and 30).

NOTES 40 · 11

- On average, what percentage of hosts are fully patched at any given time? Percentage of high impact hosts? Moderate impact? Low impact?
- What percentage of patches are applied fully automatically, versus partially automatically, versus manually?

Examples of possible impact measures include:

- What cost savings has the organization achieved through its patch management processes?
- What percentage of the agency's information system budget is devoted to patch management?

40.6 FURTHER READING

- Jang, M. (2006). *Linux® Patch Management: Keeping Linux® Systems Up To Date*. Prentice-Hall, 2006.
- Nicastro, F. M. *Security Patch Management*, 2nd ed. CRC Press, 2013.
- Taylor, J. R., J. H. Allen, G. L. Hyatt, and G. H. Kim. *Change and Patch Management Controls: Critical for Organizational Success*. Global Technology Audit Guide: The Institute of Internal Auditors, Inc., 2005. www.aicpa.org/interestareas/informationtechnology/resources/itassuranceservices/guidance/downloadabledocuments/gtag2patchmgmtcontrols.pdf
- USGAO. *Information Security: Continued Action Needed to Improve Software Patch Management*. U.S. Government Accountability Office, 2011.
- USGAO. *Information Security: Effective Patch Management is Critical to Mitigating Software Vulnerabilities*. U.S. Government Accountability Office, 2011.

40.7 NOTES

1. This chapter is based on the National Institute of Standards and Technology (NIST) Special Publication 800-40 Version 2.0, *Creating a Patch and Vulnerability Management Program*, and Special Publication 800-40 Revision 3, *Guide to Enterprise Patch Management Technologies*, which are available at <http://csrc.nist.gov/publications/nistpubs/>
2. NIST SP 800-53 Revision 4 is available for download at <http://csrc.nist.gov/publications/PubsSPs.html#800-53-rev4>
3. The PCI DSS is available for download at www.pcisecuritystandards.org/security_standards/
4. NIST SP 800-55 Revision 1 is available for download at <http://csrc.nist.gov/publications/PubsSPs.html#800-55>

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 41

ANTIVIRUS TECHNOLOGY

Chey Cobb and Allysa Myers

41.1 INTRODUCTION	41·1			
41.1.1 Antivirus Terminology	41·2	41.3.5 Reputation-Based Scanning	41·9	
41.1.2 Antivirus Issues	41·3			
41.2 ANTIVIRUS BASICS	41·3	41.4 CONTENT FILTERING	41·10	
41.2.1 Early Days of AV Scanners	41·4	41.4.1 How Content Filters Work	41·11	
41.2.2 Validity of Scanners	41·4	41.4.2 Efficiency and Efficacy	41·11	
41.2.3 Scanner Internals	41·5	41.5 ANTIVIRUS DEPLOYMENT	41·12	
41.2.4 Antivirus Engines and Antivirus Databases	41·6	41.5.1 Desktops Alone	41·12	
		41.5.2 Server-Based Antivirus	41·12	
		41.5.3 Mobile Devices	41·13	
41.3 SCANNING METHODOLOGIES	41·7	41.6 POLICIES AND STRATEGIES	41·13	
41.3.1 Specific Detection	41·7			
41.3.2 Generic Detection	41·8	41.7 CONCLUDING REMARKS	41·13	
41.3.3 Heuristics	41·8			
41.3.4 Intrusion Detection and Prevention	41·9	41.8 FURTHER READING	41·14	

41.1 INTRODUCTION. For over three decades, malware has been a persistent, annoying, and costly threat, and there is no end in sight to the problem. There are many vendors offering to provide a cure for viruses and malware, but the mere existence of these software pests is understandably vexing to those charged with system security.

Initially, most viruses were not designed to cause harm but were created more to gain notoriety for the creator or as a prank. Because these early viruses were designed to subvert legitimate program operations across multiple systems, they were more likely to cause unexpected problems because the virus writers didn't do exhaustive testing. These viruses, and later some Trojans, often damaged data and caused system downtime. The cleanup required to recover from even a minor virus infection was expensive in terms of lost productivity and unbudgeted labor costs.

Viruses and Trojan behavior have merged, and now both are considered as part of the larger family referred to as malware. No longer is malware just written for a virus writer's 15 minutes of fame; today, malware is created primarily for financial gain. Malware can still cause damage, but now it is more likely to have been created

41 · 2 ANTIVIRUS TECHNOLOGY

to steal valuable information or the resources of an infected computer. Malware can steal valuable data for sale on underground markets, send spam from a compromised machine, or install software designed to eavesdrop on an affected user's activities, to name a few examples. Some malware programs can give full access to a hacker to control infected machines, and malware continues the trend of controlling affected machines in a "botnet." These collections of hijacked computers are used for a variety of activities, such as launching attack or sending spam. A botnet makes use of the increased power of networked computers, while making it more difficult to trace a suspected attacker. For more information about malware, see Chapters 16, 17, and 20 in this *Handbook*.

Because of this shift to a financial motivation for malware writing, something of an open-source community has sprung up to develop new malware. There have been two major areas where development has been most notable: evasion of traditional antivirus defenses and distribution by exploiting vulnerabilities in common software. In their efforts to evade traditional antivirus defenses, criminals often slightly modify and compress malware with *packers*, which frequently use sophisticated obfuscation or anti-reverse-engineering techniques. The primary focus of today's malware writers is to avoid the newsworthy outbreaks of the past, staying under the radar to accomplish targeted attacks (*phishing* and *spear phishing*). Some of the most common malware families now have as many as tens of thousands of variants with hundreds of variants being released each day. Once they are in a system, they often employ stealth techniques such as kernel-mode rootkits to hide themselves from the operating system and from many antivirus (AV) products.

New malware programs are being developed every day because vulnerabilities in software programs are so numerous, and it is often weeks between the time a vulnerability is discovered and the time that the vendor issues a security patch. Even more often, patches are announced by vendors but not implemented by users. Malware writers take advantage of these vulnerabilities and time lags, so research into vulnerabilities is imperative in order to create effective defenses with antivirus programs. A layered defense against malware infection is essential and should employ a modern antivirus scanner as well as a firewall and data encryption.

Although the threat of new malware remains, and its growth has become exponential, the technology deployed to defend computers against viruses and malware is one of the least understood aspects of security architecture. As a result, antivirus defenses are often set up improperly. This chapter describes available antivirus technologies and how they work and outlines methods for their effective application within a comprehensive program of computer security.

41.1.1 Antivirus Terminology. The acronym AV is widely used to describe the industry, the products (also sometimes abbreviated AVPs), and the programs that have been developed to defeat computer viruses.

Early AV programs used simple hard- and floppy-drive scans to search for a specific text string hidden within a specific file or a boot sector. This is the origin of the term *AV scanner*, which is now widely used as a generic term for all AV programs. That is how the term is used in this chapter, although it should be pointed out that many of today's AV scanners do much more than merely scan for known malware, and detection may now include unwanted programs as well as suspicious file activities.

When exploring AV product literature, it is common to see statements about the number of viruses and Trojans that exist. It is not unusual to see estimates of the total number of malware variants in the hundreds of millions, and six-digit estimates

ANTIVIRUS BASICS 41 · 3

of the number of new malware added every day. AV vendors' methods for counting differ from one to the next, and these numbers tell us very little about their detection capability. It is far more instructive to rely on the reports of reputable AV testing labs to get a clear idea of a product's effectiveness, than to rely on these numbers.

41.1.2 Antivirus Issues. In the early days of viruses, aside from having the ability to spread to other machines, there was little benefit in keeping a virus on a machine for any length of time. During the early days, viruses were relatively easy to find, because a virus usually caused noticeable damage that made users aware that their systems were infected. At other times, an unmistakable message or image appeared on the screen, alerting users to the infection. As virus writers discovered a monetary benefit to remaining on a system as long as possible, the number of affected machines increased tremendously. Viruses and malware have become complex, are adept at disguising themselves, and are not as easy to find and eradicate. To counter this problem, AV scanners have become sophisticated logic machines. This new level of complexity, of both the problem and the solution, has made it even more confusing for users. Too many people fail to understand the importance of having up-to-date AV products. Users do not know what to look for, and do not understand that viruses and malware do not announce their presence.

Modern AV software has an automatic-update capability, which has made updating a less labor-intensive process. Updates are the key to keeping systems malware-free because updates give the AV program the information necessary to find the newest viruses and malware. Unfortunately, because malware is released at such a prodigious pace, even daily updates of a traditional string-based scanner may not be enough to keep systems totally clear of threats. Many AV products now support essentially constant contact with the developers' servers and verify checksums on the servers to identify changes requiring a download. These methods allow updates within minutes of any change with no inconvenience to the users.

Most AV scanners look at what has happened before (i.e., known behavior of known malware), in order to detect infections. Although AV products now attempt to anticipate new threats, that is not their strong point. Additionally, many consumers who buy a new computer receive a working AV product with their purchase, but they do not choose to buy a subscription for AV updates once their free use period expires. They do not understand why they have to pay for something that seems to be working, and they are lulled into a false sense of security by thinking that their AV product is doing its job. Consumers who bring work from unprotected home computers into the office can quickly infect an entire network.

Security products historically have suffered from a lack of upper-management support because they are often viewed as high-cost/low-return items, and they are given a low priority in the security budget. The information technology (IT) team of an organization needs resources for things like identifying and patching vulnerable systems, monitoring and managing machines entering their network, monitoring suspicious behavior identified by firewalls, AV products, and security scanners. Considering the wide array of AV products from which to choose, and the large number of patches introduced each day, many system administrators despair. They will install and patch just what they have time for, and consider it good enough.

41.2 ANTIVIRUS BASICS. AV scanners are a bit like police officers walking the beat. They try to watch everything that is going on around them, look out for suspicious behavior, and attempt to intercede when they think something bad is happening or about

41 · 4 ANTIVIRUS TECHNOLOGY

to happen. Both the police and AV scanners look for certain patterns and behaviors, and they leap into action when a suspect crosses a predetermined threshold of acceptability. Like the police, AV scanners sometimes reach the wrong conclusions. These errors are usually caused by insufficient data or by new and unexpected behavioral patterns.

Malware detection is an inexact science, and it is impossible to create an AV scanner with a 100 percent success rate. It is simply not possible to know the intent of every bit of code that enters a computer, and it is not feasible to test every bit of code before it executes. To do so would require that the AV scanner demand so much of the processing power of the CPU that valid programs would not be able to execute. Malware behaviors are subject to broad variations, and many use stealth or cloaking techniques to hide from the operating system and even from the AV scanner itself. There are no longer hard-and-fast rules that a user can apply to determine if a system harbors malware.

41.2.1 Early Days of AV Scanners. When viruses first started appearing with regularity in the late 1980s, their detection and eradication was relatively straightforward, but not necessarily easy. The viruses were quite simple and normally did not spread very quickly. The AV community quickly researched them, determined what made them work, and published effective fixes in short order. These fixes tended to be written for a specific virus and could not disinfect other viruses, even if they were of similar types. Most of the work fell on users to identify which virus (or type of virus) they thought they had and then search for a program that would fix it. Because Internet connectivity was not as prevalent as it is now, users frequently spent much time calling friends and associates in the hope that they could forward a disk copy of the necessary AV program. Additionally, there were no naming conventions for viruses, and it was difficult to determine with any conviction that the fix obtained would actually work.

Viruses of that period generally inserted their code in predictable sections of a program. The early scanners ran a search for a specific string of characters. If they found it, they would delete the virus code and attempt to restore the host program to its original uninfected form. Failing that, the scanner would usually advise the user that disinfection was incomplete and that they should delete the infected application and reinstall it.

As the number of viruses began to climb, software companies that had ventured into the AV market began to realize that creating and distributing individual fixes was no longer feasible. Instead, they began to develop more comprehensive scanners that could look for more viruses, both old and new. The new generations of scanners were comprised of two components: the scanning engine and the signature files. Each component was entirely dependent on the other to work. The engine consisted of the user interface and the application that scanned the system for viruses. The signature files were a database of the *fingerprints* (unique segments of code) of known viruses. Although some of these early scanners did a good job, many did not. None of the early AV scanners was able to catch all known viruses.

41.2.2 Validity of Scanners. The vendors of software scanners in the late 1980s and early 1990s faced a number of obstacles. It seemed there was a new AV vendor appearing every month, and the market became highly competitive as user awareness of the virus problem grew. Given this competitive state, there was vast dissension among the AV community as to how viruses for research should be stored and tested. Many AV vendors kept a library of viruses for their own use, and this fact was used in their marketing. Claims that one program worked better than another because it checked for more viruses were misleading because no one knew how many

ANTIVIRUS BASICS 41 · 5

viruses existed. There was simply no method of commercial or independent testing to check the validity of claims made by the AV product vendors. Additionally, there was a problem of naming the viruses. Each vendor created its own names for viruses, and it was not uncommon for one virus to be known by several names.

The AV vendors also disagreed on how AV scanners should operate in principle. Some vendors felt that AV scanners should only look for new viruses, and others felt that a good product should search for both old and new viruses. While this argument raged, viruses looked as though they would eventually gain the upper hand, especially as virus writers began to use underground bulletin boards, and later the Internet, to share and distribute virus code.

With no standards for the AV products, the public had little to go by other than the vendors' marketing copy and the advice of other users. However, if a recommendation was made by a friend for Brand X Antivirus because no viruses were found on the friend's system, it was possible that no viruses had ever been introduced in the friend's system at all, and Brand X could not find old viruses, new viruses, or any viruses at all.

Two things happened that revolutionized the AV scanner market. In 1993, Joe Wells, a research editor with a business magazine, began collecting viruses and virus reports from experts around the world and began assembling a library of these viruses. He named this library of viruses the *WildList* and made it available to legitimate AV researchers. His list divided viruses into those known to have infected systems (in the wild) and those that had been written but were not actively infecting (in the zoo). A naming convention of viruses also began to emerge in order to maintain an efficient and searchable database.

The other notable event was the development of commercial AV testing and certification by a company known as the National Computer Security Association (NCSA), which is now known as ICSA Labs. The NCSA started a consortium of AV vendors that, for a fee, submitted their products to be tested. NCSA and Joe Wells began collaboration for the use of his WildList, and Dr. Richard Ford, a noted virus expert, created a virus-testing laboratory for NCSA. Dr. Ford fashioned an environment in which AV scanners were put through their paces to see if they could detect all of the viruses in the WildList. AV vendors submitted their products every time a new version of their product was about to be released. Although the original test results were dismal (many scanners could not detect more than 80 percent of the viruses in the list), an environment had been created in which measurable improvements in the effectiveness of AV technology could be achieved. Naturally, the public and the press began to look for AV products that had been certified by NCSA. Eventually, other commercial and independent test laboratories independently developed their own certification and testing schemes to help users find reliable AV products.

Before long, other testing organizations joined the fray and AV products increased in complexity. It was clear that there was a need for testing standards to help users distinguish effective, unbiased tests from those that did not clearly demonstrate the capabilities of AV products. A group of AV researchers, academics, testers, and other interested parties created an organization to develop such standards, called the Anti-Malware Testing Standards Organization (AMTSO). This group has published a number of papers outlining considerations for how to produce or distinguish a good test from an inadequate one. These papers are meant to help both testers and those people reading tests.

41.2.3 Scanner Internals. As was noted earlier, an AV scanner cannot simply put each program into a computer's RAM and test it for malware before the program is

41 · 6 ANTIVIRUS TECHNOLOGY

allowed to execute. To do that would require almost all the resources of the CPU, and users would have a system that operated at a snail's pace.

In order to operate efficiently, AV scanners have had to go to great lengths to improve their speed in order to check for the vast number of existing malware without bringing the entire system to a halt. Using sophisticated elimination methods to exclude signatures whose characteristics do not match those of the file being scanned, helps keeps speeds reasonable.

There are five basic types of detection within a modern AV product. Most people are familiar with signature-based detection, which is also known as specific detection, but this is not the only technique in place anymore. AV products may use some or all of these five basic methods of operation:

1. **Specific detection.** Looking for known malware
2. **Generic detection.** Looking for infections by variants of known malware
3. **Heuristics.** Scanning for previously unknown viruses by noting suspicious behavior or file structures, as described more fully in Section 41.3.3
4. **Intrusion prevention.** Monitoring known-suspicious system changes and behaviors to prevent suspected malware
5. **Reputation.** Detecting malware by a file or Website's reputation, based on a number of factors including the number of times the file or Website has been downloaded, customer reports, or age of a Website

41.2.4 Antivirus Engines and Antivirus Databases. The AV engine and its signature database work in concert to prevent and detect malware trying to enter a system. The engine generally provides a library of commonly used functions. It consists of dozens of complex searching algorithms, CPU emulators, and various forms of programming logic. The engine determines which files to scan, which functions to run, and how to react when suspected malware is found. However, the engine knows absolutely nothing about the malware themselves and is almost useless without the signature database.

The signature database contains the fingerprints (snippets of distinctive code) of hundreds of millions of malware variants. As new malware and variants appear at an accelerating rate, it is imperative that the signature database be updated often. In 1995, the experts advised updating the database files at least once a month, but with so many viruses appearing each day, users today are advised to update at least daily—and ideally, to allow constant updates controlled by their AV product. AV manufacturers now provide products that check for updates automatically and download changes whenever a user is connected to the Internet. Some vendors also provide cloud-based signature databases so that computers that are connected to the Internet can check against a more-frequently updated set of signatures that are hosted on vendors' own remote systems.

The signature database also contains the rule sets used in heuristic scans. These types of scans can be slower and more intrusive than simple signature scans, and their design and implementation vary greatly between products. Most products now give users configurable options to lessen or increase heuristics as desired. Although signature scans can be considered a heuristic in themselves, the term is more commonly used to identify the more complex AV functions that attempt to locate viruses by identifying suspicious behavior and/or file structure.

SCANNING METHODOLOGIES 41 · 7

Because the distinction between a scanning engine and a signature database is not obvious to many system administrators, many religiously update the database but are unaware that the engine also may need updating. This is a poor strategy that can result in many viruses slipping by the scanner undetected.

41.3 SCANNING METHODOLOGIES. AV products are configurable by the user or the system administrator to scan upon startup, constantly, or on demand. To be at its most effective, a scanner should be set to a continuous or “on access” scan, with a periodic, scheduled scan set to occur when the system is on but not in use. Users running less-optimized AV programs may find that this degrades system performance. Some scanners need to use much of the system’s memory on continuous scans in order to be able to test sections of code, which may make the applications noticeably slower when they first run. Therefore, a happy medium must be found—the AV scanner must be able to protect the system, while the user must be able to have full use of the system.

There is no one scanning method that is superior to the others. All of the scanning methods have their advantages and disadvantages, but none is able to detect malware with unfailing accuracy. A scan is looking for code and behaviors that have been noticed in other malware, and if a new malware exhibits new, previously unknown behaviors, it can pass by undetected. Therefore, most AV scanners do not rely on only one scanning method to detect malware, but have several included in their design.

41.3.1 Specific Detection. Each malware uses different code to perform its functions. A sequence of code that is specific to each malware is referred to as the fingerprint, or signature of that malware. To detect the presence of malware, the scanner looks for the signature, removes its code from the host file or system, and attempts to restore the infected program or system to an uninfected state. In early viruses, it was discovered that the signatures were usually found within specific areas of a program, specific to each virus. The scanners set out to inspect only those areas of a file rather than scanning an entire program from top to bottom. This saved vast amounts of time and processing power.

As malware is often static, rather than parasitic in nature, AV researchers now must get more creative to quickly identify known-bad files. But as malware authors have armored their files to make their creations difficult for researchers to analyze, this can also make a malicious file look very different from a valid document or application. As such, researchers can key in on these differences and swiftly exclude benign files.

Every vendor’s AV product has a different implementation of scanner and database, although the signature scanning technique is the most common. Signature scans can identify whether a program contains one of the many signatures contained in the database, but it cannot say for certain whether a system has actually been affected by malware (e.g., the malware may be present but not yet executed). Users can only trust the guess of the AV scanner, because the odds are in the scanner’s favor. It is possible, however, that a program that is suspected of being infected actually contains random data that only coincidentally looks like a virus signature. The legitimate program could contain instructions that by sheer chance matched the search string in the virus database. However, when there is a possibility that the code is actually from a virus, the scanner reports it as a positive hit.

False-positive reports are a possible problem of signature scanners. If users notice that their scanner falsely reports the presence of viruses too often, they view this as an annoyance, and will likely seek to disable the software or find ways of circumventing the scans.

41 · 8 ANTIVIRUS TECHNOLOGY

41.3.2 Generic Detection. As was discussed earlier, malware is now often created for financial purposes, and it makes fiscal sense for its authors to get as much use as possible out of successful creations. They often make open-source code that is shared widely in the malware-author community, so that it is often updated with new functionality such as exploit-code or password-stealing capabilities, to target new games, applications, or online banking sites.

Likewise, it makes sense for AV scanners to look for common properties of popular malware families or known malicious behavior in order to proactively detect variants based on those codebases. These generic detections can vary from being more heuristic in nature to being fairly specific. For example, it could include protection as broad as buffer overflow detection to prevent certain types of exploits, to specific homegrown packers used by only one malware family.

Generic detection caused a great deal of controversy in the early days of AV products. There was concern, when most viruses spread by infecting clean files parasitically, that generic detection would require generic cleaning, which could leave host files mangled beyond repair or usability. As the vast majority of malware now infects systems rather than host files, simple deletion of malicious files and cleaning of registry entries can remove malware. But because malware is also frequently component-based, removing one threat may be only the tip of the iceberg. There may be other components that are as yet undetected, and it's best to thoroughly examine any system which has had malware detected. As such, more and more system administrators have adopted a policy of simply re-imaging affected systems rather than relying on AV products' removal procedures.

41.3.3 Heuristics. By adding heuristics to their AV scanners, vendors looked to increase the efficacy of their products. The scanners could now look for malware that are new and unknown and not contained within the signature database.

The word *heuristic* comes from a Greek word meaning “to discover.” The term is used today in computer science to describe algorithms that are effective in solving complex questions quickly. A heuristic algorithm makes certain assumptions about the problem it is trying to solve. In the case of an AV scanner, it analyzes a program’s structure, its attributes, and its behavior to see if these meet the rules that have been established for identifying malware, even without its specific signature being known. Basically, a heuristic algorithm works on the assumption that if it “looks like a duck, walks like a duck, and sounds like a duck, it must be a duck.”

The drawback to heuristic scanning is that it makes intelligent assumptions but is nevertheless bound to make mistakes. Another problem with heuristic scanning is that, on slower systems, it may take longer to run and may require user interaction. Some users consider this intrusive and may turn off the feature. By combining both signature scanning and heuristic scanning in their products, AV vendors have increased their effectiveness and speed.

Heuristic scanners use a rule-based system to verify the existence of malware. It applies all the rules to a given program and gives the program an overall score. If the score is high, there is a good likelihood that it is malicious. Generally, the scanner first looks for features of a file that are more common to malicious files. A well-designed heuristic scanner will limit the regions of the program to be examined in order to scan the highest number of suspects in the shortest possible time. The scanner then examines the logic of the suspected program to determine if it might be trying to perform known, likely malicious actions. This type of scanning is considered to be a *static* scan. The

SCANNING METHODOLOGIES 41 · 9

static method applies the rules and gives a pass/fail score to the program—whether the program has actually executed or not.

The other type of heuristic scanning is called the *dynamic method*. This method applies basically the same rules as the static method, and if the score is high, it attempts to emulate the program. Rather than examining the logic of the suspected code, the dynamic scanner runs a simulation of the virus in a virtual environment. This technique has come to be known as *sandbox* emulation, and is effective for attempting to identify new malware that does not appear in the signature database.

Neither of these heuristic scanning methods is necessarily better than the other, but in concert they give fairly good results. Although a static heuristic scan may miss some malware because they have not yet executed, the dynamic heuristic scan can catch previously unknown malware before they are allowed to execute on the system. Between heuristic and generic detections, it is becoming increasingly common for at least one AV vendor to identify any new malware as soon as it is released.

41.3.4 Intrusion Detection and Prevention. As more complex malware appears at an increasingly prodigious pace, scanning programs solely for signatures has become less effective at finding viruses. Virus authors use encryption or obfuscation techniques, or release a large number of individual variants in the hopes that the AV scanner will not find them. Today's operating systems and legitimate programs have bloated to millions of lines of code, so that finding a virus signature may be resource intensive. Because malware may steal valuable data or damage systems, it is not a good strategy to let them execute and then attempt to clean up the mess. A better strategy is to try to find malware before it has had a chance to affect a system and to prevent it from doing harm.

The use of a cyclical redundancy check (CRC), or checksums, was originally added to some security products, including AV suites, to aid in the detection and prevention of virus infections. This method tracks unauthorized file and system changes, as when malware or a hacker enters and alters a system. To track those changes, a fingerprint of each executable program and data file is computed and stored in a database when the AV product is first installed. These fingerprints are quite small, usually consisting of less than 100 bytes of information—this is the “sum” or checksum. Because malware must add or change files on a system in order to affect it, the checksums of the fingerprints are compared with any newer version. If the checksums vary, then the AV scanner runs other routines to investigate further.

In the case of intrusion-prevention systems, a behavior-based scan is executed on each new file when it is run. If certain suspicious behaviors are observed, a user may be asked whether the behavior should be allowed to continue. If a sequence of behaviors is observed that is sufficiently malicious, the program may be halted entirely. This technique has been found to be remarkably effective in preventing execution of new, unknown malware, although it shares the same sort of difficulties found with using heuristic scans. Again, as part of a complete security arsenal, it can be a valuable tool. For more information on intrusion detection and intrusion prevention, see Chapter 27 in this *Handbook*.

41.3.5 Reputation-Based Scanning. As new malware is released at such a prodigious rate, it is no longer practical for researchers to investigate each sample individually. While tens or even hundreds of thousands of samples are investigated daily, there are still a number of threats that slip through the cracks. In order to deal with this constant deluge, rather than relying solely on a handful of researchers to

41 · 10 ANTIVIRUS TECHNOLOGY

look for known-malicious files or Websites, companies are beginning to rely on the general public to increase their view of the malware landscape. Reputation scanning gathers telemetry and reputation data provided by a product's users, plus information about the age of files and Websites, and then rates them as trusted, bad, or somewhere in-between. This can help a user determine whether to proceed, based on the relative risk level.

This sort of scanning has limitations. Certain types of files are unique by design, especially data files such as documents, spreadsheets, and presentations. And if a user chooses to run a file that turns out to be malicious, they're on their own when it comes to figuring out how to remove the threat from their system. But this is why AV products include a variety of different scanning technologies; to provide protections that will ideally overlap in such a way that they can catch a larger percentage of brand-new, malicious code.

41.4 CONTENT FILTERING. In the early days of virus infections, computer security experts often allayed the fears of computer users by telling them that they could never catch a computer virus from email. This assurance was based on the fact that email was almost exclusively composed of ASCII text documents, with no ability to execute program code. At the same time, the skeptics were saying “never say never.” The skeptics won.

First, there were several waves of macro virus-infected documents sent as file attachments to email. This led to the modified assurance from security experts that no one could ever catch a computer virus from an email message attachment, if the attachment was not opened. Then virus writers started embedding commands that use HTML and the scripting capability of email programs. This led to the further modified assurance from experts that a computer virus could not be caught from unopened email. This assurance in turn proved unwarranted, because email preview capabilities were exploited to trigger malicious code even without user intervention. At one point, merely highlighting a message subject in MS Outlook was enough to execute an attachment, although this default was later changed.

Virus writers also began to exploit the user’s email facility by forwarding copies of the virus to entries in the user’s email address book. The Melissa virus was the first virus that really leveraged email to spread rapidly. Since the Melissa virus, users have been advised to suspect just about any unsolicited email. Malware writers are always looking for new delivery methods, and they were richly rewarded when email programs began to allow executable code within the email. Although users enjoy the point-and-click convenience of this feature, it allowed new viruses to proliferate at a rate not seen before.

The Web also has seen an explosion in malicious code distribution, particularly since the advent of “Web 2.0”—social networking and collaboration sites. Through Adobe Flash, Java, and JavaScript objects that exploit vulnerabilities in Internet browsers or media-viewer software, malware can be automatically downloaded and installed on users’ machine without their being aware. This is commonly referred to as a *drive-by download*.

Content filtering is one popular way of controlling Web and email threats. It consists of a server-based application that interrogates all incoming and outgoing traffic, according to its configuration and rule sets. Early versions were cumbersome to configure, due to a text-based interface that required all rules to be composed in a laborious text editor. Misconfigurations were commonplace because administrators often were not sure which of the text files was causing a failure.

CONTENT FILTERING 41 · 11

The new generation of content filters has increased user friendliness, using interactive graphic user interfaces to set and adjust policies. Administrators are able to fine-tune policies so that they meet the specific needs of their organizations. For example, all email containing executable attachments may be blocked, quarantined, or simply deleted. This may be established as a rule, for some users or all users.

Content filters were particularly effective in preventing infections of email-borne viruses, even before the specific virus signature was released. For example, when the security officer of one government office heard of the Love Bug virus on the morning news, that person set the content filters to block all email attachments containing the extension ".vbs," thus averting infection of a large network. No user interaction was required, and most users were not aware that the block had been placed. The costs and downtime of a virus attack were prevented.

The majority of threats now come from compromised legitimate Websites rather than the dark, seedy underbelly of the Internet, which changes the way products view unwanted content. Consequently, much of this content filtering technology has moved from a focus on blocking potentially malicious content to blocking Websites that are deemed inappropriate for viewing at work or by children. Many content filters come with extensive, frequently updated block-lists for porn, "hate-speech," and other politically sensitive subjects.

41.4.1 How Content Filters Work. These applications work in the same general manner that AV scanners do. They scan all incoming data on specific ports on the server, and they compare the traffic to rules and strings in the database. Because content filters are capable of blocking more than one type of file or program, they have the ability to scan text files, graphics, zipped files, self-extractors, and various executables. Many content filters contain AV scanning components, so if traffic does contain known-malicious code, it can be intercepted and disinfected before it is sent to the recipient.

The standards for formatting email for transmission using standard protocols have long been in place and detail the type of information in every section. It is easy for a program to look for particular information within these sections to determine what is included in the message—attachments, for example. Content filters first begin by disassembling a message to look at its various parts before scanning the message for the items to be allowed or denied into the system. Before sending the message onward, it is reassembled and checked for the conditions specified in the configurations. For example, a condition may state that any attachments be stripped and deleted, that the message body be sent to the recipient, and that an outward email be sent to the sender stating that attachments are not allowed in email.

In terms of data security, a content filter adds several elements that are beyond the traditional AV scanner. For example, message and attachment content can be scanned for inappropriate material. This might be proprietary company information that employees should not be sending out via email, or it could be offensive material such as pornography, that should not be coming into the company's email system. Content filtering also can stop the spread of common forms of spam mail, such as chain letters and get-rich-quick schemes.

41.4.2 Efficiency and Efficacy. Speed of operations is a concern with content-filtering mechanisms, given the large volume of network traffic in most organizations. However, because the operations are all contained within the server, the users will not notice any change in performance of their desktop systems; that is, the processing will

41 · 12 ANTIVIRUS TECHNOLOGY

be completed before the messages are received by the client systems. For instance, in the case of email, if filtering causes mail to queue up, delivery of mail may lag for a period of time. If Websites are blocked for inappropriate content, the traffic will simply not be delivered to the user.

Content filters are also subject to the same failures as traditional AV scanners. New malware can be missed if the data is not present in the scanning database. Additionally, the configuration of the product and the application of patches and updates are crucial to its successful operation. False positives are also a problem, where a legitimate message inadvertently includes content that triggers a block. It is possible to quarantine questionable messages and have a system administrator follow up with the sender. This can lead to refinement of the filters or to the detection of serious offenses. Before a content-filtering system is deployed, it is important to put in place the response mechanisms, so that abuses of policy can be addressed appropriately. This may well involve several departments besides security, including legal and human resources.

41.5 ANTIVIRUS DEPLOYMENT. AV scanners can be installed on the desktop or on servers. Each strategy has its advantages and disadvantages. For example, if the system is server based, malware on USB thumb drives, DVDs, and CDs on the desktop will not be scanned. The consensus of most experts, however, is to use both. With the advances in AV products and network management systems, it is entirely possible to install scanners both on the desktop and on servers while still maintaining an acceptable level of control and performance.

41.5.1 Desktops Alone. If an organization's computer security policy allows unrestricted use of thumb drives, floppies, DVDs, and CDs, it imperative that AV scanners be deployed to the desktop. Unless these drives are locked or disabled, there is no way, other than scanning, to prevent users from accidentally introducing malware. The preferences of a desktop AV scanner can be set to automatically scan external media.

Updates to desktop AV scanners can now be distributed via a central server. This is particularly effective when new signature files are needed to prevent infiltration by newly discovered malware. The updates can be pushed to the desktop, and the users need not be present at the workstation, although the desktop system must be on and connected to the network at the time. If the updates are scheduled after working hours, it is important to verify that systems have updates pushed to them as soon as they log back into the network. Some AV products have management consoles that will allow this to occur automatically, rather than having someone check each system individually.

In order to prevent unauthorized AV-setup changes, it is possible to prevent the users from changing the configuration of their desktop AV scanners. Since that may not be the default installation, it must be checked. Again, with the use of a management console, it is possible to enforce security-software policies remotely.

41.5.2 Server-Based Antivirus. Many companies have sought to bolster the defenses at the perimeter of their network by installing AV products on the server where downloads are frequently stored and traffic is high. A server-based AV scanner can be configured to send alerts to administrators when suspected malware is detected. Like the desktop-based scanners, the response to malware detection can be predetermined. Many system administrators set the program to erase all infected files rather than to send them to quarantine. This strategy works to lessen the possibility that quarantined malware can be "released" by mistake.

CONCLUDING REMARKS 41 · 13

41.5.3 Mobile Devices. With the increase in use of smartphones and tablets by individuals, organizations, and through bring-your-own-device policies, securing these devices against malware has grown in importance. Host-based intrusion detection and signature-based tools have been developed to protect such devices, but there are problems based on the security model of some operating systems (e.g., Android) because of the strict partition of memory areas controlled by applications.

41.6 POLICIES AND STRATEGIES. In the battle against malware, the promulgation of appropriate policies and the implementation of realistic plans of action are equally as important as the installation of an AV scanner. The policies should spell out in detail what actions are allowed or denied, and they should also be specific about the users' responsibilities. The policies and responsibilities should be regularly updated to reflect the realities of the changing malware landscape, especially within a company's ecosystem.

End-user AV awareness training should be high on the list of priorities in every organization. Users are more likely to cooperate in preventing and quarantining infections if they are aware of the types of malware that may affect their system and of the damage they can cause. A simple security-incident bulletin board in a central location is an easy and effective way to communicate with users. Email is probably not an effective method of distributing malware awareness information because users become confused between the education effort, actual and legitimate malware alerts, and bogus malware alerts.

The roles and responsibilities of each person within an organization should be clearly defined and communicated to the general populace. For example, the responsibilities of an average user will be different from those of a system administrator, and the responsibilities should be reflected in their roles. An individual user's role may describe the actions required of the user if malware is detected on a workstation, while the system administrator's role may describe how to handle the report from the user and prepare for removal.

Problems and catastrophes will usually occur when they are least expected, so every organization should have an emergency response plan in its policies. The emergency plan should detail the list of persons to be called in an emergency and the priority order in which they should be called.

For every major malware incident within an organization, care must be taken that a "lessons learned" session be undertaken as soon after the event as possible. No matter how well-written a policy may be, it cannot be proven effective until it is put into use. An actual infection will highlight the failures of a policy in action, which should be rectified before the next attack. For more on security policy guidelines see Chapter 4 in this *Handbook*.

For more details of computer security incident response team management, see Chapter 56 in this *Handbook*.

Management's support for such policies is vital. Support is required not only to approve the AV budget and the policies, but also to ensure that everyone abides by the policies. It is highly unlikely that users will follow a policy that upper management routinely flouts.

41.7 CONCLUDING REMARKS. Both AV technology and malware technology have progressed rapidly over time, and AV technology largely remains in step with malware technology. The success of malware is now primarily due to financial

41 · 14 ANTIVIRUS TECHNOLOGY

motivation, but a large part of the reason it is such a lucrative business is that people continue to dismiss the threat of malware. However good AV technology gets, it will not make a serious dent in the malware problem unless it is appropriately implemented by organizations and properly employed by users who act responsibly. As AV technology continues to improve and becomes better understood, we hope that it will continue to be used more widely and more wisely.

41.8 FURTHER READING

- AMTSO. “Documents and Principles.” Anti-Malware Testing Standards Organization, 2013. www.amtsot.org/documents.html
- AVG. www.avg.com/us-en/homepage
- Avast. www.avast.com/index
- Aycock, J. D. *Computer Viruses and Malware*. New York: Springer, 2006.
- Bitdefender. www.bitdefender.com
- Drake, J. J., Z. Lanier, C. Mulliner, P. Oliva, S. A. Ridley, and G. Wicherski. *Android Hacker’s Handbook*. Wiley, 2013.
- Elisan, M. *Malware, Rootkits & Botnets A Beginner’s Guide*. McGraw-Hill Osborne Media, 2012.
- ESET. www.eset.com/us
- Johnston, J. R. *Technological Turf Wars: A Case Study of the Antivirus Industry*. Philadelphia: Temple University Press, 2008.
- Kaspersky Lab. <http://usa.kaspersky.com>
- Li, Q., and G. Clark. “Mobile Security: A Look Ahead.” *IEEE Security and Privacy* 11, no. 1 (January 2013): 78–81.
- Ligh, M., S. Adair, B. Hartstein, and M. Richard. *Malware Analyst’s Cookbook and DVD: Tools and Techniques for Fighting Malicious Code*. Wiley, 2010.
- Malin, C. H., E. Casey, and J. M. Aquilina. *Malware Forensics Field Guide for Windows Systems: Digital Forensics Field Guides*. Syngress, 2012.
- McAfee. www.mcafee.com/us
- Miller, C., D. Blazakis, D. DaiZovi, S. Esser, V. Iozzo, and R.-P. Weinmann. *iOS Hacker’s Handbook*. Wiley, 2012.
- Nazario, J. *Defense and Detection Strategies against Internet Worms*. Artech House, 2003.
- Scott, W. *Virus and Malware Removal Made Easy*. Amazon Digital Services, 2013.
- Sikorski, M., and A. Honig. *Practical Malware Analysis: The Hands-On Guide to Dissecting Malicious Software*. No Starch Press, 2012.
- Skoudis, E., and L. Zeltser. *Malware: Fighting Malicious Code*. Prentice Hall, 2003.
- Slade, R. “Antiviral Software Evaluation FAQ.” Doug’s Home on the Web, 2006. www.dmath.org/virus/faqs/evaluation.html
- Symantec | Norton. <http://us.norton.com>
- Szor, P. *The Art of Computer Virus Research and Defense*. Addison-Wesley, 2005.
- Vasudevan, A., J. M. McCune, and J. Newsome. *Trustworthy Execution on Mobile Devices*. Springer, 2013.
- Virus Bulletin. www.virusbtn.com/index

Computer Security Handbook, Sixth Edition, Volume 1
Edited by Seymour Bosworth, Michel E. Kabay and Eric Whyne
Copyright © 2014 by John Wiley & Sons, Inc.

CHAPTER 42

PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

**Robert Guess, Jennifer Hadley,
Steven Lovaas, and Diane E. Levine**

42.1 INTRODUCTION	42·1	42.4.2 Application Examples	42·13
42.1.1 Digital Rights	42·2	42.4.3 Examples	42·13
42.1.2 Patent, Copyright, and Trademark Laws	42·2		
42.1.3 Piracy	42·2		
42.1.4 Privacy	42·3		
42.2 SOFTWARE-BASED ANTIPIRACY TECHNIQUES	42·3	42.5 PRIVACY-ENHANCING TECHNOLOGIES	42·14
42.2.1 Organizational Policy	42·4	42.5.1 Network Proxy	42·14
42.2.2 Software Usage Counters	42·4	42.5.2 Hidden Operating Systems	42·14
42.3 HARDWARE-BASED ANTIPIRACY TECHNIQUES	42·5	42.6 POLITICAL AND TECHNICAL OPPOSITION TO DRM	42·15
42.3.1 Dongles	42·5	42.6.1 Political Opposition	42·15
42.3.2 Specialized Readers	42·6	42.6.2 Technical Countermeasures	42·17
42.3.3 Evanescent Media	42·10		
42.3.4 Software Keys	42·11		
42.4 DIGITAL RIGHTS MANAGEMENT	42·12	42.7 FUNDAMENTAL PROBLEMS	42·17
42.4.1 Purpose	42·12	42.8 SUMMARY	42·18
		42.9 GLOSSARY	42·19
		42.10 FURTHER READING	42·22
		42.11 NOTES	42·22

42.1 INTRODUCTION. Ever since publishing and commerce were introduced to the digital world, the risks to intellectual property and to personal privacy in cyberspace have steadily escalated on comparable but separate paths. These paths have now converged. Unfortunately, many times, antipiracy efforts lead to possible breaches in personal privacy.

Efforts to stem the flow of pirated software worldwide remain mediocre in efficacy; piracy is still proving to be big business in the new millennium. According to the Business Software Alliance (BSA), the “global piracy rate hovered at 42 percent in 2011 while a steadily expanding marketplace in the developing world drove the commercial

42 · 2 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

value of software theft to \$63.4 billion.” Further, “these startling findings come from a survey of approximately 15,000 computer users in 33 countries that together make up 82 percent of the global PC market.”¹ Continuing piracy means lost jobs, wages, tax revenues, and a potential barrier to success for software start-ups around the globe.

At the same time, freedoms inherent in the Internet have made maintaining privacy of personal information a true challenge. Identity theft keeps rising, as more and more companies actively engage in the accumulation of customer data through e-commerce. By paying bills, checking medical insurance accounts, or completing taxes online, people are making their personal information available to be stockpiled, shared, and regurgitated to the point that simply “Googling” one’s identity can be a real eye opener. Technologies that aim to prevent piracy have the potential to use this wealth of available personal information in a way that significantly erodes personal privacy. In particular, some of these technologies send personally identifiable information to servers on a routine basis (see, e.g., Section 42.4.2). The idea that a corporation—or a government—could be scanning what a specific person reads, listens to, or views is grounds for concern to civil libertarians.

42.1.1 Digital Rights. In this environment of rapid change in both piracy and privacy, even the term *digital rights* is ambiguous. When software companies and music producers talk about digital rights, they mean the kinds of rights long protected by copyright, trademark, and patent law. When privacy advocates argue about digital rights, however, they may be talking about a completely different thing: that an individual does not forfeit personal rights, including the right to privacy, merely by turning on a computer.

This chapter’s primary focus is on technologies designed to protect traditional rights of content producers, but it also enumerates areas where those technologies threaten personal privacy.

42.1.2 Patent, Copyright, and Trademark Laws. There are differences among the applicable laws and the materials they protect.

Patents give owners exclusive rights to use and license their ideas and materials; patents generally protect nonobvious inventions in mechanical and electrical fields, as well as those that can be embodied in computer software and hardware.

Copyrights give owners the exclusive rights to create derivative works, to reproduce original works, and to display, distribute, and conduct their works. Copyrights apply to original works of authorship including paintings, photographs, drawings, writings, music, videos, computer software, and any other works that are fixed in a tangible medium. Copyrights, their infringement, and remedies are described in the Copyright Act of 1976.

Trademarks give owners the right to restrict the use of distinctive marks in certain contexts. These rights may apply to words, sounds, distinctive colors, symbols, and designs.

For an extensive discussion of intellectual property law, see Chapter 11 in this *Handbook*.

42.1.3 Piracy. Once thought of as a mere copyright infringement of printed matter or production of a counterfeit audiotape, piracy has grown with technology and has expanded to encompass intellectual property, digital data, DVDs, CDs, VHS, analog and high-definition TV, and streaming media.

SOFTWARE-BASED ANTIPIRACY TECHNIQUES 42 · 3

There are several types of piracy. End user piracy occurs when end users use a single copy of software to run on several different systems, or when they distribute copies of software to others without permission of the software manufacturer. Reseller piracy occurs when unscrupulous resellers distribute multiple copies of a single software package to multiple customers, preload the same software on multiple systems, or knowingly sell counterfeit software to customers. Internet and bulletin board (BBS) piracy occurs when users download and upload copyrighted materials and use it or make it available for use by others without proper licenses.

To understand why and how piracy occurs and the enormous impact on society worldwide, we need to have a clear understanding of what we mean by the word “piracy.” Whenever information is created and published in print, on the Internet, or incorporated into software, that information may be protected by copyright, patent, or trademark law. This principle applies to a broad spectrum of material that includes, for example, the Wright Brothers’ specifications for their “Flying Machine”; Microsoft Windows software; the icon Mickey Mouse and all related materials; and television shows, plays, movies, and music created and performed live and on recordings. Making unauthorized copies of such material in any medium is referred to as *piracy*.

Within the last three years, the Software & Information Industry Association² has “filed more than 100 lawsuits in the U.S. against illegal eBay sellers as well as sellers on other websites dealing in counterfeit, OEM, academic, region-specific, and other illegal software and publications. Defendants have paid millions of dollars in damages, and, in some cases, criminal charges were pursued and defendants sentenced to jail time.”³

42.1.4 Privacy. As in many aspects of security, end users are now being called on to protect their online (and offline) identities through diligent monitoring of financial activities and through awareness programs focused on personal rights, laws, and obligations of Internet use. The average Web user today can create and publish a blog or personal video to the Web faster than that same user can apply software patches to a desktop PC. So widespread is this new facility that the law has yet to catch up with the needs and expectations of the users in online society. It is difficult to keep data private when—with a mouse click—it is easy to share personal information with the whole world.

Beyond bad personal habits that reduce privacy, however, a whole new class of applications termed *digital rights management* (DRM) collects more personal information than ever before in an effort to reduce improper use of copyrighted material. DRM products may record and report on an individual’s Web-browsing habits, types of files created and accessed by a particular program, number of uses of a particular file or program, source IP address of the user’s system, and presence (or absence) of a license for a program. In the name of protecting digital rights for content producers, the consumers of this content are being cataloged and tracked in a way that the framers of copyright laws would have been hard-pressed to imagine.

For extensive coverage of privacy issues on the Internet, see Chapter 69 in this *Handbook*.

For a glossary of terminology used in discussing digital rights management, see Section 42.9.

42.2 SOFTWARE-BASED ANTIPIRACY TECHNIQUES. A variety of software-based technical approaches are used to prevent inappropriate use of copyrighted and otherwise protected material in an organization’s networks and on the

42 · 4 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

public Internet. Current methods include ensuring proper configuration of operating systems, monitoring of installed software, encryption of content, and insertion of some sort of key or identifier in the digital product itself.

42.2.1 Organizational Policy. Controls in free-standing commercial applications represent only one of the technical means to protect digital content. Existing system controls are also useful in this regard. Operating system access controls can specify who can access particular content, while encryption supported by the OS and other applications can limit access to users who possess the appropriate key. More generally, organizational policy should specify good practices in configuring operating systems and applications in order to help protect content. Such policy includes:

- Allow users to install only software that is necessary.
- Encrypt information that should not be publicly viewable.
- Install software with the lowest possible privilege consistent with the ability to do its job.
- Disable active content (Java, JavaScript, ActiveX, cookies, etc.) wherever feasible.
- Use network operating system access controls to limit access to shared, copyrighted media to members of the organization for whom licenses are purchased.

42.2.2 Software Usage Counters. Software metering has been popular for several years. Special software monitors system usage and inventories the software on the system or network. This type of software also can be used to block or limit the use of specific software, such as browsers and games. In addition to fighting piracy, it can reduce the load on IT personnel by reducing complications due to use of unauthorized software.

42.2.2.1 Controlling Concurrent Installations. Software metering products can monitor concurrent installations even on networks used by people in different geographic areas who have different requirements and different software installed. The metering software permits an administrator to maintain a live and updated inventory of the software installed at different locations on the network. The logs show where installation has taken place, when the licenses expire, and when updates are necessary. Alerts can be set to notify system administrators when a license is about to expire or when an update has been accomplished.

42.2.2.2 Controlling Concurrent Usage. Software metering allows network administrators to identify and resolve cases of illegal installation of unauthorized copies of authorized software and also to catch people who install unauthorized software on the organization's computers. Metering software also allows a company to report and analyze logon and logout times, track software usage, and meter software licenses to keep those in the company legal. In addition to avoiding legal entanglements, monitoring can reduce the demand on system resources, network bandwidth, and technical support staff.

42.2.2.3 Examples and Implementation. Microsoft announced in 2000 that to combat piracy, new releases of the Office 2000 program would include a counting feature making the programs malfunction if the owner had not registered the software after launching it 50 times.⁴ Prior to this announcement, Microsoft published

HARDWARE-BASED ANTIPIRACY TECHNIQUES 42 · 5

antipiracy literature and provided a great deal of consumer education regarding the effects of software piracy on society; it established, as well, a piracy hotline (1-800-RU-LEGIT). Novell (1-800-PIRATES), Adobe, and Xerox are other companies that have vigorous antipiracy programs in place, although they have not yet announced that they are building metering into their products.

Metering software requires a bona fide license in order for it to be implemented. Typically, companies establish CD-ROM keys that are printed on legitimate copies of their product installation discs or jewel cases. The keys include checksums or message authentication codes that can be checked by the installation routines. The algorithms for the checksums are intended to cause difficulty for people trying to create counterfeit keys. The security of such measures depends on the cryptographic strength of the validation keys.

Software counters for controlling concurrent usage need to store information securely so that each load operation increments a counter and each unload operation decrements it. However, the security problem is to store this information in such a way that unauthorized people cannot modify it easily. Encrypting the data in a complex sequence of operations can stop most abuses of the system by making the effort required for circumventing the mechanisms more costly than buying a license.

More recently, starting with Vista, copies of the Windows operating system must be registered (Microsoft calls this *activating* a copy) soon after installation, or they will go into *reduced functionality mode*, in which basic functions are available for only an hour of operation before logging on again.⁵ For single-copy home use, Windows relies on a license key encoded on the installation medium. For enterprise deployments of Vista, an organization needs to run key servers to allow the registration to occur.⁶

42.3 HARDWARE-BASED ANTIPIRACY TECHNIQUES. Working on the theory that software is prone to misconfiguration and compromise, antipiracy groups and researchers have experimented with a variety of hardware-based approaches to preventing inappropriate use of protected content. These techniques include dongles and specialized readers attached to the reading hardware, evanescent media designed for viewing or playing only a limited number of times, and software keys incorporated into the media and the reading hardware.

42.3.1 Dongles. Dongles are hardware lock devices or modules that connect to a computer and communicate with software running on the computer. Without a dongle in place, the external device or regulated software does not work fully, or at all.

Initially, dongles controlled printing. With a dongle installed on the computer, no one could print data from the computer without authorization. However, there is now a necessity for protecting all types of devices. Now dongles are used to protect scanners; external drives (e.g., ZIP drives); CD-ROMs and rewritable CD-ROMs; DVDs and DVD-Rs; VHS recorders; PlayStation, Nintendo, and Sega video gaming systems; and even personal digital assistants (PDAs).

The most common type of dongle provides a pass-through port to connect a device cable. Generally, a dongle incorporates some type of algorithmic encryption in its on-board microelectronic circuitry. The sophistication of the encryption varies depending on the manufacturer and the device. Many dongles provide additional onboard non-volatile memory for the software to access. Some models even have real-time checks that keep track of date and time information, including when an application's license (temporary or leased) is set to expire.

42 · 6 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

Dongles provide some definite advantages:

- Because a dongle is an external device, it is fairly simple to install and uninstall. Early dongles used serial or parallel ports, but USB has become the norm in recent years. In most cases, since manufacturers support their devices, an ordinary user can install and use a dongle without help from an IT department.
- Dongles also require registration, which provides adequate control over the use of the dongle and thus provides legitimacy to both the device and the users. Registration (dependent on the contract in place) may provide support for both the software and the hardware.
- Dongles that support encryption provide an extra layer of protection by making transmitted data indecipherable until it reaches its destination, unless the hardware is in place.

There are also disadvantages to using dongles:

- Consumers resist the requirement for installation, maintenance, and additional cost. Most large corporations do not use dongles for their products.
- Dongles can be lost or stolen, and they also may fail.
- Sometimes a dongle will work well with a slow computer but cause errors when installed on a faster computer.
- Since not every manufacturer automatically replaces lost or stolen dongles without charge, there may be additional costs involved in getting replacements.
- Dongles can present a serious risk-management problem for critical applications where delays in obtaining replacements or registering them may be unacceptable.
- As with any device, there can be a serious problem if the dongle manufacturer ceases to support the model of dongle a company has installed or if the manufacturer goes out of business entirely.
- Laws regarding encryption usage differ in various countries. Specialized dongles that may be legal to use in the United States may be illegal in another country.

42.3.2 Specialized Readers. One of the impediments to illegal copying used to be the difficulty and cost of obtaining specialized hardware and software for reading and copying proprietary materials with fidelity. However, today such copying equipment is inexpensive and easy to find. In addition, the media for distributing illegal copies are less expensive than ever.

42.3.2.1 Audio. According to the Recording Industry Association of America (RIAA), the global audio industry loses in excess of \$4 billion every year to piracy worldwide.⁷ RIAA says that \$1 million a day, in just physical product, is lost in the United States alone. The loss of ancillary revenues drives the figures higher. But the RIAA claims that these figures are low, since it estimates that in some countries up to 98 percent of the music in use comes from illegal copies.

As part of an industry-wide, organized approach to highlighting and reducing the music piracy problem, the RIAA has taken a very active role in pursuing legal action against suspected pirates. Through the 1990s and continuing into the current decade, the biggest problem with audio piracy was illegally copied CDs. In 1998, for example, the RIAA confiscated 23,858 illegal CDs in the first half of the year. In that same year,

HARDWARE-BASED ANTIPIRACY TECHNIQUES 42 · 7

Operation Copycat—a joint investigation by RIAA, the Motion Picture Association of America (MPAA), and the New York Police Department—saw the arrest of 43 CD pirates and the shutdown of 15 illegal manufacturing locations. Many of the CDs seized in these types of operations apparently came from Asia and Eastern Europe, financed by organized crime operations with ties to drugs and prostitution.⁸ By 2002, even corner convenience stores were contributing to the problem, providing coin-operated CD copying machines akin to photocopiers, along with the familiar posted warning transferring liability to the user.⁹ Given the enormous profits involved, the problem has been simply too large and widespread for law enforcement to control. In 2005, the RIAA seized approximately 5 million illegal CDs.¹⁰

More recently, the RIAA has focused on the problem of music files illegally shared across the Internet. Using free software, sometimes already bundled into commercial operating systems, anyone can download music tracks and burn CDs. The MP3 music file format has become a ubiquitous way to share music with others. The RIAA originally protested against MP3 players, but the phenomenal success of personal digital music players, particularly the Apple iPod, has rendered such efforts futile. Some musicians and independent record labels have adopted the MP3 format to promote their records, showing that the technology itself has no inherent ties to piracy. These musicians and recording studios claim that they are happy with consumers downloading the music, and they are at odds with the RIAA.

Some musical groups have even experimented with severing ties with the traditional music industry altogether, using Websites and social networking services to advertise and distribute digital music files directly to their listeners.¹¹ As more musicians begin to use the Internet either in addition to or in lieu of traditional distribution models, the very model of music distribution is in flux.

Meanwhile, the music industry continues to face serious financial loss due to illegal downloading. The industry, through the RIAA, has been pursuing legal remedies. Some major lawsuits have achieved significant press coverage and have been instrumental in enforcing or changing existing laws or in helping develop new laws. For instance, deliberate copyright violations resulted in an award of \$50 million in statutory damages and \$3.4 million in legal fees to Universal Studios in a suit against the MyMP3.com music service. MyMP3 created a database of over 80,000 albums, which, when combined with the MyMP3 software, let users access and store music digitally, without paying a fee.

Perhaps the most recognizable name in the music field in regard to piracy at one time was Napster, a site that enabled individuals to share tracks of music via the Internet. The site provided free downloadable software for downloading and playing MP3 files. Essentially, Napster software turned a user's PC into part of a distributed server network that published available music files. It did not take long for the site to acquire a user group of millions of people who after "sampling" the music might then go to a store and buy the entire CD.

However, since the Napster site did not limit the length of the download, many users simply downloaded the entire track. Most never bought the commercial version of the music. Adding to successful piracy attempts was the development and ready availability of rewritable CD drives. More and more Napster users decided to download the music tracks they wanted and then burn their own CDs without ever purchasing the CDs made by the recording artists and music companies.

Creating a stir in the industry and ultimately a landmark judicial case, Napster was forced to radically alter operations in March 2000 after protracted court proceedings. Upon the verdict, Jack Valenti, president and chief executive officer of MPAA,

42 · 8 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

commented that the consumer would benefit most from the court's decision because "You cannot take for free what belongs to someone else."¹² But the subject of Napster and audio piracy remains highly controversial. Although some people argued that little-known artists received exposure that they may never have gotten without the free file-sharing service, others, especially large music companies and recording artists, argue that they were being denied the royalties they deserve.

Napster attempted to re-create itself as a pay-for-subscription music download service, but found record labels unwilling to work with it. In 2002, Napster folded; the name was eventually purchased by Roxio, Inc. to rebrand its own subscription service. In the meantime, other peer-to-peer (P2P) file-sharing protocols and applications appeared to fill the void left by Napster. By the time the iTunes Music Store emerged as a powerful contender with music company backing and with copyright protection, the world of free file sharing was reinvigorated by names such as Gnutella, FastTrack, Grokster, Limewire, and Kazaa.

Over the last few years, the music industry has identified universities as fertile ground for pursuing illegal downloading activities. In 2007, the RIAA launched a new round of attempts to bring music pirates to justice, sending letters offering to settle with students identified as probably sharing copyrighted files, in advance of any trial.¹³ The RIAA's tactics have been raising hackles in the higher education community, with opponents criticizing the letters as bordering on extortion. Some universities are refusing to forward the letters to students; one university agreed to forward the letters, but promised to bill the RIAA \$11 for every letter to pay for its staff time.¹⁴

42.3.2.2 Video. On the video side, Scour, Inc., provided free downloads of digital movies as well as software allowing users to share the downloaded files among themselves without the use of a central server. With an easy-to-use interface and quick response time, Scour.com became quite popular in a short period of time. Launched in 1997, with a Web search feature added in 1998 and a subsequent P2P tool, Scour eventually attracted negative attention from the movie and music industries. In July 2000, MPAA, RIAA, and the National Music Publishers Association (NMPA) sued Scour, accusing it of large-scale theft of copyrighted material and of trafficking in stolen works. By November 2000, the company was out of business.¹⁵

This case did not stop the counterfeiting of video media. Despite increasingly steep fines, judicial rulings, and even raids by various law enforcement agencies, counterfeit video is readily available. Along Fifth Avenue in New York City, for \$5 to \$10 anyone can buy the latest films and DVDs; in markets in Hong Kong, Southeast Asia, and India, the copies are even cheaper. It is true that some of the copies available may have been "legally produced," but it is more than likely that the counterfeit or bootleg copies were made illegally from a master copy that was either borrowed or stolen.

Advances in consumer electronics help the trend, as many PCs now come with a recordable/rewritable DVD player as a standard component. Arguments about legality of time-shifting and space-shifting that once defended the practice of making personal mixes on audiocassette have now moved to the realm of digital video.

42.3.2.3 Television (Analog). Broadcast television has been one of the most successful technologies in history, and financial interests are still huge, despite the growth of cable and satellite services. In January 2000, major television companies, the National Football League, and the National Basketball League all filed complaints against iCraveTV, a Canadian company that had been in existence for only a year.

HARDWARE-BASED ANTIPIRACY TECHNIQUES 42 · 9

According to the complaints, iCraveTV was illegally using broadcast television signals without authorization or payment and streaming the signals to the iCrave Internet site for viewing free of charge. Although this practice apparently did not violate Canadian copyright laws at the time, U.S. judges ruled in February 2000 that the unauthorized transmissions of broadcast signals into the United States via the Internet were a direct violation of U.S. copyright law, and iCraveTV was ordered to stop the practice. Shortly after iCraveTV agreed to an out-of-court settlement, the Web-site was shut down and iCraveTV went out of business.¹⁶

Hacking cable decoders is another technique for obtaining services without paying for them. Although it is not illegal to buy, install, or modify equipment for converting encoded cable TV signals from pay-per-view or other commercial suppliers, it is illegal to use such set-top decoders to obtain services without paying for them.

In the United States, Congress mandated that after February 17, 2008, all TV stations must transmit in digital format (DTV) only.

42.3.2.4 Television (HDTV). The first television image was created in 1884 when Paul Nipkow created a mechanical scanning disk. With only 18 lines of resolution, the picture was poor. Current National Television System Committee (NTSC) standard TV transmissions are done with bandwidth that does not exceed 6 MHz. The current analog system broadcasts 30 frames per second and 525 lines per frame.

High-definition television (HDTV) is a digital television system that offers twice the horizontal and vertical resolution of the current TV system. HDTV has the ability to deliver a video composed of approximately 1,125 lines per frame and 60 frames per second. Viewers then see a picture quality close to that of 35-mm film. Obviously, transmitting images containing that large amount of audio and video information requires wide bandwidth, actually about 18 MHz. Such bandwidth would permit the transmission of 1,050 lines of 600 pixels per line. However, the Federal Communications Commission (FCC) decided to limit HDTV to a 6-MHz maximum bandwidth. In order to meet that requirement, MPEG compression would be used.

MPEG compression applies algorithms to pixel groups and records information that changes within a frame, rather than all of the information in all of the frames. Audio is synchronized to the video. Using MPEG saves storage space and transmission requirements while retaining high image and sound quality. According to the Advanced Television Committee Standard (ASTC), the FCC will require that audio and video compression as well as the transmission of HDTV terrestrial signals follow this standard.

As with all other transmissions and media, there are serious concerns about piracy of HDTV transmissions and programs. At the present time, even though many TV transmissions are scrambled in order to thwart reception, it is fairly simple, although illegal, to purchase a descrambler and unscramble the transmissions. It is, however, legal for home viewers to record programs for their own personal use.

The HDTV market space has been evolving, and consumer demand for HD-capable devices has exploded over the past several years. Attempts by U.S. television producers to protect themselves and their content using a variety of scrambling and encryption schemes, including content scrambling systems (CSSs), have been made difficult by the frequent changes in hardware and signal formatting that have accompanied this rapid market expansion.

Encrypting terrestrial broadcast television programming would secure the transmissions, but according to the Home Recording Rights Coalition (HRRC), such encryption will threaten established home recording rights. The HRRC contends that Section 1202

42 · 10 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

(k) of the Digital Millennium Copyright Act provides a carefully balanced approach to analog home recording rights and stipulates that mandated technology may not be applied to interfere with consumer recording of free, over-the-air terrestrial broadcasts.

Furthermore, the HRRC contends that encrypting the free television broadcast content will create very little incentive for consumers to switch from regular analog to digital television. Instead of thwarting digital pirates, the HRRC contends that strong encryption will impose unfair and even illegal restrictions on consumers.

42.3.2.5 Consumer Acceptance of Specialized Readers. Illegal sharing of copyrighted content is common across the Internet. When the Software Publishing Association (SPA) was first formed, the group, together with law enforcement agencies, raided the physical premises of companies believed to be using pirated software. As a result of finding quantities of the pirated software, SPA won many legal actions and related settlements.

Some people see encryption as a challenge and work at breaking the algorithms so they can pirate the data—digital, video, or audio. In addition, the lack of standardization of laws throughout industries and countries has led to controversy and ongoing piracy.

Although average consumers do not think of themselves as intellectual property pirates, many otherwise honest citizens do get and use illegal programs, applications, games, audio tracks, CDs, DVDs, VHS tapes, and television signals. This situation might be attributed to a lack of ethical education, but many people like to save money and simply do not believe they will be caught and punished for such pilfering. A 2001 study by the Pew Internet & American Life Project, based on phone interviews with 4,205 adults 18 and over, some 2,299 of whom were Internet users, suggested that around 30 million U.S. residents had downloaded music from the Internet.¹⁷ At that time, the phrase “had downloaded music from the Internet” was basically equivalent to “had *illegally* downloaded music from the Internet,” since legal means of doing so had yet to evolve.

Since that report, Apple’s 2001 announcement of its iTunes and iPod products, and the 2003 launch of the iTunes Music Store, began a movement to provide consumer-friendly ways of downloading music that also give compensation to copyright owners. iTunes uses device authentication and proprietary encoding formats to limit redistribution of downloaded songs and videos. Despite—or perhaps because of—Apple’s attempts to comply with copyright laws, it is clear that Apple filled a perceived need in the market, as it has generated both a host of competitors and substantial sales. Consumers downloaded the first million songs from the iTunes Store in five days, and the overall market for digital music has grown to at least \$790 million per year. At \$0.99 per song, downloads from the iTunes Music Store passed the \$1 billion mark on January 23, 2006.¹⁸

42.3.3 Evanescent Media. There are many interpretations of the term *evanescent media*. The broad interpretation includes digital imaging, optics, multimedia and other electronic art, and data that are short-lived or transitory. When such media are original, creative works, society has an interest in protecting them against piracy.

Since most evanescent media involve some visual aspects as well as text, antipiracy techniques now being used or considered for other types of data may be applicable. Such techniques include previously discussed dongles, software keys, watermarks, encryption, and digital rights management. Part of the problem in electing and implementing a solution is the lack of existing standards that specifically deal with this new area of art and science.

HARDWARE-BASED ANTIPIRACY TECHNIQUES 42 · 11

42.3.4 Software Keys. Software keys of various kinds are used to secure data and equipment. A software key is generally a string of numbers that is used for identification purposes, either to allow access to the use of equipment or to permit authorized printing, processing, or copying of data. As described earlier in the discussion of dongles, most antlicing hardware devices are accompanied by software that works in tandem with the hardware. A software key activates or deactivates the hardware lock. When the software is working perfectly, there are generally no difficulties. However, all software can malfunction, and when that happens, there can be serious problems in getting equipment to work. Additional problems occur when the computer containing the software key malfunctions, and the software key cannot be made to work on a replacement machine.

42.3.4.1 Videocassettes versus Copy Machines. Watermarking is one of the techniques being seriously considered for protecting videocassettes and DVDs. In 1995, the ASTC formed a Copyright Protection Technical Working Group, which spun off a special Watermarking and Embedded Data Encoding subgroup. The group has broad representation, including representatives from the PC market, the Macintosh market, the MPAA, the Consumer Electronics Manufacturers Association (CEMA) and related manufacturers, technicians, and users. Their task is to look for technologies and services that might use hidden data clues as a means of inhibiting or barring digital piracy. Using a hidden watermark that can be embedded in the content would then prevent machines from making copies or would alert the operator that the videocassette is marked and that unauthorized copies would be considered pirated.

42.3.4.2 DVD Area Encoding. Digital video requires very large storage space—too large for a single CD to hold. However, by applying compression techniques, the digital video can be compressed to fit into the digital videodisc's maximum capacity of 17 gigabytes. Two different types of compression are used for encoding audio and video content for DVD: constant bit rate (CBR) and variable bit rate (VBR) compression.

In order to prevent piracy of the content of the DVD, many companies are turning to encryption. The compressed data are encapsulated through a mathematical algorithm that can be decrypted only through the use of a decryption key.

42.3.4.3 Implementation. For shorter programs, CBR is ideal. Based on MPEG-2 encoding, CBR compresses each frame of audio and video by a user-selected amount. This degree of compression is then applied to the entire program. Using VBR, it is possible to create a database of video content based on the amount of change in each frame or scene. This is particularly useful in programs with a long format. To construct the database, the encoding software does several analytical passes of the master footage and then makes a final digitizing pass. From the created database, the computer can encode the video with a variable data rate, allowing a higher bit rate for scenes with pans, zooms, and fast motion and giving scenes with little or no motion low data rates. By greatly compressing the areas of lower detail, areas of higher details can be allocated more space and use less compression.

42.3.4.4 Watermarks. Watermarking involves embedding one set of data inside a larger set of data. The embedded set of data identifies the origins or ownership of a specific work, just as a watermark does on paper.

42 · 12 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

Using digital watermarks can help copyright owners track the use of anything digital, including music, movies, photographs, and clip art. Digital watermarking is widely used for protecting images. For instance, photographers often post low-resolution (low-res) versions of their photos on public Websites and use visible digital watermarks to clearly label the low-res images as copyrighted. Upon payment of the appropriate fee, the customer receives the high-resolution version of the photo, presumably with at least any visible watermarks removed. The use of invisible watermarks to prevent undetected sharing after purchase of digital content is more controversial and more prone to questions about reliability of detection; for instance, how many false positives and false negatives will occur? Additionally, there is the question of survivability of the mark itself as it is run through various transformations.

The music industry flirted with digital watermarking to protect music files beginning in 1998 with the formation of the Secure Digital Media Initiative (SDMI), a consortium of technology, security, and music organizations. The SDMI developed several watermarking schemes, and in 2000, it offered a reward to anyone who could crack the code and remove the watermark from a song protected by SDMI's technologies. The Electronic Frontier Foundation asked the Internet community to boycott the contest, stressing that the use of DMAT (Digital Music Access Technology) would mean that manufacturers and users would be forced to adopt the DMAT format in equipment and would create additional costs for manufacturers and consumers. A team of researchers, led by Princeton professor Ed Felten, was able to remove the invisible watermarks. When Felten attempted to publish the results of his process, attorneys for SDMI threatened to sue him under the Digital Millennium Copyright Act (DMCA). SDMI never filed suit, but Felten himself sued for a declaratory judgment to clarify the matter. Felten's suit was dismissed by a federal judge, but not before the government and the RIAA agreed that researchers should not be punished under the DMCA for testing technologies to protect copyright.¹⁹ The SDMI has been inactive since 2001.²⁰

42.4 DIGITAL RIGHTS MANAGEMENT. Recognizing that piracy is a huge moral and financial problem, software developers have adopted and modified another type of system that can be applied to print, audio, video, and streaming media. Called *Digital Rights Management* (DRM), the system was originally devised to protect proprietary information and military information. The idea behind the system is to protect all types of intellectual digital content from anyone who would take it without the consent of the developer(s) or owners. Major companies like Microsoft, Adobe, and IBM are developing and marketing DRM systems, and dozens of smaller companies are springing up.

42.4.1 Purpose. The purpose of DRM is to protect all digital content that originators or owners want protected. DRM permits distributors of electronic content to control viewing access to that content. The content can be text, print, music, or images. Basically, DRM systems use a form of customized encryption. When an end user purchases viewing, listening, or printing rights, an individual "key" is provided. The system works on rules, meaning that although a key is provided, it generally comes with limitations regarding the copying, printing, and redistribution.

Unfortunately, there is no agreement on a DRM solution. Lack of standards is hampering businesses from moving forward with online business initiatives. Because there are so many companies promoting their own incompatible forms of DRM, customers will have to download megabytes of code for each version. Maintaining, upgrading, and managing all of those different versions are major headaches for customers. There

DIGITAL RIGHTS MANAGEMENT 42 · 13

does not appear to be a simple solution; rather than being a technology issue, it is really a matter of business and politics.

42.4.2 Application. Typically, when users become prospective owners of digital rights, they download a content file. The DRM software does an identity check of users, contacts a financial clearinghouse to arrange for the payment to be made, and then decrypts the requested file and assigns users a key. The key is used for future access to the content.

Because the system works on rules, it is possible to impose restrictions. One user might pay just to view material, while another user might want to have printing privileges. A third user might want to download the content to his or her own machine, and, finally, a fourth user might want to have viewing privileges for a specified time. The four different authorized users would thus use the same content, and each would pay according to a rate scale established by the content distributor. Throughout all of the transactions, each user would need a mechanism that allows secure transmissions and identifies that user and the associated level of access privileges.

Although this approach to publishing may sound fairly simple, it is really quite complex. In addition to arranging for different users to access material according to the set rules and to pay according to a rate schedule, it is also necessary for content distributors to handle the back end of the application. Everyone involved in the creation, production, and distribution of the content has to be paid fairly for the use of the content.

Payment is especially important as more and more content providers digitize materials that they can show or print on demand. Many users will read books online, but some physical printing on paper will continue. However, publishers will be able to print precisely those volumes that are requested. This approach will provide customized printing (e.g., large-print editions) as well as saving paper and physical warehouse storage space.

42.4.3 Examples. Several different types of DRM systems exist. Experts agree that the best DRM systems combine both hardware and software access mechanisms. With the advent of the eBook, the digital pad, PDA modems, Internet access devices, and increasingly smaller laptop computers, tying access rights directly to storage media gives publishers and distributors control of where the content is being used as well as by whom.

With the passage of the *Electronic Signatures in Global and National Commerce Act*, referred to as the *E-Sign Bill* and the increasing use of digital signatures, original documents (e.g., legal, medical, or financial) will be stored digitally. President Clinton signed the E-Sign Bill on June 30, 2000, in Philadelphia's Congress Hall using a ceremonial pen and a digital smart card. The bill went into effect on October 1, 2000. The E-Sign Bill gives an online signature the same legal status as a signature etched on paper and makes the digital document the original. Any printout will be considered a copy, so the ability to view documents and videos (e.g., living wills) digitally will actually give the viewer access to the original. Eventually, records for medical treatment and documents for trials may be totally digital and may require submission and viewing digitally. When this becomes the norm rather than the exception, strict adherence to DRM in order to maintain privacy, as well as to provide restitution, will be paramount. In addition, such content-protection schemes will prevent unauthorized modifications of the digital data that would otherwise contribute to fraud.

For example, IBM has released antipiracy technology called the *Electronic Media Management System* that allows for downloading music tracks but puts controls on

42 · 14 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

how many copies can be made or allows for a copy length limitation to be inserted. Thus, a minute of music could be downloaded to give a listener a taste, but not a chance to pirate the entire music track. To obtain the entire track, the user would be required to pay a fee.

Microsoft distributes free software that embeds metatags in each audio file. The metatags refer back to a central server in which the business rules are stored. This approach requires that material be tagged as it is created; otherwise, if it is released without the embedded tags, it can be illegally copied.

Major companies like Xerox, Microsoft, IBM, and Adobe got heavily involved producing and using this software in the 1990s, and many smaller firms opened shop. As with other new technology launches, eventually many of the small entrants went out of business or were bought by larger firms. Some of the small companies, such as ContentGuard,²¹ continue to exist independently, as of this writing.

42.5 PRIVACY-ENHANCING TECHNOLOGIES. Although DRM may seem to be a valid solution for the piracy problem, dissenters feel that DRM and other antipiracy measures give producers and distributors too much control. The rationale is that excessively restrictive rights management may undermine the fair use rights of consumers and academics. Partly as a result of these conflicting viewpoints, many consumers are making increasing use of privacy-enhancing technologies (PET), a broad term for a range of technologies designed to hide the identity and activities of individual users and computers as their traffic traverses the Internet.

42.5.1 Network Proxy. One broad class of tactics and tools used to enhance privacy is based on the concept of a network proxy. In its most general form, a proxy takes a network connection request from a client and redirects it to the ultimate destination, changing the address headers to make it look like the original request came from the proxy itself. When the destination server responds, the proxy returns the results to the requesting client. Proxies have long been used within organizations, both to protect internal users as they make requests to untrusted networks and to track and sometimes block access to undesirable content. For more details on the use of proxies in Web content filtering and monitoring, see Chapter 31 in this *Handbook*.

More recently, proxies have been employed outside the bounds of corporate networks to allow anonymous connections across the Internet. These so-called anonymizing proxies sometimes use encryption to hide the traffic in transit and so also provide protection from traffic analysis. Anonymizing proxies are a serious threat to organizations that desire (or are required by law) to monitor and block users from accessing certain kinds of information. More advanced versions of this concept use multiple routers to hide the path that a request takes through the public network. As a general class, these are known as mixing networks. One example described as early as 1981 is the Chaum Mix.²² Recently, a concept known as onion routing has become popular in its incarnation as Tor (the onion router), which uses nested layers of traffic encryption as a session travels from one router to the next.²³ See Chapter 70 in this *Handbook* for further information about anonymity on the Internet.

42.5.2 Hidden Operating Systems. Rather than relying on network technologies, some users are choosing to make their actions on the network private by using hidden operating systems. Two basic approaches are common: the virtual machine and the bootable system.

POLITICAL AND TECHNICAL OPPOSITION TO DRM 42 · 15

A virtual machine is a system that runs within another operating system. Long used for cross-platform compatibility and the ability to run multiple systems on a single hardware platform, virtual machines, such as Java Virtual Machine, VirtualPC, and VMware, are also used to hide activity from those who would disapprove of it (system administrators, parents, law enforcement, etc.), since the activities of the virtual machine can be made invisible to the host machine.

The bootable system approach, however, stores an entire operating system on some sort of bootable medium, such as a CD or USB device. If a computer can boot from such media and store downloaded content to a peripheral device rather than the host operating system's hard drive, then no record of the usage will remain when the host is next booted up.²⁴

42.6 POLITICAL AND TECHNICAL OPPOSITION TO DRM. Vocal opponents of DRM have organized using the Web to exert pressure on vendors using the techniques; criminal hackers and others have developed and circulated software for disabling DRM.

42.6.1 Political Opposition. The *Electronic Frontier Foundation* (EFF) is a widely respected organization, which describes itself as follows:

From the Internet to the iPod, technologies are transforming our society and empowering us as speakers, citizens, creators, and consumers. When our freedoms in the networked world come under attack, the Electronic Frontier Foundation (EFF) is the first line of defense. EFF broke new ground when it was founded in 1990—well before the Internet was on most people's radar—and continues to confront cutting-edge issues defending free speech, privacy, innovation, and consumer rights today. From the beginning, EFF has championed the public interest in every critical battle affecting digital rights.

Blending the expertise of lawyers, policy analysts, activists, and technologists, EFF achieves significant victories on behalf of consumers and the general public. EFF fights for freedom primarily in the courts, bringing and defending lawsuits even when that means taking on the US government or large corporations. By mobilizing more than 140,000 concerned citizens through our Action Center, EFF beats back bad legislation. In addition to advising policymakers, EFF educates the press and public.²⁵

The EFF has steadfastly argued against DRM:

Digital Rights Management (DRM) technologies attempt to control what you can and can't do with the media and hardware you've purchased.

- Bought an ebook from Amazon but can't read it on your ebook reader of choice? That's DRM.
- Bought a DVD or Blu-Ray but can't copy the video onto your portable media player? That's DRM.
- Bought a video-game but can't play it today because the manufacturer's "authentication servers" are off-line? That's DRM.
- Bought a smart-phone but can't use the applications or the service provider you want on it? That's DRM.

Corporations claim that DRM is necessary to fight copyright infringement online and keep consumers safe from viruses. But there's no evidence that DRM helps fight either of those. Instead DRM helps big business stifle innovation and competition by making it easy to quash "unauthorized" uses of media and technology.

42 · 16 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

DRM has proliferated thanks to the Digital Millennium Copyright Act of 1998 (DMCA) which sought to outlaw any attempt to bypass DRM.

Fans shouldn't be treated like criminals and companies shouldn't get an automatic veto over user choice and innovation. EFF has led the effort to free the iPhone and other smart phones, is working to uncover and explain the restrictions around new hardware and software, has fought for the right to make copies of DVDs, and sued Sony-BMG for their "rootkit" CD copy protection scheme.²⁶

The Free Software Foundation created the *Defective by Design* campaign in May 2006 and organized the first International Day Against DRM. The eighth event took place on May 3, 2013. A recent campaign is *Stop DRM in HTML5*, which is described as follows:

The World Wide Web Consortium (W3C) is considering a proposal to weave Digital Restrictions Management (DRM) into HTML5—in other words, into the very fabric of the Web. Millions of Internet users came together to defeat SOPA/PIPA, but now Big Media moguls are going through non-governmental channels to try to sneak digital restrictions into every interaction we have online. Giants like Netflix, Google, Microsoft, and the BBC are all rallying behind this disastrous proposal, which flies in the face of the W3C's mission to "lead the World Wide Web to its full potential."²⁷

Defective by Design is

... [A] participatory and grassroots campaign exposing DRM-encumbered devices and media for what they really are: Defective by Design. We are working together to eliminate DRM as a threat to innovation in media, the privacy of readers, and freedom for computer users. Our actions involve identifying and targeting defective products, pressuring media retailers and hardware manufacturers to stop supporting DRM, exposing the immense concentration of power over media created by DRM, and raising awareness of DRM to libraries, schools, and individuals around the world.²⁸

Industry supporters of DRM refer to it as "digital rights management" as if they are the ultimate authority to grant us our rights, as if they are the ones who should have complete and total control over how we use and interact with our media. What they are really doing is managing the restrictions they impose on our media and devices that we would normally have control over in the absence of DRM. We should own our media, not be at the mercy of media companies. For that reason, we refer to it as "Digital Restrictions Management."²⁹

Defective by Design argues that DRM is grossly intrusive and unfair:

Amazon's new movie download service is called Unbox and it outlines what DRM implies. The user agreement requires that you allow Unbox DRM software to monitor your hard drive and to report activity to Amazon. These reports would thus include a list of: all the software installed; all the music and video you have; all your computer's interaction with other devices. You will surrender your freedom to such an extent that you will only be able to regain control by removing the software. But if you do remove the software you will also remove all your movies along with it. You are restricted even geographically, and you lose your movies if you ever move out of the USA. You of course have to agree that they can change these terms at any time. Microsoft's newly upgraded Windows Media Player 11 (WMP11) user agreement has a similar set of terms.²⁹

Stop DRM Now! argues along the same lines as *Defective by Design*:

DRM exists for the exclusive benefit of content producers and providers. Companies such as Disney, Sony, and Lion's Gate argue that DRM is needed to prevent people from pirating

FUNDAMENTAL PROBLEMS 42 · 17

music, movies, and other works on P2P networks or by other means. What they won't tell you is that they are really trying to control who, what, when, where, and how you access your music and videos. For instance, when you purchase a song from Apple's iTunes Music Store, you are purchasing a song that can only be played by media applications that support QuickTime or an Apple iPod. Suppose next year you want to buy a new portable music player? Instead of having choice, you must now buy another iPod if you expect to play music you've already purchased.³⁰

42.6.2 Technical Countermeasures

42.6.2.1 Reverse Engineering. Reverse engineering allows a programmer to work backward from the finished program or product. Encryption keys can be extracted by reverse engineering playback software. Reverse engineering can circumvent most anti-piracy solutions. As a result, manufacturers of anti-piracy software and hardware are strongly opposed to permitting reverse engineering. The DMCA does allow reverse engineering, but the provisions of DMCA were not intended to enable the circumvention of technical protection measures (TPMs) in order to gain unauthorized access to or to make unauthorized copies of copyrighted works.

42.6.2.2 Published Attacks. The most notable attack on a software key took place when a licensee of CSS neglected to encrypt a decryption key. Obtaining a key by reverse engineering the XingDVD from Xing Technologies, the hackers were then able to guess many other keys. This left the hackers with a collection of decryption keys; even if the XingDVD key was removed, they could still copy DVDs by using the other keys. Using the results of this compromise, a group of people including the Norwegian teenager Jon Lech Johansen developed a program called *DeCSS* to decrypt CSS-encrypted DVDs and play them on Linux machines.

A variety of groups, including the DVD CCA, sued Johansen for publishing tools to subvert copyright protection. Johansen was acquitted twice in Norwegian courts. As of 2007, all remaining lawsuits against him were dropped,³¹ and a number of programs like DeCSS are freely available across the Internet. The next section discusses the widespread availability of DRM-cracking tools.

42.6.2.3 Tools for Cracking DRM. Software for breaking DRM controls abounds on the Web. Methods exist to defeat restrictions on:

- CD copying
- DVD copying
- Blu-ray copying
- iTunes copying
- PDF content extraction
- Software activation

We decline to include specific links to such tools in this *Handbook*.

42.7 FUNDAMENTAL PROBLEMS. A number of experts have pointed out that there are fundamental flaws in all the methods for preventing illegal copying of digital materials as described in this chapter. Bruce Schneier, a respected cryptographer, has repeatedly explained that all digital information must be converted to a cleartext

42 · 18 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

(unencrypted) form before it is displayed or otherwise used. Schneier calls this “the Achilles’ heel of all content protection schemes based on encryption.”³² Because the cleartext version has to reside somewhere in volatile or nonvolatile memory for at least some period of time to be usable, it is theoretically possible to obtain a copy of the cleartext version regardless of the complexity of the methods that originally concealed or otherwise limited access to the data. For example, if a DVD movie using complex regional encoding is to be seen on a monitor, at some point the hardware and software that decoded the DVD have to send that data stream to a monitor driver. The decoded data stream is vulnerable to interception. By modifying the low-level routines in the monitor driver, it is possible to divert a copy of the raw data stream to a storage device for unauthorized replay or reconstitution. Similarly, a system may be devised to prevent more than one copy of a document from being printed directly on a printer; however, unless the system prevents screen snapshots, a user can circumvent the restrictions by copying the screen as a bit image and storing that image for later, unauthorized use. Although hardware devices such as dedicated DVD or CD players may successfully interfere with piracy for some time, the problem is exacerbated under the current popular operating systems that have no security kernel and thus allow any processes to access any region of memory without regard to security levels.

42.8 SUMMARY. Piracy is a rapidly growing societal problem affecting a multitude of people and industries. Although producers may suffer the greatest financial losses, there is a substantial impact on consumers. Pirated copies are generally inferior in quality and are sometimes defective. If anything goes wrong, pirated copies, being illegal, are unsupported. Additionally, producers’ financial losses due to pirated copies may push the cost of legitimate copies up. Retailers and distributors also suffer due to the loss of sales to pirates. Illegal copies generally are sold more inexpensively than legitimate copies, so retailers and distributors cannot compete on price.

Creative talent, whether software developers, writers, musicians, artists, or performers, plus all the people who helped create the book, magazine, record, performance, painting, concert, or other media, are cheated out of their royalties by pirates. Frequently, because of the amount of time and effort needed to create the end product, the creators depend on the royalties for their livelihood. In addition, poor quality of stolen concepts can irreparably damage the reputation of the creative talent.

Publishers, record companies, art dealers, and other individuals and companies that invest artistic and technical skill along with money and effort to create an original work also lose revenues when that work is pirated. Because of the expenses already laid out to create the original product, companies frequently have to recoup their losses by raising prices for the consumer.

Due to the sophistication of systems and the increased use of the Internet, piracy has become more widespread and has an even greater financial impact worldwide. Many different types of antipiracy systems and techniques have been developed and implemented in an effort to cut down on the ever-increasing instances of piracy. One drawback to all of the systems is the lack of standards applied to software, audio, video, and other media. One of the most promising antipiracy systems is digital rights management. However, even DRM systems are not yet standardized, thus creating even more confusion regarding which is best and what to use.

The media industry is still working out which DRM solution it likes best, or whether it likes DRM at all. In 2007, Apple’s Steve Jobs announced that a large portion of EMI’s song catalog will be available for DRM-free download from the iTunes Music Store, at a price per song just \$0.30 higher than the standard \$0.99. Previously downloaded songs,

GLOSSARY 42 · 19

with DRM, will be upgradeable to a DRM-free version for the difference in the two prices.³³ This move echoes Jobs's recent comments encouraging the music industry to move away from DRM. Apple's actions notwithstanding, some consumers and privacy organizations continue fighting what they see as serious threats to individual privacy in nascent DRM efforts. The use of network proxies and hidden operating systems is adding to the difficulty of discovering inappropriate use of content, much less preventing or prosecuting it. As the balance between content producers and consumers works itself out, it seems likely that many more technologies will come and go, and DRM may be the harbinger of things to come or an unfortunate choice on the way to a bold, new business model for the media industry.

42.9 GLOSSARY. Discussions of piracy and privacy often are riddled with a confusing array of terms and acronyms. Some of the most commonly used terms, organizations, and acronyms are listed in this glossary.

AAC—Advanced Audio Coding. A standardized, lossy compression and encoding scheme for digital audio. Most commonly used format for compressing audio CDs for Apple's iPod and iTunes.

ACATS. Advisory Committee on Advanced Television Service.

Anti-Bootleg Statute (Section 2319A). A U.S. federal statute that criminalizes the unauthorized manufacture, distribution, or trafficking in sound recordings and music videos of “live” musical performances.

ATSC. Advanced Television Systems Committee.

ATV. Advanced Television.

Bootleg recordings. The unauthorized recording of a musical broadcast on radio or television, or at a live concert or performance. These recordings are also known as *underground recordings*.

BSA—Business Software Alliance. A consortium of major software developers, including IBM, Microsoft, Novell, Apple, Dell, and Sun Microsystems, that is attempting to stem lost revenues from pirated computer software. BSA educates computer users on software copyrights and fights software piracy. Individual members, such as Microsoft and Adobe, have their own anti-piracy programs in addition to belonging to the BSA.

CEA. Consumer Electronics Association.

CEMA. Consumer Electronics Manufacturers Association.

CSS—Content Scrambling System. A form of data encryption used to discourage reading media files directly from the disc, without a decryption key. Descrambling the video and audio requires a 5-byte, 40-bit key.

DeCSS—Descrambling Content Scrambling System. A utility developed by Norwegian programmers via reverse engineering and posted on the Web. This utility decrypts CSS and allows individuals to make illegal copies of DVD movies.

DFAST. Dynamic Feedback Arrangement Scrambling Technique.

DMAT. Digital Music Access Technology.

DMCA. The Digital Millennium Copyright Act signed into law October 28, 1998. Designed to implement World Intellectual Property Organization (WIPO) treaties

42 · 20 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

(signed in Geneva in December 1996); the DMCA strengthens the protection of copyrighted materials in digital formats.

DivX. A brand name of products created by DivX, Inc. (formerly DivXNetworks, Inc.), including the DivX Codec. Known for its ability to compress lengthy video segments into small sizes, it has been the center of controversy because of its use in the replication and distribution of copyrighted DVDs. Many newer DVD players are able to play DivX movies.

DRM—Digital Rights Management. Refers to any of several technologies used by publishers or copyright owners to control access to, and usage of, digital data or hardware and to restrictions associated with a specific instance of a digital work or device.

DVD CCA—DVD Copy Control Association. A not-for-profit corporation that owns and licenses CSS. DVD CCA has filed numerous lawsuits against companies and individuals that make pirated copies of films.

EFF—Electronic Frontier Foundation. A nonprofit organization working in the public interest to protect fundamental civil liberties, including privacy where computers and the Internet are concerned. The organization frequently disagrees with the steps that other organizations and corporations want to take to protect copyrighted materials.

EPIC—Electronic Privacy Information Center. A public interest research group established in 1994 to focus public attention on emerging civil liberties issues and to protect privacy, the First Amendment, and constitutional values.

FairPlay. A digital rights management (DRM) technology created by Apple Inc. (formerly known as Apple Computer), built in to the QuickTime multimedia technology, and used by the iPod, iTunes, and the iTunes Store. Every file bought from the iTunes Store with iTunes is encoded with FairPlay. It digitally encrypts AAC audio files and prevents users from playing these files on unauthorized computers.

FAST—Federation against Software Theft. A group headquartered in Great Britain represents software manufacturers and works with law enforcement agencies in finding and stopping software pirates in Europe.

FCC. Federal Communications Commission.

grpff. A 7-line, 526-character program of Perl code developed by two students at the Massachusetts Institute of Technology. More compact than DeCSS, the program descrambles DVDs but does not contain a decryption key. The code is readily available on the MIT campus via hats, T-shirts, business cards, and bumper stickers.

Hard disk loading. PCs with unlicensed software preinstalled. Use of a single copy of a software program but installed illegally on many machines. The original disks and the documents that should come with the PC are often missing or incomplete.

HDTV—High-definition television. Digital television transmissions are mandated in the U.S. to be the standard.

HRRC—Home Recording Rights Coalition. A coalition representing consumers, retailers, and manufacturers of audio and audiovisual recording products and media. The HRRC dedicates itself to keeping products free of government-imposed charges or restraints on the products' distribution or operation.

GLOSSARY 42 · 21

IFPI—International Federation of the Phonographic Industry. An organization that promotes the interests of the international recording industry worldwide. Its mission is to fight music piracy; promote fair market access and good copyright laws; help develop the legal conditions and the technologies for the recording industry to prosper in the digital era; and promote the value of music.

IIPA—International Intellectual Property Alliance. A private-sector coalition formed in 1984. The organization represents the U.S. copyright-based industries in efforts to improve international protection of copyrighted materials.

MP3. A technology for downloading music files using the MPEG format via the Internet.

MPAA—Motion Picture Association of America. Composed of member companies that produce and distribute legitimate films and videos. This organization serves as the official voice and advocate of the American motion picture industry. MPAA also assists law enforcement in raids and seizure of pirated videocassettes and DVDs.

MPEG—Moving Pictures Experts Group. A generic means of compactly representing digital video and audio signals for consumer distribution. MPEG video syntax provides an efficient way to represent image sequences in the form of more compact-coded data.

NMPA—National Music Publishers Association. A trade association that represents 700 U.S. businesses that own, protect, and administer copyrights in musical works.

NTSC. National Television System Committee.

PET—Privacy-Enhancing Technologies. The general term for a variety of new technologies and Internet protocols designed to enhance online privacy; includes anonymizing proxies, mixing networks, and onion routing.

Pirated recordings. Unauthorized duplicates of the sounds of one or more legitimate recordings.

RIAA—Recording Industry Association of America. The group has an antipiracy unit that handles initial examination of product on behalf of the recording industry. Pirates can be turned in by calling RIAA at 1-800-BAD-BEAT.

SAG—Screen Actors Guild. The union for screen actors. Members do not get residuals from pirated films, as they do from authorized copies.

SDMI—Secure Digital Music Initiative. A forum of 200 companies from the electronics, music, telecommunications, and information technology industries and the RIAA. The group was active from 1998 to 2001.

SIIA—Software & Information Industry Association. A trade organization of the software and information content industries representing over 800 high-tech companies that develop and market software and electronic content. The organization provides policies and procedures for dealing with software and Internet use within businesses. SIIA also provides guidelines for telling if software is pirated or counterfeited. The SIIA Anti-Piracy Hotline is 1-800-388-7478.

SPA—Software Publishers Association. This association, a division of SIIA, assists in enforcement in dealing with software piracy and also provides education about software piracy. SPAudit Software is one of the first software audit and inventory

42 · 22 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

tools made available for use by companies (in the 1980s). Improved versions of the software are now available.

Trademark Counterfeiting—Title 18 U.S.C. Section 2320. A federal statute that deals with sound recordings that contain the counterfeit trademark of the legitimate manufacturer or artists.

Trafficking in Counterfeit Labels—Title 18 U.S.C., Section 2318. A federal statute that covers counterfeit labels printed for use on a sound recording.

U.S. Copyright Law (Title 17 U.S.C.). A federal law that protects copyright owners from the unauthorized reproduction or distribution of their work.

WGA—Writers Guild of America. The union for writers of television, video, and film scripts.

WMA—Windows Media Audio. A proprietary compressed audio file format developed by Microsoft Corporation.

42.10 FURTHER READING

Business Software Alliance—USA Homepage. *Anti-Piracy Information*, 2007, www.bsa.org

Cohen, Julie E. “DRM and Privacy.” *Berkeley Technology Law Journal*. 2003. www.law.berkeley.edu/institutes/bclt/drm/papers/cohen-drmandprivacy-btlj2003.html

Copyright Act of 1976 (Public Law 94-553); Title 17 U.S.C. Sections 101–120, www4.law.cornell.edu/uscode/html/uscode17/usc_sec_17_00000101-000-.html

Electronic Freedom Foundation home page: www.eff.org

Electronic Privacy Information Center (EPIC) home page: www.epic.org

Gross, T. “The Music Industry, Adapting to a Digital Future: Terry Gross interviews Eliot Van Buskirk.” *Fresh Air* podcast, March 13, 2008, www.npr.org/templates/story/story.php?storyId=88145070

Harte, L. *Introduction to Digital Rights Management (DRM); Identifying, Tracking, Authorizing and Restricting Access to Digital Media*. Fuquay Varina, NC: Althos, 2006.

May, C. *Digital Rights Management: The Problem of Expanding Ownership Rights*. Oxford, UK: Chandos Publishing, 2006.

Recording Industry Association of America: www.riaa.com

Schneier.com: www.schneier.com/index.html

Zeng, W., H. Yu, and C-Y. Lin, eds. *Multimedia Security Technologies for Digital Rights Management*. New York: Academic Press, 2006.

42.11 NOTES

1. *Shadow Market: 2011 BSA Global Software Piracy Study*, 9th ed., May 2012, http://globalstudy.bsa.org/2011/downloads/study_pdf/2011_BSA_Piracy_Study-Standard.pdf
2. www.siia.net
3. “SIIA Shuts Down Notorious Software Pirate, Gaining Large Financial Sum and Cooperation with Further Investigations,” April 23, 2012, http://siia.net/index.php?option=com_content&view=article&id=1057:siia-shuts-down-notorious-software-pirate-gaining-large-financial-sum-and-cooperation-with-further-investigations&catid=62:press-room-

NOTES 42 · 23

4. Microsoft, "Microsoft Incorporates New Anti-Piracy Technologies in Windows 2000, Office 2000," www.microsoft.com/Presspass/press/2000/feb00/apfeaturespr.mspx
5. Microsoft, "Description of the Behavior of Reduced Functionality Mode in Windows Vista," <http://support.microsoft.com/kb/925582>
6. Microsoft, "Microsoft Product Activation," www.microsoft.com/licensing/resources/vol/default.mspx
7. RIAA, "Anti-Piracy," www.riaa.com/issues/piracy/default.asp (URL inactive).
8. GrayZone Digest, "Worldwide Update," October 1997, www.grayzone.com/1097.htm
9. P. Mercer, "Copycat CDs in an Instant," BBC News, April 16, 2002, <http://news.bbc.co.uk/2/hi/entertainment/1930923.stm>
10. Recording Industry Association of America, "2005 Commercial Piracy Report," <http://76.74.24.142/6BE200AF-5DDA-1C2B-D8BA-4174680FCE66.pdf>
11. E. VanBuskirk, "Fans Pay Whatever They Want for Radiohead's Upcoming Album," Underwire Blog from Wired.com, October 1, 2007, <http://blog.wired.com/underwire/2007/10/fans-to-determi.html>
12. J. Valentini, Napster Statement, February 12, 2001, via Music Industry News Network (mi2n.com), www.mi2n.com/press.php3?press_nb=18419
13. S. Butler, "RIAA Sends Another Wave of Settlement Letters," Billboard.biz, September 20, 2007, www.billboard.biz/bbbiz/content_display/industry/e3i39f76c017d89e0747eaafdf53d458f14b (URL inactive).
14. "University of Nebraska Will Bill RIAA \$11 for Each Threatening Letter Received," Consumerist.com, March 26, 2007.
15. S. Musil, "Scour to End File-Swapping Service," CNET News.com, November 14, 2000, <http://cnet.news.com/2100-1023-248631.html>
16. J. Borland, "Broadcasters Win Battle Against iCraveTV.com," CNET News.com, January 28, 2000, <http://cnet.news.com/2100-1033-236255.html>
17. M. Graziano and Lee Rainie, "The Music Downloading Deluge: 37 Million American adults and youths have retrieved music files on the Internet," April 24, 2001, <http://www.pewinternet.org/Reports/2001/The-Music-Downloading-Deluge.aspx>
18. J. Silverstein, "iTunes: 1 Billion Served," ABC News Website, February 23, 2006, <http://abcnews.go.com/Technology/story?id=1653881>
19. C. Cohn, "Security Researchers Drop Scientific Censorship Case," www.eff.org/IP/DMCA/Felten_v_RIAA/20020206_eff_felten_pr.html (URL inactive).
20. www.sdmi.org/whats_new.htm
21. www.contentguard.com
22. D. Chaum, "Untraceable Electronic Mail, Return Addresses, and Digital Pseudonyms," *Communications of the ACM* 24, No. 2 (February 1981); available online at <http://freehaven.net/anonbib/cache/chaum-mix.pdf>
23. Tor: Anonymity Online, <http://tor.eff.org>
24. B. Schneier, "Anonym.OS," www.schneier.com/blog/archives/2006/01/anonymos.html
25. <https://www.eff.org/about>
26. <https://www.eff.org/issues/drm>
27. www.fsf.org/campaigns

42 · 24 PROTECTING DIGITAL RIGHTS: TECHNICAL APPROACHES

28. www.defectivebydesign.org
29. www.defectivebydesign.org/what_is_drm_digital_restrictions_management
30. <http://stopdrmnow.org>
31. C. Cohn, "DVD Descrambling Code Not a Trade Secret," http://www.eff.org/IP/Video/DVDCCA_case/20040122_eff_pr.php (URL inactive).
32. B. Schneier, "The Futility of Digital Copy Prevention," *Crypto-Gram Newsletter*, May 15, 2001, www.schneier.com/crypto-gram-0105.html
33. American Public Radio, *Marketplace*, April 2, 2007, <http://marketplace.publicradio.org/shows/2007/04/02/PM200704024.html> (URL inactive).