

CLOUD COMPUTING

THEORY AND PRACTICE

DAN C. MARINESCU

MK
Morgan Kaufmann

THIRD EDITION

Cloud Computing

Theory and Practice

Cloud Computing

Theory and Practice

Third Edition

Dan C. Marinescu

Department of Computer Science
University of Central Florida
Orlando, FL, United States



MORGAN KAUFMANN PUBLISHERS

ELSEVIER

AN IMPRINT OF ELSEVIER

Morgan Kaufmann is an imprint of Elsevier
50 Hampshire Street, 5th Floor, Cambridge, MA 02139, United States

Copyright © 2023 Elsevier Inc. All rights reserved.

MATLAB® is a trademark of The MathWorks, Inc. and is used with permission.

The MathWorks does not warrant the accuracy of the text or exercises in this book.

This book's use or discussion of MATLAB® software or related products does not constitute endorsement or sponsorship by The MathWorks of a particular pedagogical approach or particular use of the MATLAB® software.

No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or any information storage and retrieval system, without permission in writing from the publisher. Details on how to seek permission, further information about the Publisher's permissions policies and our arrangements with organizations such as the Copyright Clearance Center and the Copyright Licensing Agency, can be found at our website: www.elsevier.com/permissions.

This book and the individual contributions contained in it are protected under copyright by the Publisher (other than as may be noted herein).

Notices

Knowledge and best practice in this field are constantly changing. As new research and experience broaden our understanding, changes in research methods, professional practices, or medical treatment may become necessary.

Practitioners and researchers must always rely on their own experience and knowledge in evaluating and using any information, methods, compounds, or experiments described herein. In using such information or methods they should be mindful of their own safety and the safety of others, including parties for whom they have a professional responsibility.

To the fullest extent of the law, neither the Publisher nor the authors, contributors, or editors, assume any liability for any injury and/or damage to persons or property as a matter of products liability, negligence or otherwise, or from any use or operation of any methods, products, instructions, or ideas contained in the material herein.

Library of Congress Cataloging-in-Publication Data

A catalog record for this book is available from the Library of Congress

British Library Cataloguing-in-Publication Data

A catalogue record for this book is available from the British Library

ISBN: 978-0-323-85277-7

For information on all Morgan Kaufmann publications
visit our website at <https://www.elsevier.com/books-and-journals>

Publisher: Katey Birtcher

Acquisitions Editor: Steve Merken

Editorial Project Manager: Naomi Robertson

Production Project Manager: Nadhiya Sekar

Designer: Matthew Limbert

Typeset by VTeX

Printed in United States

Last digit is the print number: 9 8 7 6 5 4 3 2 1



Working together
to grow libraries in
developing countries

www.elsevier.com • www.bookaid.org

To Vera Rae and Luke Bell

Contents

Preface to third edition	xiii
Abbreviations	xv
CHAPTER 1 Introduction	1
1.1 Cloud computing, an old idea whose time has come	3
1.2 Energy use and ecological impact of cloud computing	7
1.3 Ethical issues in cloud computing	8
1.4 Factors affecting cloud service availability	9
1.5 Network-centric computing and network-centric content	10
CHAPTER 2 The cloud ecosystem	13
2.1 Cloud computing delivery models and services	14
2.2 Amazon Web Services	18
2.3 Google Clouds	25
2.4 Microsoft Windows Azure and online services	28
2.5 IBM clouds	30
2.6 Cloud storage diversity and vendor lock-in	30
2.7 Cloud interoperability	31
2.8 Service-level Agreements and Compliance-level Agreements	33
2.9 Responsibility sharing between user and service provider	34
2.10 User challenges and experience	34
2.11 Software licensing	36
2.12 Challenges faced by cloud computing	37
2.13 Cloud computing as a disruptive technology	39
2.14 Exercises and problems	40
CHAPTER 3 Parallel processing and distributed computing	41
3.1 Computer architecture concepts	42
3.2 Grand architectural complications	48
3.3 ARM architecture	55
3.4 SIMD architectures	58
3.5 Graphics processing units	60
3.6 Tensor processing units	62
3.7 Systems on a chip	65
3.8 Data, thread-level, and task-level parallelism	66
3.9 Speedup, Amdahl's law, and scaled speedup	68
3.10 Multicore processor speedup	70
3.11 From supercomputers to distributed systems	72
3.12 Modularity. Soft modularity versus enforced modularity	74
3.13 Layering and hierarchy	79
3.14 Peer-to-peer systems	82

3.15	Large-scale systems	84
3.16	Composability bounds and scalability (R)	86
3.17	Distributed computing fallacies and the CAP theorem	88
3.18	Blockchain technology and applications	89
3.19	History notes and further readings	91
3.20	Exercises and problems	94
CHAPTER 4	Cloud hardware and software	95
4.1	Cloud infrastructure challenges	96
4.2	Cloud hardware; warehouse-scale computer (WSC)	98
4.3	WSC performance	100
4.4	Hypervisors	104
4.5	Execution of coarse-grained data-parallel applications	104
4.6	Fine-grained cluster resource sharing in Mesos	106
4.7	Cluster management with Borg	107
4.8	Evolution of a cluster management system	110
4.9	Shared state cluster management	111
4.10	QoS-aware cluster management	113
4.11	Resource isolation	116
4.12	In-memory cluster computing for Big Data	120
4.13	Containers; Docker containers	128
4.14	Kubernetes	130
4.15	Further readings	131
4.16	Exercises and problems	132
CHAPTER 5	Cloud resource virtualization	135
5.1	Resource virtualization	136
5.2	Performance and security isolation in computer clouds	137
5.3	Virtual machines	138
5.4	Full virtualization and paravirtualization	140
5.5	Hardware support for virtualization	141
5.6	QEMU	144
5.7	Kernel-based Virtual Machine	145
5.8	Xen—a hypervisor based on paravirtualization	148
5.9	Optimization of network virtualization in Xen 2.0	154
5.10	Nested virtualization	156
5.11	A trusted kernel-based virtual machine for ARMv8	159
5.12	Paravirtualization of Itanium architecture	161
5.13	A performance comparison of virtual machines	164
5.14	Open-source software platforms for private clouds	167
5.15	The darker side of virtualization	169
5.16	Virtualization software	170
5.17	History notes and further readings	171
5.18	Exercises and problems	172

CHAPTER 6	Cloud access and cloud interconnection networks	175
6.1	Packet-switched networks and the Internet	176
6.2	Internet evolution	181
6.3	TCP congestion control	183
6.4	Content-centric networks; named data networks (R)	185
6.5	Software-defined networks; SD-WAN	187
6.6	Interconnection networks for computer clouds	189
6.7	Multistage interconnection networks	193
6.8	InfiniBand and Myrinet	194
6.9	Storage area networks and the Fibre Channel	197
6.10	Scalable data center communication architectures	200
6.11	Network resource management algorithms (R)	204
6.12	Content delivery networks	207
6.13	Vehicular ad hoc networks	211
6.14	Further readings	212
6.15	Exercises and problems	212
CHAPTER 7	Cloud data storage	215
7.1	Dynamic random access memories and hard disk drives	216
7.2	Solid-state disks	217
7.3	Storage models, file systems, and databases	220
7.4	Distributed file systems; the precursors	223
7.5	General parallel file system	228
7.6	Google file system	231
7.7	Locks; Chubby—a locking service	233
7.8	RDBMS—cloud mismatch	238
7.9	NoSQL databases	239
7.10	Data storage for online transaction processing systems	241
7.11	BigTable	243
7.12	Megastore	245
7.13	Storage reliability at scale	246
7.14	Disk locality versus data locality in computer clouds	250
7.15	Database provenance	252
7.16	History notes and further readings	254
7.17	Exercises and problems	255
CHAPTER 8	Cloud security	257
8.1	Security—the top concern for cloud users	258
8.2	Cloud security risks	259
8.3	Security as a service (SecaaS)	264
8.4	Privacy and privacy impact assessment	264
8.5	Trust	267
8.6	Cloud data encryption	268
8.7	Security of database services	270
8.8	Operating system security	272

8.9	Virtual machine security	273
8.10	Security of virtualization	275
8.11	Security risks posed by shared images	278
8.12	Security risks posed by a management OS	281
8.13	Xoar—breaking the monolithic design of the TCB	283
8.14	Mobile devices and cloud security	286
8.15	Mitigating cloud vulnerabilities in the age of ransomware	287
8.16	AWS security	289
8.17	Further readings	290
8.18	Exercises and problems	291
CHAPTER 9	Cloud resource management and scheduling	293
9.1	Policies and mechanisms for resource management	294
9.2	Scheduling algorithms for computer clouds	296
9.3	Delay scheduling (R)	298
9.4	Data-aware scheduling (R)	303
9.5	Apache capacity scheduler	306
9.6	Start-time fair queuing (R)	307
9.7	Borrowed virtual time (R)	311
9.8	Cloud scheduling subject to deadlines (R)	315
9.9	MapReduce application scheduling subject to deadlines (R)	320
9.10	Resource bundling; combinatorial auctions for cloud resources	322
9.11	Cloud resource utilization and energy efficiency	325
9.12	Resource management and dynamic application scaling	328
9.13	Control theory and optimal resource management (R)	329
9.14	Stability of two-level resource allocation strategy (R)	333
9.15	Feedback control based on dynamic thresholds (R)	334
9.16	Coordination of autonomic performance managers (R)	336
9.17	A utility model for cloud-based web services (R)	338
9.18	Cloud self-organization	342
9.19	Cloud interoperability	344
9.20	Further readings	346
9.21	Exercises and problems	346
CHAPTER 10	Concurrency and cloud computing	349
10.1	Enduring challenges	350
10.2	Communication and concurrency	353
10.3	Computational models; communicating sequential processes	358
10.4	The bulk synchronous parallel model	360
10.5	A model for multicore computing	361
10.6	Modeling concurrency with Petri nets	363
10.7	Process state; global state of a process or thread group	369
10.8	Communication protocols and process coordination	374
10.9	Communication, logical clocks, and message delivery rules	376
10.10	Runs and cuts; causal history	381

10.11	Threads and activity coordination	385
10.12	Critical sections, locks, deadlocks, and atomic actions	392
10.13	Consensus protocols	397
10.14	Load balancing	399
10.15	Multithreading in Java; FlumeJava; Apache Crunch	405
10.16	History notes and further readings	407
10.17	Exercises and problems	408
CHAPTER 11	Cloud applications	411
11.1	Cloud application development and architectural styles	412
11.2	Coordination of multiple activities	415
11.3	Workflow patterns	419
11.4	Coordination based on a state machine model—zookeeper	422
11.5	MapReduce programming model	425
11.6	Case study: the GrepTheWeb application	428
11.7	Hadoop, Yarn, and Tez	431
11.8	SQL on Hadoop: Pig, Hive, and Impala	435
11.9	Current cloud applications and new applications opportunities	440
11.10	Clouds for science and engineering	442
11.11	Cloud computing and biology research	446
11.12	Social computing, digital content, and cloud computing	448
11.13	Software fault isolation	450
11.14	Further readings	451
11.15	Exercises and problems	452
CHAPTER 12	Big Data, data streaming, and the mobile cloud	453
12.1	Big Data	454
12.2	Data warehouses and Google databases for Big Data	456
12.3	Dynamic data-driven applications	463
12.4	Data streaming	466
12.5	A dataflow model for data streaming	470
12.6	Joining multiple data streams	473
12.7	Mobile computing and applications	475
12.8	Energy efficiency of mobile computing	478
12.9	Alternative mobile cloud computing models	479
12.10	System availability at scale (R)	482
12.11	Scale and latency (R)	484
12.12	Edge computing and Markov decision processes (R)	488
12.13	Bootstrapping techniques for data analytics (R)	492
12.14	Approximate query processing (R)	495
12.15	Further readings	498
12.16	Exercises and problems	499
CHAPTER 13	Emerging clouds	501
13.1	A short-term forecast	502
13.2	Machine learning on clouds	503

13.3	Quantum computing on clouds	508
13.4	Vehicular clouds	518
13.5	Final thoughts	527
APPENDIX A	Cloud projects	529
A.1	Cloud simulation of a distributed trust algorithm	529
A.2	A trust management service	534
A.3	Simulation of traffic management in a smart city	540
A.4	A cloud service for adaptive data streaming	545
A.5	Optimal FPGA synthesis	550
A.6	Tensor network contraction on AWS	552
A.7	A simulation study of machine-learning scalability	559
A.8	Cloud-based task alert application	561
A.9	Cloud-based health-monitoring application	565
APPENDIX B	Cloud application development	571
B.1	AWS EC2 instances	572
B.2	Connecting clients to cloud instances through firewalls	575
B.3	Security rules for application- and transport-layer protocols in EC2	577
B.4	How to launch an EC2 Linux instance and connect to it	581
B.5	How to use S3 in Java	582
B.6	How to manage AWS SQS services in C#	585
B.7	How to install SNS on Ubuntu 10.04	586
B.8	How to create an EC2 placement group and use MPI	588
B.9	StarCluster—a cluster computing toolkit for EC2	590
B.10	An alternative setting of an MPI virtual cluster	590
B.11	How to install hadoop on eclipse on a windows system	592
B.12	Exercises and problems	595
Literature	597	
Glossary	621	
Index	635	

Preface to third edition

This text targets a broad audience, including the large number of individuals without a formal training in computer science, yet addicted to cloud computing, who want to optimize their applications and wish to better understand the inner workings of intricate systems such as computer clouds. The text could also be useful to software engineers tempted to exploit concurrency and take advantage of the computing power of multicore processors, graphics processing units, and tensor processing units. Graduate students who conduct research on different aspects of organization and management of computer clouds and complex systems could benefit from research sections marked (R) in several chapters.

Cloud computing has had a transformative impact not only on how we compute today but also on the algorithms we use and on computer systems' hardware and software. The text tracks the evolution of ideas in the field striving to maintain a delicate balance between what a gardener would call perennials, i.e., theoretical concepts mostly invariant in time and critical for understanding and exploiting concurrency, and annuals, i.e., practical recipes with a limited lifespan, helpful for using computer clouds now.

Any text dedicated to cloud computing must be periodically updated to keep up with the breathtaking evolution of cloud computing engines, the cloud software stack, cloud interconnection networks, and applications. The third edition of this book includes a reorganization of chapters, major revisions of most chapters, cloud computing applications to machine learning, Big Data, and other important areas, and a new chapter on emerging clouds.

Chapters 1 and 2 give an informal introduction to computer clouds and to the cloud ecosystem. Cloud computing is intimately tied to parallel and distributed processing, and Chapter 3 introduces important theoretical and practical concepts related to parallel and distributed computing. Chapters 4 and 5 discuss the cloud infrastructure and virtualization critical for facilitating access to cloud resources. Chapter 6 is dedicated to communication and cloud access and presents the cloud networking infrastructure, interconnection networks, and a scalable data-center communication architectures. Chapter 7 covers storage models, storage reliability at scale, and database services. Chapter 8 is dedicated to cloud security, and Chapter 9 covers resource management and scheduling.

The next three chapters are focused on the other important component of the cloud ecosystem, cloud applications. Chapter 10 covers subjects critical for understanding and using concurrency, including computational models, atomic actions, load balancing, and consensus protocols. Chapter 11 dissects the components of the cloud software stack and the frameworks supporting processing of large data sets. Chapter 12 analyzes challenges posed by Big Data, data streaming, mobile applications, and edge computing. Chapter 13 covers emerging clouds including cloud services supporting machine learning, quantum computing on clouds and vehicular clouds; the vehicular clouds section is written by one of the pioneers of this field, Prof. Stefan Olariu from the Old Dominion University.

Appendix A presents several cloud projects in large-scale simulations and cloud services, including applications when multiple design alternatives are evaluated concurrently, as well as Big Data applications in computational sciences and machine learning applications. Appendix B covers cloud application development. Some 540 references are cited in the text. A glossary and a list of abbreviations are also provided. The author is grateful for the advice and support of Stephen Merken and the Morgan Kaufmann team.

Teaching ancillaries for this book, including solutions manual, PowerPoint lecture slides, and image bank, are available online to qualified instructors. Visit <https://inspectioncopy.elsevier.com/book/details/9780323852777> for more information and to register for access.

Abbreviations

ABI	Application Binary Interface
ACID	Atomicity, Consistency, Isolation, Durability
ACL	Access Control List
AMI	Amazon Machine Image
API	Application Program Interface
ARM	Advanced RISC Machine
ASIC	Application Specific Integrated Circuit
AVX	Advanced Vector Extensions
AWS	Amazon Web Services
AWSLA	Amazon Web Services Licensing Agreement
BASE	Basically Available, Soft state, Eventually consistent
BCE	Basic Core Equivalent
BIOS	Basic Input Output System
BPD	Bootstrap Performance Diagnostic
BSP	Bulk Synchronous Parallel
CCN	Content Centric Network
CDN	Content Delivery Network
CISC	Complex Instruction Set Computer
COMA	Cache Only Memory Access
CORBA	Common Object Request Broker Architecture
CPI	Cycles per Instruction
CPU	Central Processing Unit
CRM	Custom Relation Management
CSP	Cloud Service Provider
CUDA	Compute Unified Device Architecture
DAT	Dynamic Address Translation
DBaaS	Database as a Service
DDoS	Distributed Denial of Service
DHCP	Dynamic Host Configuration Protocol
DMA	Direct Memory Access
DRAM	Dynamic Random Access Memory
DVFS	Dynamic Voltage and Frequency Scaling
EBS	Elastic Block Store
ECS	EC2 Container Service
EC2	Elastic Cloud Computing
EMR	Elastic Map Reduce
EPIC	Explicitly Parallel Instruction Computing
FC	Fiber Channel
FCFS	First Come First Serve
FIFO	First In First Out
FPGA	Field Programmable Gate Array
FSM	Finite State Machine
GFS	Google File System
GiB	Gibi Byte
GIMP	GNU Manipulation Program
GPFS	General Parallel File System
GPU	Graphics Processing Unit

HDD	Hard Disk Drive
HDFS	Hadoop File System
HTML	Hypertext Markup Language
HTTP	Hypertext Transfer Protocol
IaaS	Infrastructure as a Service
IDE	Integrated Drive Electronics
IoT	Internet of Things
IPC	Instructions per Clock Cycle
IRTF	Internet Research Task Force
ISA	Instruction Set Architecture
JDBC	Java Database Connectivity
JMS	Java Message Service
JSON	Javascript Object Notation
KVM	Kernel-based Virtual Machine
MAC	Medium Access Control
MFLOPS	Million Floating Point Operations per Second
MIMD	Multiple Instruction Multiple Data
MIPS	Million Instructions per Second
MLMM	Multi-level Memory Manager
MMX	Multi Media Extension
MPI	Message Passing Interface
MSCR	Map Shuffle Combine Reduce
NAT	Network Address Translation
NDN	Named Data Networks
NFS	Network File System
NTFS	New Technology File System
NUMA	Non-Uniform Memory Access
NV-RAM	Non-Volatile Random Access Memory
OCCI	Open Cloud Computing Interface
OGF	Open Grid Forum
OLTP	On Line Transaction Processing
OLAP	On Line Analytical Processing
PaaS	Platform as a Service
PHP	recursive acronym for PHP: Hypertext Preprocessor
PN	Petri Net
QOS	Quality of Service
RAID	Redundant Array of Independent Disks
RAM	Random Access Memory
RAR	Read After Read
RAW	Read After Write
RDD	Resilient Distributed Dataset
RDBMS	Relational Database Management System
RDS	Relational Database Service
REST	Representational State Transfer
RFC	Remote Frame Buffer
RISC	Reduced Instruction Set Computer
RMI	Remote Method Invocation
RPC	Remote Procedure Call
RTT	Round Trip Time
SaaS	Software as a Service
SAN	Storage Area Network
SDK	Software Development Kit
SDN	Software Defined Network

SHA	Secure Hash Algorithm
SIMD	Single Instruction Multiple Data
SLA	Service Level Agreement
SNS	Simple Notification Service
SOAP	Simple Object Access Protocol
SPMD	Single Program Multiple Data
SQL	Structured Query Language
SQS	Simple Queue Service
SSD	Solid State Disk
SSE	Streaming SIMD Extensions
S3	Simple Storage System
SWF	Simple Workflow Service
TCP	Transport Control Protocol
TLB	Translation Lookaside Buffer
UDP	User Datagram Protocol
UFS	Unix File System
UMA	Uniform Memory Access
vCPU	Virtual CPU
VLIW	Very Long Instruction Word
VM	Virtual machine
VMCS	Virtual Machine Control Structure
VMM	Virtual Machine Monitor
VMM	Virtual Memory Manager
VNC	Virtual Network Computing
VPC	Virtual Private Cloud
VPN	Virtual Private Network
WAN	Wide Area Network
WAR	Write After Read
WAW	Write After Write
WSC	Warehouse Scale Computer
WWN	World Wide Name
XML	Extensible Markup Language
YARN	Yet Another Resource Negotiator

Introduction

1

Computer clouds are utilities providing computing services. In *utility computing*, the hardware and the software resources are concentrated in large data centers, and users of computing services pay as they consume computing, storage, and communication resources. While utility computing often requires a cloud-like infrastructure, the focus of cloud computing is on the business model for providing computing services.

NIST, the US National Institute of Standards and Technology, defines cloud computing as “a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources that can be rapidly provisioned and released with minimal management effort or service provider interaction.” Cloud computing is characterized by five attributes: *on-demand self-service, broad network access, resource pooling, rapid elasticity, and measured service*.

More than half a century ago, at the centennial anniversary of MIT, John McCarthy, the 1971 Turing Award recipient for his work in artificial intelligence, prophetically stated: “... If computers of the type I have advocated become the computers of the future, then computing may someday be organized as a public utility, just as the telephone system is a public utility... The computer utility could become the basis of a new and important industry.” McCarthy’s prediction is now a technological and social reality.

The cloud computing movement is motivated by the idea that data processing and storage can be done more efficiently on large farms of computing and storage systems accessible via the Internet. Computer clouds support a paradigm shift from local to network-centric computing and network-centric content where distant data centers provide the computing and storage resources. In this new paradigm, users relinquish control of their data and code to Cloud Service Providers (CSPs).

The cloud computing age started in 2006 when Elastic Cloud Computing (EC2) and Simple Storage Service (S3) were offered by Amazon Web Services (AWS). Six years later, in 2012, EC2 was used by businesses in 200 countries. S3 has surpassed two trillion objects, and routinely runs more than 1.1 million peak requests per second. The range of services offered by CSPs and the number of cloud users have increased dramatically during the last few years. Cloud computing offers scalable and elastic computing and storage services. The resources used for these services can be metered, and the users can be charged only for the resources they have used. Cloud computing is a business reality, as a large number of organizations have adopted this paradigm.

Internet users have discovered the appeal of cloud computing either directly or indirectly, through a variety of services, without knowing the role the clouds play in their lives. The number of cloud users will continually increase in the years to come as the vast computational resources provided by the cloud infrastructure and the exabytes of data stored in the clouds are streamed, downloaded, accessed and used for deep learning, designing and engineering complex systems, scientific discovery, education, business, analytics, art, and virtually all other aspects of human endeavor.

Data analytics, data mining, computational financing, scientific and engineering applications, gaming, and social networking, as well as other computational and data-intensive activities benefit from cloud computing. Content previously confined to personal devices, such as workstations, laptops, tablets, and smartphones, no longer need to be stored locally. Data stored on computer clouds can be shared among all these devices, and it is accessible whenever a device is connected to the Internet. Clouds continually evolve in predictable, as well as unpredictable ways; for example, edge computing proposes to minimize network traffic and response time by preprocessing data locally.

Almost half a century after the dawn of the computing era, an eternity in the age of the silicon, disruptive multicore technology and the enormous computing power of Graphics Processing Units (GPUs) and Tensor Processing Units (TPUs) challenge computer science community to develop new algorithms, programming environments, and tools and challenge application developers to exploit concurrency. There is no point now to wait for faster clock rates; we better design algorithms and applications able to use modern processors and co-processors. The new challenge is harnessing the power of millions of multicore processors and systems on a chip and allowing them to work in concert efficiently.

Cloud computing has shown that there are applications that can effortlessly exploit concurrency and, in the process, generate huge revenues. A new era in computing has begun, a time when Big Data hides nuggets of useful information and requires massive amounts of computing resources. In this era, “coarse” is good and “fine” is not good, at least as far as the granularity of parallelism is concerned.

The scale of computer clouds amplifies unanticipated benefits, as well as the nightmares of system designers. Even a slight improvement of server performance and/or of the algorithms for resource management could lead to huge cost savings and rave reviews. When engineering large-scale systems, an important lesson is to prepare for the unexpected because low probability events occur and can cause major disruptions. The failure of one of the millions of hardware and software components can be amplified, propagate throughout the entire system, and have catastrophic consequences.

Cloud computing is a *disruptive computing paradigm* and requires major changes in many areas of computer science and computer engineering, including data storage, computer architecture, networking, resource management, scheduling, and last but not least computer security. Computer clouds operate in an environment characterized by the *variability of everything* and by *conflicting requirements*. Such disruptive qualities of computer clouds ultimately demand new thinking in system design. This book covers challenges posed by the scale of the cloud infrastructure and the large population of cloud users with diverse applications and requirements.

Cloud computing is cost-effective because of *resource multiplexing*. Application data is stored closer to the site where it is used in a manner that is device and location-independent; potentially, this data-storage strategy increases reliability, as well as security. Organizations using computer clouds are relieved of supporting large IT teams, acquiring and maintaining costly hardware and software, and paying large electricity bills. CSPs can operate more efficiently due to economies of scale.

Cloud computing represents a dramatic shift in the design of systems capable of providing vast amounts of computing cycles and storage space. *Computer clouds use off-the-shelf, low-cost components*. During the previous four decades, powerful, one-of-a-kind systems were built at a high cost, with the most advanced components available at the time. In early 1990s, Gordon Bell argued that one-of-a-kind systems are not only expensive to build but also that the cost of rewriting applications for them is prohibitive. He anticipated that sooner or later massively parallel computing will evolve into *computing for the masses* [55].

Since there are virtually no bounds on the composition of digital systems controlled by software, we are tempted to build increasingly more complex systems, including systems of systems [340]. The behavior and the properties of such systems are not always well understood thus, we should not be surprised that large-scale systems will occasionally fail and computing clouds will occasionally exhibit an unexpected behavior.

Cloud computing reinforces the idea that *computing and communication are deeply intertwined*. Advances in one field are also critical for the other. Cloud computing would not have emerged as an alternative to traditional computing models before the Internet was able to support high-bandwidth, reliable, low-cost communication. High-performance switches are critical elements of supercomputers and clouds infrastructure. Internet routers use powerful processors and Artificial Intelligence (AI) to enhance security.

The architecture, the coordination mechanisms, the design methodology, and the analysis techniques for large-scale complex systems such as clouds will evolve in response to changes in technology, the environment, and the demands of user community. Some of these changes will reflect changes in communication and in the Internet itself in terms of speed, reliability, security, capacity to accommodate a larger addressing space by migration to IPv6, and so on. The complexity of the cloud computing infrastructure is unquestionable and raises questions such as: How can we manage such systems? Do we have to consider radically new ideas, such as self-management and self-repair for future clouds consisting of millions of servers? Should we migrate from a strictly deterministic view of such complex systems to a non-deterministic one? Answers to these questions provide a rich set of research topics for the computer science and engineering community.

The cloud movement is not without skeptics and critics. The critics argue that cloud computing is just a marketing ploy, that users may become dependent on proprietary systems, and that the failure of a large system such as the cloud could have significant consequences for a very large group of users who depend on the cloud for their computing and storage needs. Security and privacy are major concerns for cloud computing users.

1.1 Cloud computing, an old idea whose time has come

It is hard to pinpoint a single technological or architectural development that triggered the movement towards computer clouds. This movement is the result of a cumulative effect of developments in micro-processors, storage, and networking technologies coupled with architectural advancements in all these areas and with advances in software systems, tools, programming languages, and algorithms supporting distributed and parallel computing.

Through the years, we have witnessed the breathtaking evolution of solid-state technologies that led to the development of multicore processors. The proximity of multiple cores on the same silicon die allows cache-coherency circuitry to operate at a much higher clock rate than would be possible if signals were to travel off-chip. Systems on a chip (SoCs), like Apple's M1, with general-purpose cores, GPU cores, and TPU cores, deliver impressive computing power with lower power consumption.

Storage technology has also evolved dramatically. Solid-state disks enable systems to manage very high transaction volumes and larger numbers of concurrent users while the price of memory has dropped significantly. Optical storage technologies and flash memories are widely used nowadays.

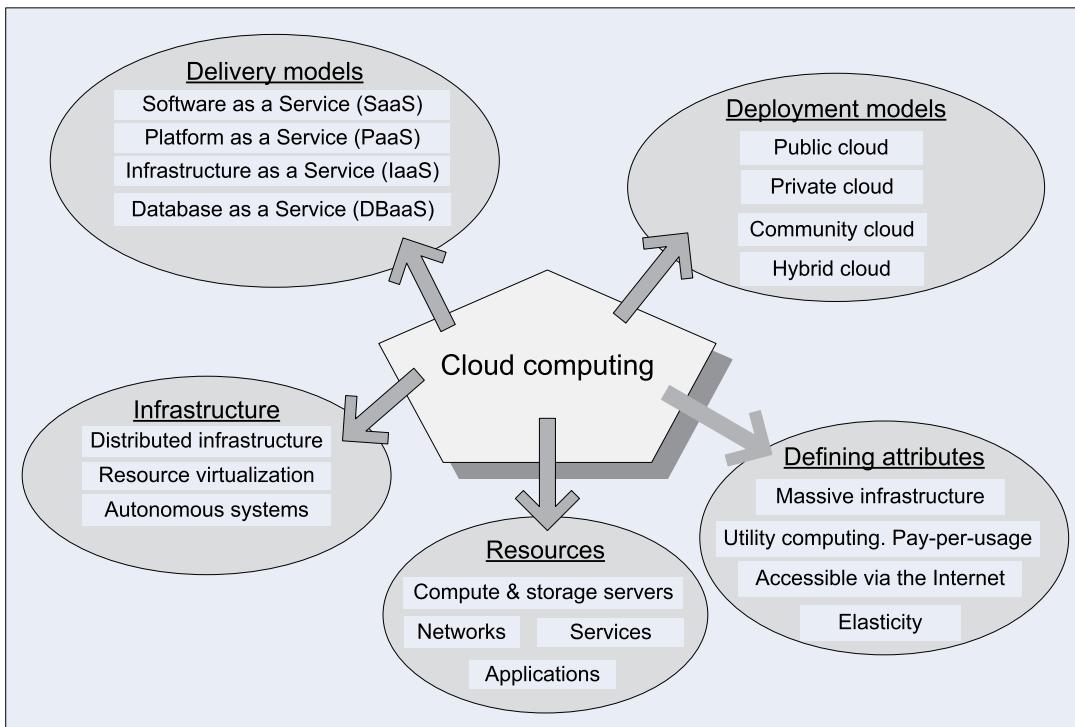
The thinking in software engineering has also evolved, and new models have emerged. A software architecture and a software design pattern, namely, the *three-tier model*, has emerged. Its components are discussed next. The *presentation tier* is the topmost level of the application; typically, it runs on a desktop or laptop, uses a standard graphical user interface, and displays information related to services. The *application/logic tier* controls the functionality of an application, may consist of one or more separate modules running on a workstation or application server, and may be multi-tiered itself. The *data tier* controls the servers where the information is stored; it runs a relational database management system on a database server or a mainframe and contains the computer data storage logic; it keeps data independent from application servers or processing logic and improves scalability and performance. Any of the tiers can be replaced independently: for example, a change of operating system in the presentation tier would only affect the user interface code.

Once the technological elements were in place, it was only a matter of time until the economical advantages of cloud computing became apparent. Due to the economy of large-scale data centers, centers with more than 50 000 systems are more economical to operate than medium-sized centers that have around 1 000 systems. Large data centers equipped with commodity computers experience a five to seven times decrease in resource consumption, including energy, compared to medium-sized data centers [32].

Several factors contribute to the success of cloud computing: (i) technological advances; (ii) a realistic system model; (iii) user convenience, and (iv) cost. A non-exhaustive list of cloud computing advantages over previous attempts to network centric computing includes:

- Cloud computing is in a better position to exploit recent advances in software, networking, storage, and processor technologies. Cloud computing is promoted by large IT companies where technological developments take place; the companies have a vested interest to promote the new technologies.
- A cloud consists of hardware and software resources in a single administrative domain where resource management, fault-tolerance, and quality of service are less challenging than distributed computing with resources in multiple administrative domains.
- Cloud computing is focused on enterprise computing [159,164]; its adoption by industrial organizations, financial institutions, healthcare organizations, and so on, has a potentially huge economic impact.
- A cloud provides the illusion of infinite computing resources; computer cloud elasticity frees applications designers from the confinements of a single system.
- A cloud eliminates the need for up-front financial commitment, and it is based on a pay-as-you-go approach; this has the potential to attract new applications and new users for existing applications, fomenting a new era of industry-wide technological advancements.

The term “computer cloud” covers infrastructures of different sizes, with different management and different user populations. Several types of clouds can be identified: (i) *public cloud*—the infrastructure is made available to the general public or a large industry group and is owned by an organization selling cloud services; (ii) *private cloud*—the infrastructure is operated solely for an organization; (iii) *hybrid cloud*—the infrastructure is a composition of two or more clouds (e.g., public and private) that remain unique entities but are bound together by standardized or proprietary technology that enables data and application portability; and (iv) *community cloud*—the infrastructure is shared by several organizations and supports a specific community with shared concerns.

**FIGURE 1.1**

A summary view of cloud computing: delivery models, deployment models, defining attributes, and organization of cloud infrastructure, discussed in Chapter 2, and cloud resources, discussed in Chapters 4, 6, 7.

A private cloud could provide the computing resources needed for a large organization, e.g., a research institution, a university, or a corporation. The argument that a private cloud does not support utility computing is based on the observation that an organization has to invest in the infrastructure and a user of a private cloud does pay as it consumes resources [32]. Nevertheless, a private cloud could use the same hardware infrastructure as a public one; its security requirements will be different from those for a public cloud, and the software running on the cloud is likely to be restricted to a specific domain.

Cloud computing delivery models, deployment models, defining attributes, resources, and organizations of the infrastructure are summarized in Fig. 1.1. The defining attributes of the new philosophy for delivering computing services are:

- Cloud computing uses Internet technologies to offer elastic services, i.e., the ability of dynamically acquiring computing resources and supporting a variable workload. A cloud service provider maintains a massive infrastructure to support elastic services.
- Resources used for these services can be metered, and users can be charged only for the resources they used.
- Maintenance and security are ensured by service providers.

- Economy of scale allows service providers to operate more efficiently due to specialization and centralization.
- Cloud computing is cost-effective due to resource multiplexing; lower costs for the service provider are passed on to cloud users.
- The application data is stored closer to the site where it is used in a device- and location-independent manner; potentially, this data storage strategy increases reliability and security and, at the same time, lowers communication costs.

In spite of the technological breakthroughs that have made cloud computing feasible, there are still major obstacles for this new technology; these obstacles provide opportunities for research. We list a few of the most obvious obstacles:

- Availability of service; what happens when the service provider cannot deliver? Can a large company such as GM move its IT to the cloud and have assurances that its activity will not be negatively affected by cloud overload? A partial answer to this question is provided by Service Level Agreements (SLA)s discussed in Section 2.8.
- Performance unpredictability, unavoidable due to resource sharing.
- Elasticity, the ability to scale up and down quickly and *overprovisioning*, maintaining pools of resources considerably larger than the average need.
- Vendor lock-in; once a customer is hooked to one provider, it is hard to move to another. The standardization efforts at NIST attempt to address this problem.
- Data confidentiality and auditability, a serious concern analyzed in Chapter 8.
- Data transfer bottlenecks critical for data-intensive applications. Transferring 1 TB of data on a 1 Mbps network takes 8 000 000 seconds or about 10 days; it is faster and cheaper to use courier service and send data recorded on some media than to send it over the network. High-speed networks will alleviate this problem, e.g., a 1 Gbps network would reduce this time to 8 000 seconds, or slightly more than two hours.

Cloud computing is a technical and social reality and an emerging technology. At this time, one can only speculate how the infrastructure for this new paradigm will evolve and what applications will migrate to it. The economical, social, ethical, and legal implications of this shift in technology, when the users rely on services provided by large data centers and store private data and software on systems they do not control, are likely to be significant.

Scientific and engineering applications, deep learning, data mining, computational financing, and gaming and social networking, as well as many other computational and data-intensive activities, can benefit from cloud computing. A broad range of data from the results of high-energy physics experiments to financial or enterprise management data, to personal data such as photos, videos, and movies, can be stored on the cloud.

The obvious advantage of network-centric content is the accessibility of information from any site where one can connect to the Internet. Clearly, information stored on a cloud can be shared easily, but this approach also raises major concerns: Is the information safe and secure? Is it accessible when we need it? Do we still own it?

In the near future, the focus of cloud computing is expected to shift from building the infrastructure, today's main front of competition among the vendors, to the application domain. This shift in focus is reflected by Google's strategy to build a dedicated cloud for government organizations in the United

States. The Internet made cloud computing possible: We could not even dream of using computing and storage resources from distant data centers without fast communication. The evolution of cloud computing is organically tied to the future of the Internet. The Internet of Things (IoT) has already planted some of its early seeds in computer clouds, and Amazon already offers services such as Lambda and Kinesis.

In a discussion of the technology trends, Jim Gray emphasized that the cost of communication in a wide area network has decreased dramatically and will continue to do so. Thus, it makes economical sense to store the data near the application [207], in other words, to store it in the cloud where the application runs. This insight leads us to believe that several new classes of cloud computing applications could emerge in the next few years [32].

The excitement due to cloud computing has translated into a flurry of publications. In this book, we attempt to sift through the large volume of information and dissect the main ideas related to cloud computing. We first discuss applications of cloud computing and then analyze the infrastructure for cloud computing. Several decades of research in parallel and distributed computing have paved the way for cloud computing. Through the years, we have discovered the challenges posed by the implementation of parallel and distributed systems and the ways to address some of them while avoiding the others.

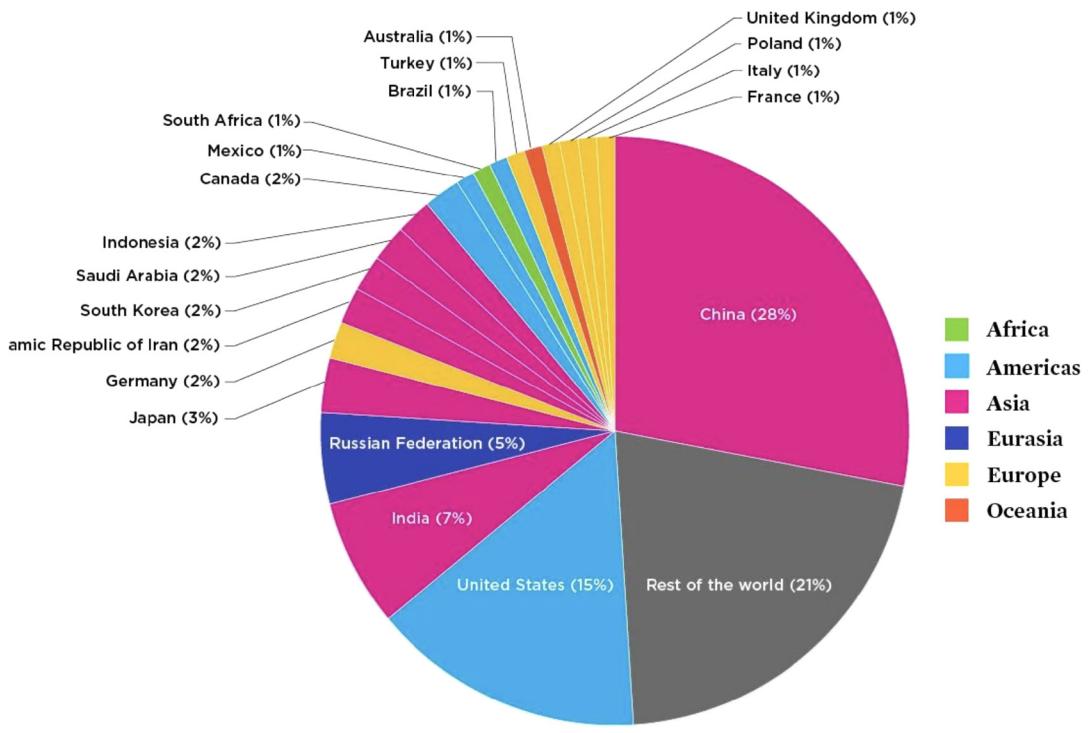
1.2 Energy use and ecological impact of cloud computing

The discussion of cloud infrastructure cannot be concluded without an analysis of the energy used for cloud computing and its impact on the environment. The energy consumption required by different types of human activities is partially responsible for greenhouse-gas emissions and has grave implications on climate change. Fig. 1.2 shows that industrialized countries, including the US and China, contribute significantly to CO₂ emissions according to data published in August 2020 by the Union of Concerned Scientists (see <https://www.ucsusa.org/resources/each-countrys-share-co2-emissions>). The data centers and the IT industries contribute to the large CO₂ footprint of these countries.

Reduction of energy consumption and thus of the carbon footprint of cloud-related activities is increasingly more important for society. Indeed, more and more applications run on clouds, and cloud computing uses more energy than many other human-related activities. Reduction of the carbon footprint can only be achieved through a comprehensive set of technical efforts. The hardware of the cloud infrastructure has to be refreshed periodically, and new and more energy efficient technologies have to be adopted; the resource management software has to pay more attention to energy optimization.

International Data Corporation (IDC) forecasts that continued adoption of cloud computing could prevent the emission of more than 1 billion metric tons of carbon dioxide (CO₂) from 2021 through 2024. This projection is based on the assumption that 60% of data centers will adopt the technology and processes underlying more sustainable “smarter” data centers by 2024 (<https://www.environmentalleader.com/2021/03/report-continued-adoption-of-cloud-computing-could-prevent-emission-of-1-billion-tons-of-co2/>). Several factors could contribute to lowering CO₂ emissions by cloud data centers: (i) shifting to cleaner sources of energy; (ii) reducing wasted energy and having more energy spent on running the IT equipment rather than on cooling; and (iii) delivering IT service wherever needed, by shifting workloads to enable greater use of renewable resources.

A 2020 paper published by the journal *Science* [336] supports this optimistic view in spite of the significant increase of data-center activity: “By 2018, global data center workloads and compute in-

**FIGURE 1.2**

The CO₂ contribution of various countries of the world.

stances had increased more than sixfold, whereas data center internet protocol (IP) traffic had increased by more than 10-fold. Data center storage capacity has also grown rapidly, increasing by an estimated factor of 25 over the same time period... In 2018, we estimated that global data center energy use rose to 205 TWh, or around 1% of global electricity consumption. This represents a 6% increase compared with 2010, whereas global data center compute instances increased by 550% over the same time period. Expressed as energy use per compute instance, the energy intensity of global data centers has decreased by 20% annually since 2010, a notable improvement.”

The website of the Center of Expertise for Energy Efficiency in Data Centers at Lawrence Berkeley National Laboratory, <https://datacenters.lbl.gov/resources/recalibrating-global-data-center-energy>, provides a wealth of information, tools, technology, and training for data-center energy optimization.

1.3 Ethical issues in cloud computing

Cloud computing is based on a paradigm shift with profound implications for computing ethics. The main elements of this shift are: (i) Control is relinquished to third party services; (ii) data is stored on

multiple sites administered by several organizations; and (iii) multiple services interoperate across the network. The complex structure of cloud services makes it difficult to determine who is responsible for each realm of action. Many entities contribute to an action with undesirable consequences, and no one can be held responsible. As a result of de-perimeterization, “not only the border of the organizations IT infrastructure blurs, also the border of the accountability becomes less clear” [477].

Unauthorized access, data corruption, infrastructure failure, and service unavailability are some of the risks related to relinquishing control to third-party services; moreover, whenever a problem occurs, it is difficult to identify the source and the entity causing it. Systems can span the boundaries of multiple organizations and cross the security borders, a process called *de-perimeterization*.

Ubiquitous and unlimited data sharing and storage among organizations test the self-determination of information, the right or ability of individuals to exercise personal control over the collection, and the use and the disclosure of their personal data by others. This tests the confidence and trust in today’s evolving information society. Identity fraud and theft are made possible by unauthorized access to personal data in circulation and by new forms of dissemination through social networks. All these factors can also pose a danger to cloud computing.

Cloud service providers have already collected petabytes of sensitive personal information stored by data centers around the world. The acceptance of cloud computing will be determined by the effort dedicated by the CSPs and the countries where the data centers are located to ensure privacy. Privacy is affected by cultural differences: While some cultures favor privacy, other cultures emphasize community, and this leads to an ambivalent attitude towards privacy on the Internet, which is a global system.

The question of what can be done proactively about the ethics of cloud computing does not have easy answers because many undesirable phenomena in cloud computing will only become apparent over time. However, the need for rules and regulations for the governance of cloud computing are obvious. Governance means the manner in which something is governed or regulated, the method of management, and the system of regulations. Explicit attention to ethics must be paid by governmental organizations providing research funding; private companies are less constrained by ethics oversight, and governance arrangements are more conducive to profit generation.

Accountability is a necessary ingredient of cloud computing; adequate information about how data is handled within the cloud and about allocation of responsibility are key elements for enforcing ethics rules in cloud computing. Recorded evidence enables us to assign responsibility, but there can be tension between privacy and accountability, and it is important to establish what is being recorded and who has access to the records. Unwanted dependency on a cloud service provider, the so-called *vendor lock-in*, is a serious concern, and the current standardization efforts at NIST attempt to address this problem. Another concern for the users is a future in which only a handful of companies dominate the market and dictate prices and policies.

1.4 Factors affecting cloud service availability

Clouds are affected by malicious attacks and failures of the infrastructure, e.g., power failures. Such events can affect the Internet domain name servers and prevent access to a cloud or can directly affect the clouds. For example, an attack at Akamai on June 15, 2004, caused a domain name outage and a major blackout that affected Google, Yahoo, and many other sites. In May 2009, Google was the target

of a serious denial of service (DNS) attack that took down services like Google News, and Gmail for several days.

Lightning caused a prolonged down time at Amazon on June 29–30, 2012; the AWS cloud in the East region of the US, which consists of ten data centers across four availability zones, was initially troubled by utility power fluctuations, probably caused by an electrical storm. Availability zones are locations within data-center regions where public cloud services originate and operate. A June 29, 2012, storm on the East Coast took down some of Virginia-based Amazon facilities and affected companies using systems exclusively in this region. Instagram, a photo sharing service, was one of the victims of this outage.

The recovery from the failure took a very long time and exposed a range of problems. For example, one of the ten centers failed to switch to backup generators before exhausting the power that could be supplied by UPS units. AWS uses “control planes” to enable users to switch to resources in a different region, and this software component also failed. The booting process was faulty and extended the time to restart EC2 and EBS services. Another critical problem was a bug in the Elastic Load Balancer (ELB), which is used to route traffic to servers with available capacity. A similar bug affected the recovery process of Relational Database Service (RDS). This event brought to light “hidden” problems that occur only under specific circumstances.

The stability risks due to interacting services are discussed in [182]. A cloud application provider, a cloud storage provider, and a network provider could implement different policies, and the unpredictable interactions between load-balancing and other reactive mechanisms could lead to dynamic instabilities. The unintended coupling of independent controllers that manage the load, the power consumption, and the elements of the infrastructure could lead to undesirable feedback and instability similar to the one experienced by the policy-based routing in the Internet BGP (Border Gateway Protocol).

For example, the load balancer of an application provider could interact with the power optimizer of the infrastructure provider. Some of these couplings may only become manifest under extreme conditions and be very hard to detect under normal operating condition but could have disastrous consequences when the system attempts to recover from a hard failure, as in the case of the AWS 2012 failure.

Clustering resources in data centers located in different geographical areas lowers the probability of catastrophic failures. This geographic dispersion of resources could have additional positive side effects, such as reduction of communication traffic, lowering energy costs by dispatching the computations to sites where the electric energy is cheaper, and improving performance by an intelligent and efficient load-balancing strategy.

1.5 Network-centric computing and network-centric content

Network-centric computing and network-centric content imply that data processing and data storage can be done on computers accessed through the Internet. The term *content* refers to any type or volume of media, be it static or dynamic, monolithic or modular, live or stored, produced by aggregation or mixed. *Information* is the result of functions applied to content.

The content should be treated as having meaningful semantic connotations rather than a string of bytes; the focus will be on the information that can be extracted by content mining when users request

named data and content providers publish data objects. Content-centric routing will enable users to fetch the desired data from the most suitable location in terms of network latency or download time. In turn, the creation and consumption of audio and visual content is likely to transform the Internet to support increased quality in terms of resolution, frame rate, color depth, stereoscopic information, etc.

There are also some challenges, such as providing secure services for content manipulation, ensuring global rights-management, control over unsuitable content, and reputation management. Network-centric computing and network-centric contents share a number of characteristics:

1. Most applications are data intensive. Data analytics enable enterprises to optimize their operations; computer simulation is a powerful tool for scientific research in virtually all areas of science, from physics, biology, and chemistry to archeology. The widespread use of sensors contributes to the increase of the volume of data. Artificial Intelligence (AI) and Machine Learning (ML) algorithms require massive amounts of data.
2. Computing and communication resources (CPU cycles, storage, network bandwidth) are shared, and resources can be aggregated to support data-intensive applications. Multiplexing leads to a higher resource utilization: When multiple applications share a system, their peak demands for resources are not synchronized, and average system utilization increases.
3. Data sharing facilitates collaborative activities. Indeed, many applications in science and engineering, as well as industrial, financial, and governmental applications, require multiple types of analysis of shared-data sets and multiple decisions carried out by groups scattered around the globe. Open software development sites are another example of such collaborative activities.
4. Virtually all applications are network intensive. Indeed, transferring large volumes of data requires high bandwidth networks. Parallel computing, computation steering, and data streaming are examples of applications that can only run efficiently on low-latency networks. Computation steering in numerical simulation means to interactively guide a computational experiment towards a region of interest.
5. The systems are accessed using *thin clients* running on systems with limited resources.

There are sources of concern regarding the paradigm shift from locally owned resources to network-centric computing: (i) the management of large pools of resources poses new challenges because such systems are vulnerable to malicious attacks that can affect a large population of users; (ii) large-scale systems are affected by phenomena characteristic to complex systems such as phase transitions when a relatively small change in the environment can lead to an undesirable system state [332]; (iii) ensuring Quality of Service (QoS) guarantees is extremely challenging because perfect performance isolation and the ability to accommodate workloads with very large peak-to-average ratios are elusive; and (iv) data sharing poses not only security and privacy challenges but also requires mechanisms for access control for authorized users and for detailed logs of the history of data changes.

The cloud ecosystem

2

The statement that cloud computing is a disruptive technology can be easily supported by the arguments presented in the final section of this chapter where cloud computing's impact on enterprise computing is quantified in billions, possibly trillions, of dollars. The argument that cloud computing is a major contributor to the large CO₂ footprint of the IT industry is more subtle,¹ as we have seen in Section 1.2.

Even more subtle is the transformative effect of cloud computing on virtually all areas of computing, including algorithms, computer software and hardware, database systems, networking, etc. This is the main theme of this text, discussed in the next eleven chapters, which track how concepts, ideas, and technologies from the precloud computing era have changed to accommodate the qualitative and quantitative challenges brought about by the cloud ecosystem on hardware, software, and applications.

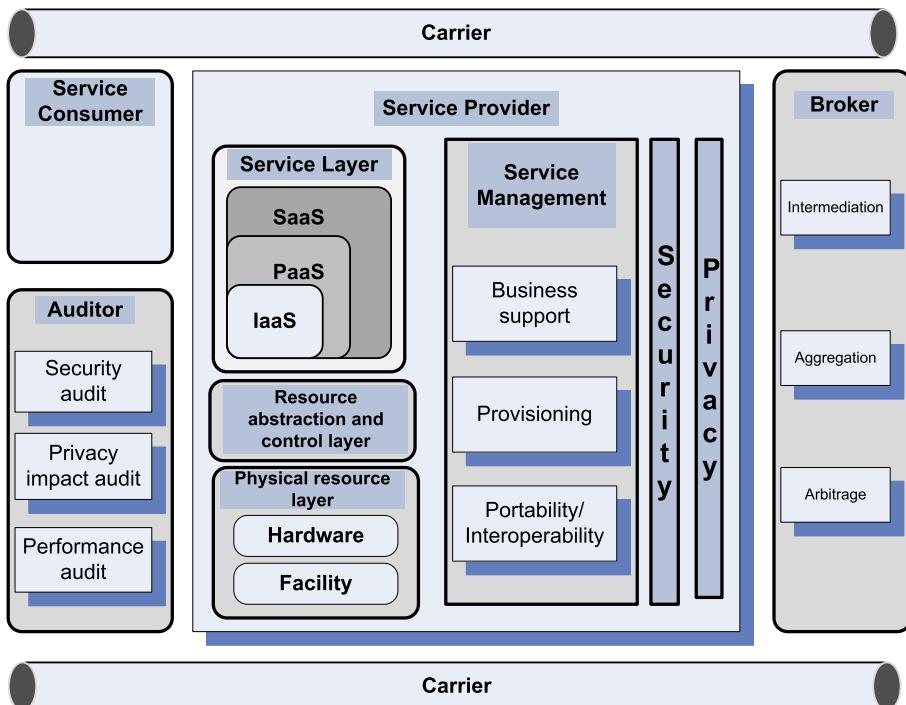
A case in point is the evolution of the *load balancing* concept used for three decades or longer to suggest that, when you have a system with a large number of components, it is best to divide the load if you can and then assign each unit about the same amount of work. The new twist in cloud computing is that load balancing should attempt to minimize energy consumption, and this means concentrating the workload on the smallest number of servers, while switching the others to an energy-saving mode.

Previously, computer architecture supported only the privileged kernel mode and the non-privileged user mode, but system resources in a virtual environment are managed by a hypervisor rather than one operating system. Consequently, the hardware needs to support different levels of protection. Moreover, the hardware must provide a trusted execution mode because clouds are target-rich environments.

The *scale of the cloud computing infrastructure* has profound implications on cloud resource management forced to satisfy conflicting requirements; on one hand, the system must be elastic, i.e., able to accommodate sudden spikes of workload and, on the other hand, minimize energy consumption. All these complications must be hidden from users who should be presented with simple and easy-to-use interfaces for low-cost, high-level services. This is the subject of this introductory chapter which presents the cloud ecosystem as of mid 2021.

The major players in this ecosystem are Amazon, Google, and Microsoft and now IBM, which is poised to challenge the dominance of the first three on this list. The chapter starts with an overview of the cloud ecosystem and an in-depth discussion of cloud delivery models and services in Section 2.1, followed by the analysis of cloud computing at Amazon, Google, Microsoft, and IBM in Sections 2.2, 2.3, 2.4, and 2.5, respectively. Sections 2.6 and 2.7 cover the pervasive issue of vendor lock-in and the prospects of cloud interoperability. The presentation of Service Level Agreements (SLAs) in Section 2.8 is followed by a discussion of responsibility sharing between cloud users and cloud service providers in

¹ The IT industry may one day end up competing with the transportation industry as far as its CO₂ footprint is concerned, unless the Bitcoind craze is stopped or at least limited, see Section 3.18.

**FIGURE 2.1**

Entities involved in service-oriented computing and, in particular, in cloud computing according to NIST. The carrier provides connectivity between service providers, service consumers, brokers, and auditors.

Section 2.9. User experience is analyzed in Section 2.10, while Section 2.11 covers software licensing. Major challenges faced by cloud computing are discussed in Section 2.12. The cloud is changing the enterprise ecosystem at a stunning speed as presented in Section 2.13.

2.1 Cloud computing delivery models and services

According to NIST reference model in Fig. 2.1 [366], the entities involved in cloud computing are: *service consumer*—maintains a business relationship with, and uses service from, service providers; *service provider*—responsible for making a service available to service consumers; *carrier*—intermediary providing connectivity and transport of cloud services between providers and consumers; *broker*—manages the use, performance, and delivery of cloud services, and negotiates relationships between providers and consumers; *auditor*—a party that can conduct independent assessment of cloud services, information system operations, performance, and security of the cloud implementation.

An *audit* is a systematic evaluation of a cloud system that measures how well it conforms to a set of established criteria. For example, a security audit evaluates cloud security, and a privacy-impact audit evaluates cloud privacy assurance, while a performance audit evaluates cloud performance. Fig. 2.1 shows activities supporting cloud delivery models: (i) service management and provisioning including: virtualization, service provisioning, the call center, operations management, systems management, QoS management, billing, accounting, asset management, SLA management, technical support, and backups; (ii) security management including: ID and authentication, certification and accreditation, intrusion prevention, intrusion detection, virus protection, cryptography, physical security, incident response, access control, audit and trails, and firewalls; (iii) customer services such as: customer assistance and on-line help, subscriptions, business intelligence, reporting, customer preferences, and personalization; and (iv) integration services including data management and development.

Cloud service providers (CSPs) support one or more cloud delivery models: SaaS (Software as a Service), PaaS (Platform as a Service), IaaS (Infrastructure as a Service), and DBaaS (Database as a Service). SaaS is a centrally hosted service based on a subscription model for licensed software. PaaS provides a computing platform and applications, without the need to manage the infrastructure required to build and run the applications. IaaS is the model where the CSP provides a menu of resource bundles and the users can choose the operating system and manage the resources in the bundle. DBaaS enables users to setup and operate several databases using a common set of abstractions, without the need to know the exact implementations of those abstractions for individual databases. Amazon is a pioneer in IaaS, Google's efforts are focused on SaaS and PaaS delivery models, and Microsoft is mostly involved in PaaS. Amazon, Oracle, and many other CSPs offer DBaaS services. Oracle Cloud is based on Java, SQL standards, and software systems such as Exadata, Exalogic, WebLogic, and Oracle Database.

Software as a Service. A wide range of stationary and mobile devices enable a large population of clients to access services provided by the applications using a thin client interface such as a web browser (e.g., web-based email). The users of services do not manage or control the underlying cloud infrastructure including network, servers, operating systems, storage, or even individual application capabilities, with the possible exception of limited user-specific application configuration settings. Examples of SaaS include: enterprise services, such as workflow management, group-ware and collaborative, supply chain, communications, digital signature, customer relationship management (CR), desktop software, financial management, geo-spatial, and search.

The SaaS delivery model is not suitable for applications requiring real-time response or those in which data is not allowed to be hosted externally. The most likely candidates for *SaaS* are applications with these characteristics: (a) Many competitors use the same product, such as Email; (b) there is a significant peak in demand periodically, such as in billing and payroll applications; (c) there is a need for the web or mobile access, such as mobile sales-management software; (d) there is only a short-term need, such as collaborative software for a project.

Platform as a Service. PaaS offers the capability to deploy consumer-created or acquired applications using programming languages and tools supported by the provider. The user does not manage or control the underlying cloud infrastructure including the network, servers, operating systems, or storage. The user has control over the deployed applications and, possibly, applications hosting environment configurations. Such services include: session management, device integration, sandboxes, instrumentation and testing, contents management, knowledge management, and Universal Description, Discovery, and Integration (UDDI), a platform-independent, Extensible Markup Language (XML)-based registry providing a mechanism to register and locate web-service applications.

PaaS is not particularly useful when the application must be portable, when proprietary programming languages are used, or when the underlaying hardware and software must be customized to improve the performance of the application. Its major application areas are in software development when multiple developers and users collaborate and the deployment and testing services should be automated.

Infrastructure as a Service. IaaS has the capability to provide processing, storage, networks, and other fundamental computing resources; the consumer is able to deploy and run arbitrary software, which can include operating systems and applications. The consumer does not manage or control the underlying cloud infrastructure but has control over operating systems, storage, deployed applications, and possibly limited control of some networking components, e.g., host firewalls. Services offered by this delivery model include: server hosting, web servers, storage, computing hardware, operating systems, virtual instances, load balancing, Internet access, and bandwidth provisioning.

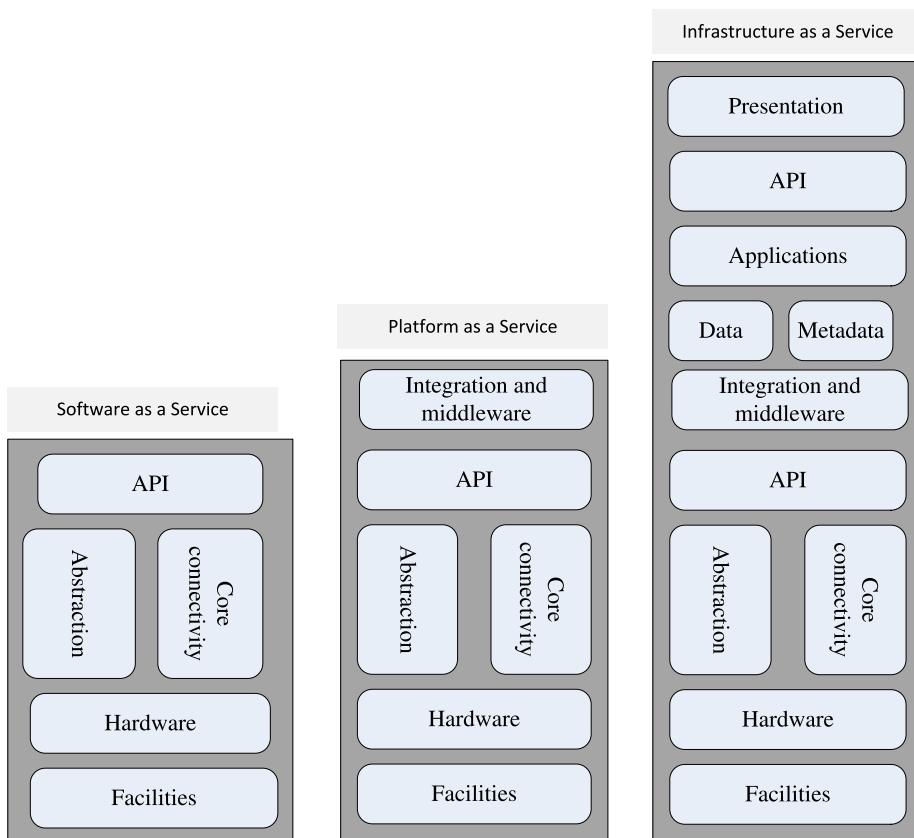
The IaaS cloud computing delivery model has a number of characteristics: The resources are distributed and support dynamic scaling; it is based on a utility pricing model and variable cost; and the hardware is shared among multiple users. This cloud delivery model is particularly useful when the demand is volatile, a new business needs computing resources, and it does not want to invest in a computing infrastructure, or when an organization is expanding rapidly.

Fig. 2.2 shows the degrees of freedom of cloud users and the interaction with cloud infrastructure, which vary from extremely limited for SaaS, to modest for PaaS, and significant for IaaS [122]. Table 2.1 shows that SaaS cloud service providers control not only the cloud infrastructure and the hypervisor, but also the network traffic, the operating system, and the applications. IaaS CSPs only provide the cloud infrastructure and the hypervisor.

Database as a Service. DBaaS is a cloud service where the database runs on the service provider physical infrastructure. Compared with on-site physical server and storage architecture, a cloud database service offers distinct advantages: instantaneous scalability, performance guarantees, specialized expertise, latest technology, failover support, and declining pricing. Most relevant features of the DBaaS model are:

1. Self-service—service provisioning without major deployment or configuration and without performance and cost penalties.
2. Device and location-independent abstract database resources without concern for hardware utilization.
3. Elasticity and scalability—automated and dynamic scaling.
4. Pay-as-you-go model—metered use of resources and cost reflecting the resources used.
5. Agility—the applications adapt seamlessly to new technology or additional requirements.

Cloud DBaaS uses a layered architecture. The *user interface layer* supports access to the service via the Internet. The *application layer* is used to access software services and storage space. The *database layer* provides efficient and reliable database service; it saves time for querying and loading data by reusing the query statements residing in the storage. *Data storage layer* encrypts the data when stored without user involvement; backup management and disk monitoring are also provided by this layer. Multi-tenancy is an integral part of the DBaaS model. In spite of its advantages, multi-tenancy poses resource management and security challenges, as discussed in Section 8.7.

**FIGURE 2.2**

Software and cloud infrastructure control versus ease of use of the service. SaaS allows no infrastructure control and does not require any sophistication to use applications supplied by CSP. PaaS supports consumer-created or acquired applications using CSP-provided tools. IaaS empowers advanced developers and uses to chose the operating system, run arbitrary software, choose the type of instances, and manage instance resources.

Table 2.1 The shared security responsibility (SSR) model between cloud service providers (in bold**) and cloud users for the three cloud delivery models.**

SaaS	PaaS	IaaS
User Access	User Access	User Access
Data	Data	Data
Applications	Applications	Applications
Operating Systems	Operating Systems	Operating Systems
Network Traffic	Network Traffic	Network Traffic
Hypervisor	Hypervisor	Hypervisor
Infrastructure	Infrastructure	Infrastructure

This discussion shows that a service-oriented architecture involves multiple subsystems and complex interactions among these subsystems. Individual subsystems can be layered; for example, we see that the service layer sits on top of a resource abstraction layer that controls the physical resource layer.

Cloud delivery models will continue to coexist for the foreseeable future. Services based on SaaS will probably be increasingly more popular because they are more accessible to lay people, while services based on the IaaS will be the domain of computer-savvy individuals, large organizations, and the government. If the standardization effort succeeds, then we may see IaaS designed to migrate from one infrastructure to another and overcome the concerns related to vendor lock-in. The popularity of DBaaS services is likely to grow.

Infrastructure as Code (IaC). It is always beneficial to automate provisioning of cloud infrastructure, i.e., manage servers, operating systems, database connections, storage, and other infrastructure elements using a high-level descriptive language. This is precisely the function of IaC that takes advantage of virtualization and cloud native development allowing developers to provision their own virtual servers or containers on demand. A cloud native application consists of discrete, reusable components, *microservices*, acting as building blocks, often packaged in containers, and designed to integrate into any cloud environment. Microservices can be independently scaled, continuously improved, and quickly iterated through automation and orchestration processes.

IaC advantages are obvious, namely, improved consistency, more efficient development, lower costs, shorter time to production, and a more secure environment. For the functional or declarative IaC approach, a skilled administrator specifies the desired final state of the infrastructure, and the IaC software handles the rest—spinning up the virtual machine (VM) or container, installing and configuring the necessary software, resolving system and software interdependencies, and managing versioning. The procedural IaC approach automates provisioning the infrastructure one step at a time.

The most used IaC tools are *Ansible* and *Terraform*. Ansible is an open-source IT automation engine that is used to improve the scalability, consistency, and reliability of the IT environment. Ansible supports: (i) provisioning, i.e., setting up the servers; (ii) configuration management, i.e., changing the configuration of an application, OS, or device, start and stop services, installing or updating applications, and implementing a security policy, of other configuration tasks; and (iii) application deployment automates the deployment of applications using DevOps, a philosophy of efficient software development, deployment, and operation.

Terraform automates resource management across multiple providers, regardless of where physical servers, DNS servers, or databases reside and provisions applications written in any language. The first step in *Terraform* is to create an execution plan; then, the system generates a graph of all resources and parallelizes creation and modification of any non-dependent resources. A number of use cases are described on <https://www.terraform.io/intro/use-cases.html> including multi-tier applications, self-service clusters, SDN (Software Defined Networks), resource schedulers, and multi-cloud environments.

2.2 Amazon Web Services

Amazon has changed the face of computing over recent decades. First, it installed a powerful computing infrastructure to sustain its core business, selling online a variety of goods ranging from books and CDs to gourmet foods and home appliances. Then, the company discovered that this infrastructure can be

further extended to provide affordable and easy-to-use resources for enterprise computing, as well as computing for the masses.

Amazon Web Services (AWS), based on the IaaS delivery model, offers an infrastructure consisting of computation and storage servers interconnected by high-speed networks and supports a set of services to access these resources. Businesses in more than 200 countries use AWS. A significant number of corporations and start-ups take advantage of Amazon cloud computing infrastructure.

AWS computing, storage, and communication services. Amazon announced a limited public beta release of its Elastic Computing platform called EC2 in 2006. AWS offered 24 services in 2008, 48 in 2009, 61 in 2010, 82 in 2011, 159 in 2012, and 280 in 2013; 449 new services and major features were released in 2014. AWS services reflect leading-edge technological developments: For example, machine learning and quantum simulation services are now supported, see Sections 13.2 and 13.3.

Elastic Compute Cloud (EC2) is a web service with a simple interface for launching instances of an application under several operating systems, such as several Linux distributions, OpenSolaris, FreeBSD, NetBSD, and Microsoft Windows Server 2008, 2008 R2, 2012, 2012 R2, 2016, and 2019. EC2 is based on the virtualization strategy discussed in detail in Chapter 5. For more than ten years AWS used the Xen hypervisor discussed in Section 5.8 and in 2017 replaced it with KVM discussed in Section 5.7.

An example of EC2 is a virtual server; the user chooses the region and the availability zone where this virtual server should be placed and also selects from a menu of instance types the one that provides the resources, CPU cycles, main memory, secondary storage, communication, and I/O bandwidth needed by the application.

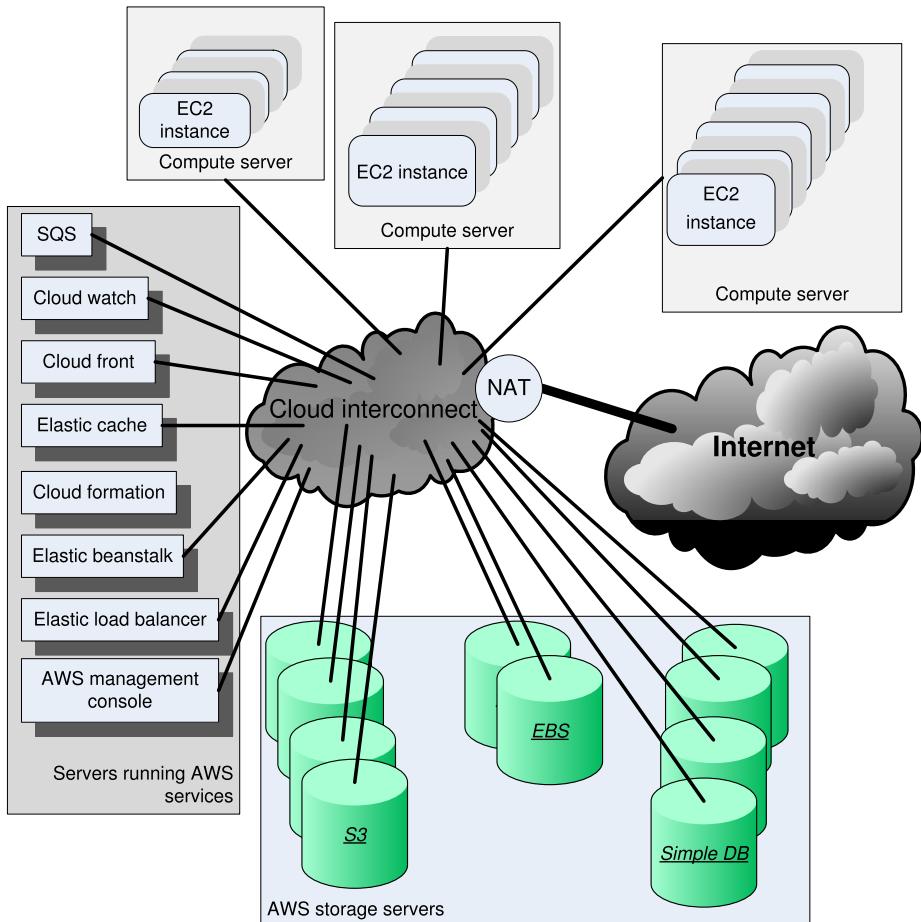
When launched, an instance is provided with a *DNS name*; this name maps to a *private IP address* for internal communication within the internal EC2 communication network and a *public IP address* for communication outside the internal Amazon network, e.g., for communication with the user that launched the instance. Network Address Translation (NAT) maps external IP addresses to internal ones.

The public IP address is assigned for the lifetime of an instance and is returned to the pool of available public IP addresses when the instance is either stopped or terminated. An instance can request an *elastic IP address*, rather than a public IP address. The elastic IP address is a static public IP address allocated to an instance from the available pool of the availability zone; an elastic IP address is not released when the instance is stopped or terminated and must be released when no longer needed.

An instance is created either from a predefined Amazon Machine Image (AMI) digitally signed and stored in S3, or from a user-defined image. The image includes the operating system, the run-time environment, the libraries, and the application desired by the user. AMI images create an exact copy of the original image but without configuration-dependent information such as *hostname* or MAC address. A user can: (i) launch an instance from an existing AMI and terminate an instance; (ii) start and stop an instance; (iii) create a new image; (iv) add tags to identify an image; and (v) reboot an instance.

An instance specifies the maximum amount of resources available to an application and the interface for that instance, as well as the cost per hour. A server may run multiple virtual machines or instances started by one or more users; an instance may use storage services, S3, EBS, and Simple DB, as well as other services provided by AWS, see Fig. 2.3.

A user can interact with EC2 using a set of SOAP messages, see Section 11.1, and can list available AMI images, boot an instance from an image, terminate an image, display the running instances of a user, display console output, etc. The user has root access to each instance in the elastic and secure computing environment of EC2. The instances can be placed in multiple locations in different regions and availability zones.

**FIGURE 2.3**

Availability zone—a cloud interconnect supports high-speed communication among compute and storage servers in the zone. Communication with servers in other availability zones and with cloud users is supported via a Network Address Translation (NAT), which maps external IP addresses to internal ones.

EC2 allows the import of Virtual Machine (VM) images from the user environment to an instance through a facility called *VM import*. It also distributes automatically the incoming application traffic among multiple instances using the *elastic load balancing* facility. EC2 associates an *elastic IP address* with an account; this mechanism allows a user to mask the failure of an instance and re-map a public IP address to any instance of the account, without the need to interact with the software support team.

Simple Storage System (S3) is a storage service designed to store large objects. It supports a minimal set of functions: write, read, and delete. S3 allows an application to handle an unlimited number of

objects ranging in size from one byte to five terabytes. An object is stored in a *bucket* and retrieved via a unique developer-assigned key; a bucket can be stored in a region selected by the user.

S3 maintains for each object: the name, modification time, an access control list, and up to four kilobytes of user-defined metadata; object names are global. Authentication mechanisms ensure that data is kept secure; objects can be made public, and rights can be granted to other users. S3 supports PUT, GET, and DELETE primitives to manipulate objects, but it does not support primitives to copy, rename, or move an object from one bucket to another. Appending to an object requires a read followed by a write of the entire object.

S3 computes MD5² for every object written and returns it in a field called *ETag*. A user is expected to compute the MD5 of an object stored or written and compare this with the *ETag*; if the two values do not match, then the object was corrupted during transmission or storage. The S3 SLA guarantees reliability. S3 uses standards-based REST and SOAP interfaces, see Section 11.1; the default download protocol is HTTP, but BitTorrent³ protocol interface is also provided at lower cost for large-scale distributions.

Elastic Block Store (EBS) provides persistent block-level storage volumes for use with EC2 instances. A volume appears to an application as a raw, unformatted, and reliable physical disk; the size of the storage volumes ranges from one gigabyte to one terabyte. The volumes are grouped together in availability zones and are automatically replicated in each zone. An EC2 instance may mount multiple volumes, but a volume cannot be shared among multiple instances. EBS supports the creation of snapshots of the volumes attached to an instance and then uses them to restart an instance. The storage strategy provided by EBS is suitable for database applications, file systems, and applications using raw data devices.

Simple DB is a non-relational data store that allows developers to store and query data items via web services requests; it supports store and query functions traditionally provided only by relational databases. Simple DB creates multiple geographically distributed copies of each data item and supports high-performance web applications; at the same time, it manages automatically the infrastructure provisioning, hardware and software maintenance, replication and indexing of data items, and performance tuning.

Simple Queue Service (SQS) is a hosted message queue. SQS is a system for supporting automated workflows; it allows multiple EC2 instances to coordinate their activities by sending and receiving SQS messages. Any computer connected to the Internet can add or read messages without any installed software or special firewall configurations. Applications using SQS can run independently and asynchronously and do not need to be developed with the same technologies. A received message is “locked” during processing; if processing fails, the lock expires and the message is available again. The timeout for locking can be changed dynamically via the *ChangeMessageVisibility* operation. Developers can access SQS through standards-based SOAP and Query interfaces. Queues can be shared with other AWS accounts and anonymously; queue sharing can also be restricted by IP address and time of day. An example showing the use of SQS is presented in Section 11.6.

CloudWatch is a monitoring infrastructure used by application developers, users, and system administrators to collect and track metrics important for optimizing the performance of applications and for

² MD5 (Message-Digest Algorithm) is a widely used cryptographic hash function; it produces a 128-bit hash value and is used for checksums. SHA-i (Secure Hash Algorithm, $0 \leq i \leq 3$) is a family of cryptographic hash functions; SHA-1 is a 160-bit hash function resembling MD5.

³ BitTorrent is a peer-to-peer (P2P) communications protocol for file sharing.

increasing the efficiency of resource utilization. Without installing any software, a user can monitor approximately a dozen pre-selected metrics and then view graphs and statistics for these metrics. When launching an Amazon Machine Image (AMI), the user can start the CloudWatch and specify the type of monitoring. The basic monitoring is free of charge and collects data at five-minute intervals for up to ten metrics, while the detailed monitoring is subject to a charge and collects data at one minute interval. This service can also be used to monitor access latency to EBS volumes, available storage space for RDS DB instances, the number of messages in SQS, and other parameters of interest.

Virtual Private Cloud (VPC) provides a bridge between the existing IT infrastructure of an organization and the AWS cloud; the existing infrastructure is connected via a Virtual Private Network (VPN) to a set of isolated AWS computing resources. VPC allows existing management capabilities, such as security services, firewalls, and intrusion-detection systems, to operate seamlessly within the cloud.

Auto Scaling exploits cloud elasticity and provides automatic scaling of EC2 instances. The service supports: grouping of instances, monitoring of the instances in a group, and defining *triggers*, pairs of CloudWatch alarms and policies, which allow the size of the group to be scaled up or down. Typically, a maximum, a minimum, and the regular size of the group are specified. An Auto Scaling group consists of a set of instances described in a static fashion by launch configurations. When the group scales up, new instances are started using the parameters for the *runInstances* EC2 call provided by the launch configuration; when the group scales down, the instances with older launch configurations are terminated first. The monitoring function of Auto Scaling carries out health checks to enforce specified policies; for example, a user may specify a health check for elastic load balancing, and then Auto Scaling will terminate an instance exhibiting a low performance and start a new one. Triggers use CloudWatch alarms to detect events and then initiate specific actions; for example, a trigger could detect when the CPU utilization of the instances in the group goes above 90% and then scale up the group by starting new instances. Typically, triggers to scale up and down are specified for a group.

AWS services introduced in 2012 include: *Route 53*—a low-latency DNS service managing user DNS public records; *Elastic MapReduce (EMR)*—service supporting a hosted Hadoop running on EC2 and based on the MapReduce paradigm discussed in Section 11.5; *Simple Workflow Service (SWS)*—supports workflow management and allows scheduling, management of dependencies, and coordination of multiple EC2 instances; *ElastiCache*—a service enabling web applications to retrieve data from a managed in-memory caching system rather than a much slower disk-based database; *DynamoDB*—a scalable and low-latency fully managed NoSQL database service; *CloudFront*—a web service for content delivery; *Elastic Load Balancer*—a cloud service to automatically distribute the incoming requests across multiple instances of the application. Two of these services, the CloudFormation and the Elastic Beanstalk, are discussed next.

CloudFormation allows the creation of a stack describing the infrastructure for an application. The user creates a template, a text file formatted as in Javascript Object Notation (JSON), describing the resources, the configuration values, and the interconnection among these resources. The template can be parameterized to allow customization at run time e.g., to specify the types of instances, database port numbers, or RDS size.

Elastic Beanstalk interacts with other AWS services including EC2, S3, SNS, Elastic Load Balance, and AutoScaling. Elastic Beanstalk handles automatically the deployment, capacity provisioning, load balancing, auto-scaling, and application monitoring functions [488]. The service automatically scales the resources as required by the application, either up or down based on default Auto Scaling settings.

Some of the management functions provided by the service are: (i) deploying a new application version or rollback to a previous version; (ii) accessing results reported by CloudWatch monitoring service; (iii) providing email notifications when application status changes or application servers are added or removed; and (iv) accessing server log files without needing to login to the application servers.

Elastic Beanstalk service is available to developers using either a Java platform, the PHP server-side description language, or .NET framework. For example, a Java developer can create an application using an Integrated Development Environment, such as Eclipse, and package the code into Java Web Application Archive file of type “.war”. The “.war” file should then be uploaded to the Elastic Beanstalk using the Management Console and then deployed, and in a short time the application will be accessible via an URL.

Lambda service: cloud computing is associated for many with Big Data applications and long-lasting computations rather than with real-time applications and short bursts of computing triggered by some external events. A few years ago, most likely anticipating services relevant to the Internet of Things, AWS introduced a *serverless* computer service. In this case, applications are triggered by conditions and/or events specified by the end user. For example, the application may run for a brief period of time at midnight to check the daily energy consumption of an enterprise or may be activated weekly to check the sales of a chain or to turn on the alarm system of a home triggered by an event generated by the smartphone of the owner. In stark contrast to EC2, when customers are billed on an hourly basis, e.g., if a C4 instance is used for one hour and ten minutes the billing is for two hours, the *Lambda* service is billed for the actual time with a resolution of milliseconds. The service seems relatively easy to use. “First you create your function by uploading your code (or building it right in the Lambda console) and choosing the memory, timeout period, and AWS Identity and Access Management (IAM) role. Then, you specify the AWS resource to trigger the function, either a particular Amazon S3 bucket, Amazon DynamoDB table, or Amazon Kinesis stream. When the resource changes, Lambda will run a function, launch and manage the compute resources as needed in order to keep up with incoming requests.” Amazon Kinesis is a data-streaming platform.

Regions and availability zones. Amazon has 28+ data centers on several continents. An *availability zone* (AZ) is a data center with 50 000–80 000 servers using 25–30 MW of power. All regions have at least two availability zones interconnected by high-speed networks. Regions do not share resources and communicate through the Internet. Storage is automatically replicated within a region. S3 buckets are replicated within an availability zone and between the availability zones of a region, while *EBS* volumes are replicated only within the same availability zone. Critical applications should replicate important information in multiple regions to be able to function when servers in one region are unavailable due to catastrophic events. Most AWS services are available in all regions, though some are not.

The billing rates differ from one region to another and can be roughly grouped into four categories: low, medium, high, and very high billing. These rates are determined by the components of the operating costs including energy, communication, and maintenance costs. Region selection is motivated by the desire to minimize costs, reduce the communication latency, and increase reliability and security.

AWS Networking. Each AWS region has redundant *transit centers* connecting private links to other AWS regions, private links to AWS Direct Connect customers, and to the Internet through peering and paid transit. Most major regions are interconnected by private fiber channels, and this avoids peering issues, buffering problems, and capacity limitations that may occur on public links.

Peak traffic between availability zones of up to 25 Tbps is supported. The communication latency between availability zones is in the one- to two-milliseconds range. Communication latency between two servers has three components: (i) application \leftrightarrow guest OS \leftrightarrow hypervisor \leftrightarrow Network Interface (NIC)—in the milliseconds range; (ii) through the NIC—in the microseconds range; and (iii) over the fiber—in the nanoseconds range. Single Root I/O Virtualization virtualizes the NICs, and each guest gets its own virtual NIC.

Users have several choices to interact and manage AWS resources: (i) the Web Management Console, but not all options may be available in this mode; (ii) command-line tools; (iii) AWS SDK libraries and toolkits provided for several programming languages including Java, PHP,⁴ C#, and Obj C; (iv) raw REST requests as shown in Section 11.1.

Amazon Web Services Licensing Agreement (AWSLA) allows the cloud service provider to terminate service to any customer at any time for any reason and contains a covenant not to sue Amazon or its affiliates for any damages that might arise out of the use of AWS. AWSLA prohibits the use of “other information obtained through AWS for the purpose of direct marketing, spamming, contacting sellers or customers.” It prohibits AWS from being used to store any content that is “obscene, libelous, defamatory or otherwise malicious or harmful to any person or entity;” it also prohibits S3 from being used “in any way that is otherwise illegal or promotes illegal activities, including without limitation in any manner that might be discriminatory based on race, sex, religion, nationality, disability, sexual orientation or age.”

AWS Continuing Evolution. Amazon is one of the most important forces driving the spectacular evolution of cloud computing in recent years. AWS infrastructure has benefited from a wealth of new technologies. At this point in time, AWS is probably the most attractive and cost-effective cloud computing environment, not only for enterprise applications but also for computational science and engineering applications.

The massive effort to continually expand the hardware and the software of the AWS cloud infrastructure is astounding. Amazon has designed its own storage racks; such a rack holds 864 disk drives and weighs over a ton. The company has designed and built its own power substations. Three of its regions, US West (Oregon), AWS GovCloud (US), and EU (Frankfurt), are 100% carbon neutral.

AWS Nitro System, the platform for next generation EC2 instances, is a combination of dedicated hardware and a lightweight supervisor. While traditional hypervisors virtualize all system resources including, CPU, storage, and networking, the Nitro System allows AWS to “break apart those functions, offload them to dedicated hardware and software, and reduce costs by delivering practically all resources of a server to your instances.” AWS offers several classes of EC2 instances powered by different processors. Representative instances in each class, with the most performant one listed first, are:

1. General purpose—provide a balance of computing, memory, and networking resource.
 - A1**—AWS Graviton with 64-bit ARM Neoverse cores and custom silicon.
 - T3, T3a**—AWS Nitro System; burstable general-purpose instance type.
 - M6g**—ARM-based AWS Graviton2.

⁴ PHP evolved from a set of Perl scripts designed to produce dynamic web pages in a general-purpose server-side scripting language. The code embedded into an HTML source document is interpreted by a web server with a PHP processor module, which generates the resulting web page.

2. Compute optimized—for computing bound applications; high-performance processors.
C6g—ARM-based AWS Graviton2.
C5n—3.0-GHz Intel Xeon Platinum with AVX-512 instruction set.
3. Memory optimized—deliver fast performance for large memory footprint workloads.
R6g—Arm-based AWS Graviton2.
X1e—Intel Xeon E7-8880 v3; up to 3 904 GiB of DRAM memory.
R5a—AMD EPYC 7000; all core turbo clock speed of 2.5-GHz AWS Nitro System.
4. Accelerated computing—instances use hardware accelerators, or co-processors, to perform functions, such as floating-point number calculations, graphics processing, or data pattern matching, efficiently.
P3—High-frequency Intel Xeon E5-2686 v4 (Broadwell) or 2.5 GHz (base) Intel Xeon P-8175M and up to 8 NVIDIA Tesla V100 GPUs, each pairing 5 120 CUDA Cores and 640 Tensor Cores.
P2—Intel Xeon E5-2686 v4 (Broadwell); NVIDIA K80 GPUs, each with 2 496 parallel processing cores and 12 GiB of GPU memory.
G4—Intel Xeon Scalable (Cascade Lake); NVIDIA T4 Tensor Core GPUs.
G3—Intel Xeon E5-2686 v4 (Broadwell); NVIDIA Tesla M60 GPUs, each with 2 048 parallel processing cores and 8 GiB of video memory.
5. Storage optimized—for workloads that require high, sequential read-and-write access to very large data sets on local storage.
I3—Intel Xeon E5-2686 v4 (Broadwell); non-volatile memory express (NVMe) SSD-backed instance storage.
I3en—3.1 GHz Intel Xeon Scalable (Skylake) processors with AVX-512 instruction set; up to 60 TB of NVMe SSD.
H1—2.3 GHz Intel Xeon E5 2686 v4; up to 16 TB of HDD storage.

Each instance packages a different combination of processors, memory, storage, and network bandwidth. The number of vCPUs, as well as the type of processor, its architecture, and clock speed, are different for different instance types. A vCPU is a virtual processor assigned to one virtual machine. Memory is given in Gibibytes, $1 \text{ GiB} = 2^{30}$ bytes or 1 073 741 824 bytes, while $1 \text{ GB} = 10^9$ bytes. There are no recent benchmarks comparing the performance of supercomputers in the top 500 list with AWS instances.

2.3 Google Clouds

Google is a major player in cloud computing; it has pioneered applications of Artificial Intelligence (AI) and Machine Learning (ML), has developed TPU (Tensor Processing Units), and populated some of its instances with TPUs. Moreover, unlike Amazon, which is a very secretive organization, Google has disseminated information about the new developments in hardware and software of its cloud infrastructure. The large number of papers published by Google Research Division contributes significantly to advancements in cloud computing.

Google Cloud infrastructure consists of a large number of clusters in multiple geographical locations. A typical cluster has around 10 000 servers, and its workload is a mix of CPU-intensive batch computations and in-memory databases for latency-sensitive applications. Google's effort is

concentrated in several areas of Infrastructure-as-a-Service (IaaS), Software-as-a-Service (SaaS), and Platform-as-a-Service (PaaS) [204]. Services such as Gmail, Google Drive, Google Calendar, Picasa, and Google Groups are free of charge for individual users and available for a fee for organizations. These services are running on Google clouds and can be invoked from a broad spectrum of devices, including smartphones, tablets, and laptops. The data for these services is stored in data centers on the cloud.

Google App Engine. App Engine (AE) is an infrastructure for building web and mobile applications and running these application on Google servers. Initially, it supported only Python, but support for Java was added later. The database for code development can be accessed with GQL (Google Query Language) with a SQL-like syntax.

App Engine is an ensemble of computer, storage, search, and networking services. The *Compute Engine* (CE) supports the creation of VMs with resources tailored to the application needs. The CE configurations range from micro instances to the ones with 32 vCPUs or 208 GB of memory. Up to 64 TB of network storage can be attached to a VM. Always-encrypted local solid-state drive (SSD) block storage and automatic scaling are also supported. Several operating systems including Debian, CentOS, CoreOS, SUSE, Ubuntu, Red Hat, FreeBSD, or Windows 2008 R2 and 2012 R2 are supported. One can create and manage CE instances in several ways, including: (i) Google Cloud Console, a web UI comparable to management tools like vCenter or XenCenter; (ii) Google Cloud SDK; and (iii) Compute Engine API.

Container Engine (CntE) is a cluster manager and orchestration system for Docker containers built on the Kubernetes system. The *Container Registry* stores private Docker images. CntE schedules and manages containers automatically according to user specifications. JSON config files are used to specify the amount of CPU/memory, the number of replicas, and other relevant information. The Cloud Container Engine SLA commitment is a monthly uptime of at least 99.5%. *Cloud Functions* (CF) is a lightweight, event-based, asynchronous system to create single-purpose functions that respond to cloud events. CFs are written in Javascript and executed in a Node.js runtime environment. *Cloud Load Balancing* supports scalable load balancing on the Google Cloud Platform.

Cloud Storage is a unified object storage enabling a multi-region operation for applications including video streaming and frequently accessed web and images sites. *Cloud SQL* is a fully-managed database service, *Cloud Bigtable* is a high performance NoSQL database service for large analytical and operational workloads, and *Cloud Datastore* is a highly-scalable NoSQL database for web and mobile applications.

Big Data applications are supported by several services. *Biggquery* is a fully-managed enterprise data warehouse for large-scale data analytics. *Cloud Dataflow* supports stream and batch execution of pipelines. *Cloud Dataproc* manages Spark and Hadoop service. *Cloud Datalab* is an interactive tool for large-scale data exploration, analysis, and visualization. *Cloud Pub/Sub* is a global service for real-time reliable messaging and data streaming.

Network functionality is managed by *Cloud Virtual Network* (CVN). AppEngine users can connect resources to each other and isolate them from one another in a Virtual Private Cloud using the CVN. Cloud routers maintain virtual routers enabling Border Gateway Protocol (BGP) to route updates between a user Compute Engine network and user non-Google network. CVN supports Virtual Private Networks (VPNs) and distributed firewalls. *Cloud CDN* is a low-latency, low-cost content delivery network. *Cloud DNS* is a resilient, low-latency DNS for Google's worldwide network.

Several App Engine services support security. *Cloud Identity and Access Management* (IAM) provides tools to manage resource permissions and to map job functions within a company to groups and roles. *Cloud KMS* is a key management service enabling users to manage encryption and generate, use, rotate, and destroy AES256 key encryption keys. *Cloud Security Scanner* is used to scan for common vulnerabilities in Google App Engine applications. App Engine provides cloud development tools. *Cloud SDK* is a set of tools including *gcloud*, *gsutil*, and *bq*, to access the Compute Engine, the Cloud Storage and other services. *Cloud Source Repositories* provides Git version control to support collaborative development for applications or services. *Android Studio*, *PowerShell*, *IntelliJ*, *Eclipse*, and *Visual Studio* tools are also available.

The array of management tools includes *Stackdriver* which supports monitoring, logging, and diagnostics tools along with *Stackdriver Monitoring* tool, which provides performance, uptime, and overall health information on cloud applications. *Stackdriver Debugger* lets users inspect the state of an application without logging statements and without stopping or slowing down the application. *Stackdriver Error Reporting* counts, analyzes, aggregates crashes, and provides a cleaned exception stack trace.

Virtually all CSPs require users to specify in one form or another the type and amount of resources used; jobs that exceed these limits are either throttled or killed. Google uses a system called Autopilot, evolved from the previous job management sub-system of Borg, discussed in Section 4.7, to configure resources automatically. Autopilot adjusts the number of concurrent tasks in a job and the CPU/memory limits for individual tasks and uses vertical and horizontal scaling to reduce the difference between user-specified CPU and RAM resource limits and the actual usage, the so-called *slack*, while making sure that tasks do not run out of resources [426]. The slack of Autopiloted jobs is half of the slack of manually managed jobs, and the number of jobs severely impacted by out-of-memory is also significantly reduced.

A range of services supporting machine learning are also provided. *Cloud Machine Learning* is a managed service based on the TensorFlow model to build machine learning models, which work on any type of data. *Cloud Natural Language API* is a text analysis tool that can be used to extract information about people, places, events, etc., while *Cloud Speech API* enables developers to convert audio to text by applying powerful neural network models and *Cloud Vision API* is used to understand the content of an image by encapsulating powerful machine learning models.

Other widely used Google services

Gmail. The service hosts emails on Google servers and provides a web interface to access them and tools for migrating from Lotus Notes and Microsoft Exchange.

Google Docs is a web-based software for building text documents, spreadsheets, and presentations. It supports features such as tables, bullet points, basic fonts, and text size; it allows multiple users to edit and update the same document, to view the history of document changes, and it provides a spell checker. The service allows users to import and export files in several formats including Microsoft Office, PDF, text, and OpenOffice extensions.

Google Calendar is a browser-based scheduler; it supports multiple calendars for a user, the ability to share a calendar with other users, the display of daily/weekly/monthly views, search for events, and synchronization with the Outlook Calendar. The calendar is accessible from mobile devices; event reminders can be received via SMS, desktop pop-ups, or emails. It is also possible to share your calendar with other Google calendar users. *Picasa* is a tool to upload, share, and edit images; it provides 1 GB

of disk space per user free of charge. Users can add tags to images and attach locations to photos using *Google Maps*. *Google Groups* enables users to host discussion forums to create messages online or via email.

Google Co-op enables users to create customized search engines based on a set of categories.

Google Base enables users to load structured data from different sources to a central repository, a very large, self-describing, semi-structured, heterogeneous database where each item follows a simple schema: item type, attribute names. Since few users are aware of this service, *Google Base* is accessed in response to keyword queries posed on *Google.com*, provided that there is relevant data in the database. To fully integrate Google Base, the results should be ranked across properties. Also, the service needs to propose appropriate refinements with candidate values in select menus; this is done by computing histograms on attributes and their values during query time.

Google Drive is an online service for data storage available for PCs, MacBooks, iPhones, iPads, and Android devices and allows organizations to purchase up to 16 TB of storage.

Specialized structure-aware search engines for several areas, including travel, weather, and local services, have already been implemented. However, the data available on the web covers a wealth of human knowledge; it is not feasible to define all the possible domains, and it is nearly impossible to discern where one domain ends and another begins.

Google adheres to a bottom-up, engineer-driven, liberal licensing, and user-application development philosophy, while Apple, a recent entry in cloud computing, tightly controls the technology stack, builds its own hardware, and requires the applications developed to follow strict rules. Apple products, including the iPhone, the iOS, the iTunes Store, Mac OS X, and iCloud, offer unparalleled polished and effortless interoperability, while the flexibility of Google results in more cumbersome user interfaces for the broad spectrum of devices running the Android OS.

Google manages vast amounts of data. In a world where users would most likely desire to use multiple cloud services from independent providers, the question of whether the traditional Data Base Management Services (DBMS) are sufficient to ensure interoperability comes to mind. A DBMS efficiently supports data manipulations and query processing but operates in a single administrative domain and uses a well-defined schema. The interoperability of data-management services requires *semantic integration* of services based on various schemas. An answer to the limitations of traditional DBMS are *dataspaces* that do not aim at data integration but at data co-existence [183].

2.4 Microsoft Windows Azure and online services

Azure and Online Services are PaaS and, respectively, SaaS cloud platforms from Microsoft. Windows Azure is an operating system, SQL Azure is a cloud-based version of the SQL Server, and Azure AppFabric is a collection of services for cloud applications.

Windows Azure has three core components (see Fig. 2.4): *Compute*, which provides a computation environment, *Storage* for scalable storage, and *Fabric Controller*, which deploys, manages, and monitors applications; it interconnects nodes consisting of servers, high-speed connections, and switches. The Content Delivery Network (CDN) maintains cache copies of data to speed up computations. The Connect subsystem supports IP connections between the users and their applications running on Windows Azure. The API interface to Windows Azure is built on REST, HTTP, and XML. The platform

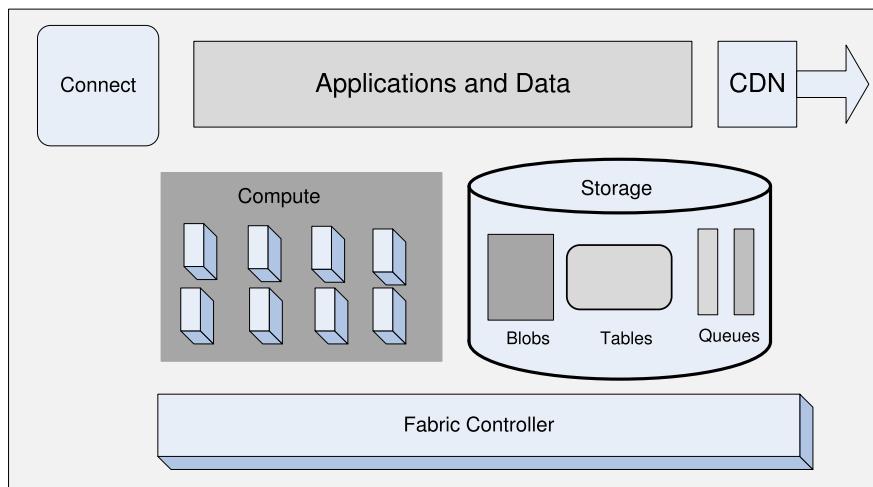


FIGURE 2.4

Azure components: Compute—runs cloud applications; Storage—uses blobs, tables, and queues to store data; Fabric Controller—deploys, manages, and monitors applications; CDN—maintains cache copies of data; and Connect—allows IP connections between the user systems and applications running on Windows Azure.

includes five services: Live Services, SQL Azure, AppFabric, SharePoint, and Dynamics CR. A client library and tools are also provided for developing cloud applications in Visual Studio.

Computations carried out by an application are implemented as one or more *roles*; an application typically runs multiple *instances of a role*. One distinguishes: (i) web role instances—to create web applications; (ii) worker role instances—run Windows-based codes; and (iii) VM role instances—run a user-provided Windows Server 2008 R2 image.

Scaling, load balancing, memory management, and reliability are ensured by a *fabric controller*, a distributed application replicated across a group of machines that owns all of the resources in its environment (computers, switches, and load balancers), and it is aware of every Windows Azure application. The fabric controller decides where new applications should run; it chooses the physical servers to optimize utilization using configuration information uploaded with each Windows Azure application. The configuration information is an XML-based description of how many web role instances, how many worker role instances, and what other resources the application needs; the fabric controller uses this configuration file to determine how many VMs to create.

Blobs, tables, queue, and drives are used as scalable storage. A blob contains binary data, and a container consists of one or more blobs. Blobs can be up to a terabyte, and they may have associated metadata, e.g., the information about where a JPEG photograph was taken. Blobs allow a Windows Azure role instance to interact with persistent storage as if it were a local NTFS⁵ file system. Queues

⁵ NTFS (New Technology File System) is the standard file system of the Microsoft Windows operating system starting with Windows NT 3.1, Windows 2000, and Windows XP.

enable web role instances to communicate asynchronously with worker role instances. Microsoft Azure platform currently does not provide or support any distributed parallel-computing frameworks, such as *MapReduce*, *Dryad* or *MPI*, other than the support for implementing basic queue-based job scheduling [213].

2.5 IBM clouds

In recent years IBM has accelerated its effort to eventually overtake Amazon, Google, Microsoft, and other CSP competitors. This effort was boosted by the 2018 acquisition of open-source software provider Red Hat and the 2020 decision to make a \$1 billion investment in cloud computing. IBM is also a leader in quantum computing and quantum information theory and provides access to several quantum computing services as discussed in Section 13.3.

IBM is making a big play in hybrid clouds and gives companies tools to more easily navigate between public and private environments. Some of the reason for focusing on hybrid clouds are: (i) flexibility—users are able to run applications wherever they perform best; (ii) increased security—based on container built-in security; and (iii) performance—use resources of the public cloud for larger workloads after developing, testing, and running smaller workloads on a private cloud.

IBM is also operating in edge computing. The motivation for this effort is that 5G technology brings computation and data storage closer to where data is generated, enabling better data control, reduced costs, faster insights and actions, and continuous operations. It is estimated that, 75% of enterprise data will be processed at the edge in 2025, instead of only 10% today.

IBM offers a cloud platform with over 170 products and services. High-performance cloud servers can be configured in hourly and monthly options, including up to 20 TB of bandwidth. Red Hat helps to build and deploy container-based applications on a fully managed, integrated, and reliable platform. The emphasis on cloud security is winning over large organizations, and IBM AI services aim to modernize, collect, organize, analyze, and infuse user data in new ways to accelerate user's journey to AI.

2.6 Cloud storage diversity and vendor lock-in

The short history of cloud computing shows that cloud services may be unavailable for short, or even for extended, periods of time. Such an interruption of service is likely to impact negatively an organization using the cloud and possibly diminish, or cancel completely, the benefits of utility computing for that organization. The potential for permanent data loss in case of a catastrophic system failure poses an even greater danger.

Last, but not least, a vendor may decide to increase the prices for service and charge more for computing cycles, memory, storage space, and network bandwidth than other cloud service providers. The alternative, switching to another cloud service provider, could be very costly due to the large volume of data to be transferred from the old to the new provider. Transferring petabytes of data over the network takes a substantial amount time and incurs substantial charges for the network bandwidth.

A solution to problems posed by vendor lock-up is to replicate data to multiple cloud service providers. Straightforward replication is very costly and, at the same time, poses technical challenges.

The overhead to maintain data consistency could drastically affect the performance of the virtual storage system consisting of multiple full replicas of organization's data spread over multiple vendors.

Another solution could be based on an extension of the design principle of a RAID-5 system used for reliable data storage. This elegant idea immediately raises several questions: How does the response time of such a scheme compare with the one of a single storage system? How much overhead is introduced by a proxy? How could this scheme avoid a single point of failure, namely, the proxy? Are there standards for data access implemented by all vendors? An experiment to answer some of these questions has been reported [8]; the RACS system uses the same data model and mimics the interface to AWS S3. The prototype implementation in [8] led the authors to conclude that the cost increases and the performance penalties of the RACS systems are relatively minor. An implementation avoiding the single point of failure uses several proxies; clients connected to several proxies can access data stored on multiple clouds.

In practice, it remains to be seen if such a solution is feasible for organizations with a very large volume of data, given the limited number of cloud storage providers and the lack of standards for data storage. A basic question is if it makes sense to trade basic tenets of cloud computing, such as simplicity and homogeneous resources controlled by a single administrative authority, for increased reliability and for freedom from vendor lock-in.

This brief discussion hints at the need for standardization and for scalable solutions, two of the many challenges faced by cloud computing in the near future. The pervasive nature of scalability dominates all aspects of cloud management and cloud applications. Solutions performing well on small systems no longer do when the system scale increases by one or more orders of magnitude. Experiments with small test-bed systems produce inconclusive results. The only alternative is to conduct intensive simulations to prove, or disprove, the advantages of a particular algorithm for resource management or the feasibility of a particular data-intensive application.

2.7 Cloud interoperability

Cloud interoperability could alleviate the concerns that users become hopelessly dependent on a single cloud service provider, the so called vendor lock-in, discussed in Section 2.6. It seems natural to ask the question if a “cloud of clouds,” a federation of clouds that cooperate to provide a better user experience, is technically and economically feasible. Since the Internet is a network of networks, a federation of clouds seems plausible [58–60].

Closer scrutiny shows that the extension of the concept of interoperability from networks to clouds is far from trivial. A network offers one high-level service, the transport of digital information from a source, a host outside a network to a destination, another host, or another network that can deliver the information to its final destination. This transport of information through a network of networks is feasible because agreements on basic questions were reached before the Internet was born. It was then decided on how to: (a) uniquely identify the source and the destination of the information; (b) navigate through a maze of networks; and (c) transport the data between a source and a destination. The three elements on which agreements were reached are, respectively, the IP address, the IP protocol, and transport protocols such as TCP and UDP.

The situation is quite different in cloud computing. First, there are no standards for either storage or processing; second, the clouds we have seen so far are based on different delivery models. Moreover, the

set of services supported by each of these delivery models is not only large, but it is open; new services are offered every few months. The question if Cloud Service Providers are willing to cooperate to build up a federation of clouds is an open one. Some CSPs may think that they have a competitive advantage due to the uniqueness of the added value of their services. Thus, exposing how they store and process information may adversely affect their business. Moreover, no CSP will be willing to change its internal operation, so a first question is if an Intercloud could be built under these conditions.

Following the concepts borrowed from the Internet, a federation of clouds that does not dictate the internal organization or the structure of a cloud, but only the means to achieve cloud interoperability, is feasible. Nevertheless, building such an infrastructure seems to be a formidable task. First, we need a set of standards for interoperability covering items such as: naming, addressing, identity, trust, presence, messaging, multicast, and time. We need common standards for identifying all objects involved, the means to transfer, store, and process information, and we also need a common clock to measure the time between two events.

An Intercloud would then require the development of an *ontology*⁶ for cloud computing. Then each cloud service provider would have to create a description of all resources and services using this ontology. Due to the very large number of systems and services, the volume of information provided by individual cloud service providers would be so large that a distributed database, not unlike the Domain Name Service (DNS), would have to be created and maintained. According to [58], this vast amount of information would be stored in Intercloud *root nodes*, analogous to the root nodes of the DNS.

Each cloud would then require an interface, a so-called Intercloud *exchange*, to translate the common language describing all objects and actions included in a request originating from another cloud in terms of its internal objects and actions. To be more precise, a request originating in one cloud would have to be translated from the internal representation in that cloud to a common representation based on the shared ontology, and then, at the destination, it should be translated into an internal representation that can be acted upon by the destination cloud. This raises immediately the question of efficiency and performance. This question cannot be fully answered now, as an Intercloud exists only on paper, but there is little doubt that the performance will be greatly affected.

Security is a major concern for cloud users, and an Intercloud could create new threats. The primary concern is that tasks will cross from one administrative domain to another and sensitive information about tasks and user could be disclosed during this migration. A seamless migration of tasks in an Intercloud requires a well-thought-out trust model.

The Public-Key Infrastructure (PKI), an all-or-nothing trust model, is not adequate for an Intercloud where the trust must be nuanced. PKI is a model to create, distribute, revoke, use, and store digital certificates. It involves several components: (1) The Certificate Authority (CA) binds public keys to user identities in a given domain; (2) the third-party Validation Authority (VA) guarantees the uniqueness of the user identity; and (3) the Registration Authority (RA) guarantees that the binding of the public key to an individual cannot be challenged, the so-called *non-repudiation*. A nuanced model for handling digital certificates means that one cloud acting on behalf of a user may grant access to another cloud to read data in storage, but not to start new instances.

A solution for trust management is based on dynamic *trust indexes* that can change over time [59]. The Intercloud roots play the role of the CA, while the Intercloud exchanges determine the trust indexes

⁶ An ontology provides the means for knowledge representation within a domain. It consists of a set of domain concepts and the relationships among the concepts.

between clouds. Encryption must be used to protect the data in storage and in transit in the Intercloud. OASIS⁷ Key Management Interoperability Protocol (KMIP) is proposed for key management. In summary, an Intercloud opens up a wide range of interesting research topics. The practicality of the concepts can only be discussed if cloud standardization efforts under way at NIST bear fruit.

2.8 Service-level Agreements and Compliance-level Agreements

A Service Level Agreement (SLA) is a negotiated contract between two parties, the customer and the service provider. The agreement can be legally binding or informal and specifies customer services rather than how the service provider delivers services. The objectives of the agreement are: (i) identify and define the customer's needs and constraints, including the level of resources, security, timing, and quality of service; (ii) provide a framework for understanding including a clear definition of classes of service and costs; (iii) simplify complex issues, for example, clarify the boundaries between the responsibilities of clients and those of service provider in case of failures; (iv) reduce areas of conflict and eliminate unrealistic expectations; and (v) encourage dialog in the event of disputes.

An SLA records a common understanding in several areas: (i) services, (ii) priorities, (iii) responsibilities, (iv) guarantees, and (v) warranties. An agreement usually covers: services to be delivered, performance, tracking and reporting, problem management, legal compliance and resolution of disputes, customer duties and responsibilities, security, handling of confidential information, and termination.

Each area of service in cloud computing should define a “target level of service” or a “minimum level of service” and specify the levels of availability, serviceability, performance, operation, or other attributes of the service, such as billing; penalties may also be specified in the case of non-compliance with the SLA. It is expected that any Service-Oriented Architecture (SOA) will eventually include middleware supporting SLA management; the Framework 7 project supported by the European Union is researching this area, see <http://sla-at-soi.eu/>.

There are two well-differentiated phases in SLA management: the negotiation of the contract and the monitoring of its fulfillment in real-time. In turn, automated negotiation has three main components: (i) the *object of negotiation* that defines the attributes and constraints under negotiation; (ii) the *negotiation protocols* that describe the interaction between negotiating parties; and (iii) the *decision models* that are responsible for processing proposals and generating counter proposals.

It is critical for a cloud user to carefully read the service-level agreement and to understand the limitations of the liability a cloud provider is willing to accept. In many instances, the liabilities do not apply to damages caused by a third party or to failures attributed either to customer's hardware and software or to hardware and software from a third party. For example, if a distributed denial of service attack (DDoS) causes the entire IaaS infrastructure to fail, the cloud service provider is responsible for the consequences of the attack. The user is responsible if the DDoS affects only several instances including the ones running the user application.

The concept of compliance in cloud computing is discussed in [70] in the context of the user ability to select a provider of service; the selection process is subject to customizable compliance with

⁷ OASIS stands for Organization for the Advancement of Structured Information Standards.

user requirements such as security, deadlines, and costs. The authors propose an infrastructure called *Compliant Cloud Computing* (C3) consisting of: (i) a language to express user requirements and the Compliance Level Agreements (CLA), and (ii) the middleware for managing CLAs. A policy-based framework for automated SLA negotiation for a virtual computing environment is described in [520].

2.9 Responsibility sharing between user and service provider

After reviewing cloud services provided by Amazon, Google, and Microsoft, we are in a better position to understand the differences between SaaS, IaaS, and PaaS. There is no confusion about SaaS, the service provider supplies both the hardware and the application software; the user has direct access to these services through a web interface and has no control on cloud resources. Typical examples are Google with Gmail, Google Docs, Google Calendar, Google Groups, and Picasa and Microsoft Online Services.

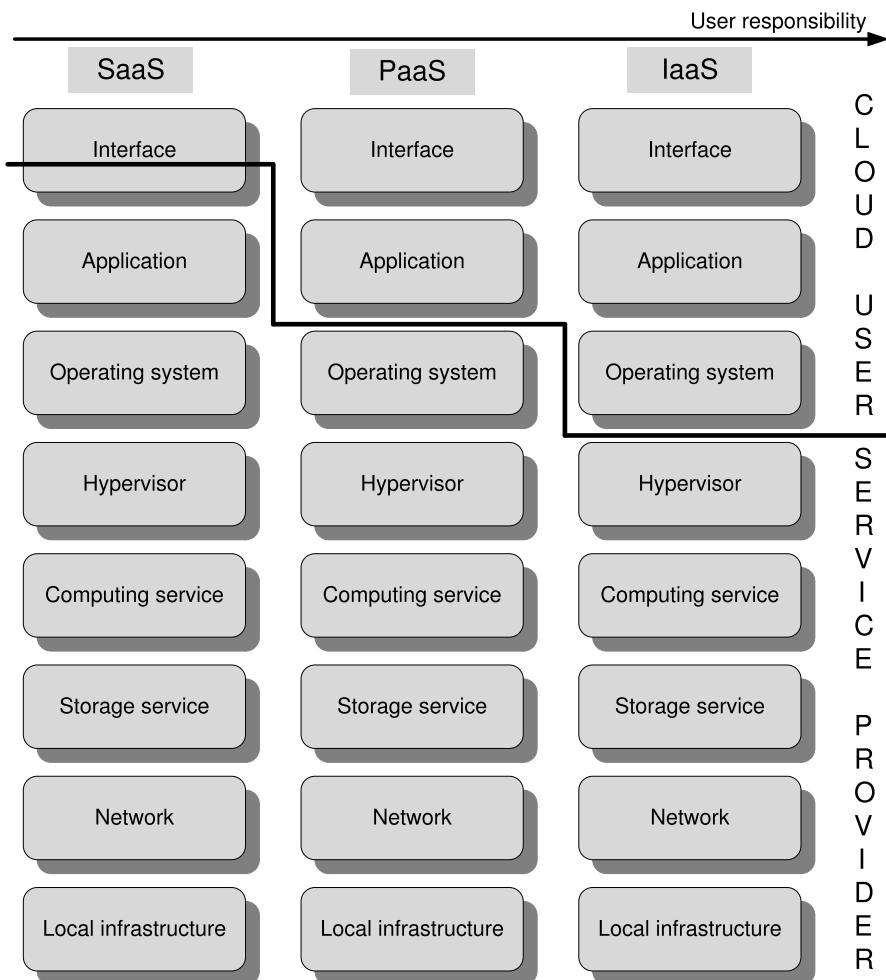
Fig. 2.5 describes the limits of responsibility between a user and CSP. An SaaS user is partially responsible for the interface; the user responsibility increases in the case of PaaS and includes the interface and the application. An IaaS user is responsible for all events in the virtual machine running the application; the service provider supplies the hardware (servers, storage, networks), and system software (operating systems, databases) and ensures system attributes such as security, fault-tolerance, and load balancing. PaaS provides only a platform including the hardware and system software such as operating systems and databases; the service provider is responsible for system updates, patches, and the software maintenance. PaaS does not allow any user control of the operating system, security features, or the ability to install applications. Typical examples are Google App Engine, Microsoft Azure, and Force.com, provided by Salesforce.com.

2.10 User challenges and experience

Computer clouds pose a fair number of challenges to software developers and cloud users, including security, compliance, managing costs, and lack of expertise. Cloud users attempt to reduce their cloud costs through a range of measures including: (i) careful resource utilization monitoring; (ii) avoiding peak hours and shutting down temporary workloads; (iii) purchasing AWS reserved instances along with effective use of spot instances; and (iv) moving workloads to regions with lower costs.

There are a few studies of user experiences based on a large population of cloud computing users. An empirical study of the experience of a small group of users of the Finish Cloud Computing Consortium is reported in [383]. The main user concerns are: security threats; the dependence on fast Internet connections; forced version updates; data ownership; and user behavior monitoring. While all users reported that trust in the cloud services is important, two-thirds raised the issue point of fuzzy boundaries of liability between cloud user and the provider, about half did not fully comprehend the cloud's functions and its behavior, and about one-third were concerned about security threats.

The security threats perceived by this group of users are: (i) abuse and villainous use of the cloud; (ii) APIs that are not fully secure; (iii) malicious insiders; (iv) account hijacking; (v) data leaks; and (iv) issues related to shared resources. Identity theft and privacy were a major concern for about half of

**FIGURE 2.5**

The limits of responsibility between a cloud user and the cloud service provider.

the users questioned; availability, liability, and data ownership and copyright was raised by a third of respondents.

The suggested solutions to these problems are: Service Level Agreements and tools to monitor usage should be deployed to prevent the abuse of the cloud; data encryption and security testing should enhance the API security; an independent security layer should be added to prevent threats caused by malicious insiders; strong authentication and authorization should be enforced to prevent account hijacking; data decryption in a secure environment should be implemented to prevent data leakage; and

compartmentalization of components and firewalls should be deployed to limit the negative effect of resource sharing.

A broad set of concerns identified by the NIST working group on cloud security includes: (i) potential loss of control/ownership of data; (ii) data integration, privacy enforcement, data encryption; (iii) data remanence after de-provisioning; (iv) multi-tenant data isolation; (v) data location requirements within national borders; (vi) hypervisor security; (vii) auditing of data-integrity protection; (viii) verification of subscriber policies through provider controls; and (ix) certification/accreditation requirements for a given cloud service.

A study conducted by IBM [249] identified barriers to public and private cloud adoption. The top workloads mentioned by users involved in this study are: data mining and other analytics (83%), application streaming (83%), help-desk services (80%), industry specific applications (80%), and development environments (80%). The study also identified workloads that are not good candidates for migration to a public cloud environment, such as: (a) sensitive data such as health care records; (b) multiple co-dependent services, e.g., online transaction processing; (c) third-party software without cloud licensing; (d) workloads requiring auditability and accountability; and (e) workloads requiring customization.

2.11 Software licensing

Software licensing for cloud computing is an enduring problem without a universally accepted solution. License-management technology is based on the old model of computing centers with licenses granted on the basis of named users or as site licenses. This technology developed for a centrally managed environment cannot accommodate the distributed service infrastructure of cloud computing. IBM reached an agreement allowing some of its software products to be used on EC2. Also, MathWorks developed a business model for MATLAB® use in Grid environments [84]. SaaS deployment model is gaining acceptance because it allows users to pay for only the services they use.

There is significant pressure to change traditional software licensing and find non-hardware-based solutions for cloud computing; the increased negotiating power of the users coupled with the increased software piracy has renewed interest in alternative schemes such as those proposed by the SmartLM research project (<http://www.smartlm.eu>). SmartLM license management requires a complex software infrastructure involving Service Level Agreement, negotiation protocols, authentication, and other management functions.

A commercial product based on the ideas developed by this research project is *elasticLM*, which provides license and billing Web-based services [84]. The architecture of the elasticLM license service has several layers: co-allocation, authentication, administration, management, business, and persistency. The authentication layer authenticates communications between the license service and the billing service, as well as the individual applications; the persistence layer stores the usage records; the main responsibility of the business layer is to provide the licensing service with the licensing prices; the management coordinates various components of the automated billing service.

When a user requests a license from the license service, the terms of the license usage are negotiated, and they are part of a Service Level Agreement (SLA) document; the negotiation is based on application-specific templates and the license cost becomes part of the SLA. SLA describes all aspects of resource usage, including application ID, duration, number of processors, and guarantees, such as

the maximum cost and deadlines. When multiple negotiation steps are necessary, the WS-Agreement Negotiation protocol is used.

To verify the authorization to use a license, an application must have the certificate of an authority. This certificate must be available locally to the application because the application may be executed in an environment with restricted network access; this opens the possibility for an administrator to hijack the license mechanism by exchanging the local certificate.

2.12 Challenges faced by cloud computing

Cloud computing inherits some of the challenges of parallel and distributed computing discussed in Chapter 3 and, at the same time, faces major challenges of its own. The complexity of hardware and software infrastructure supporting cloud service is astounding. Two billion lines of code are maintained by Google and drive applications such as Google Search, Google Maps, Google Docs, Google+, Google Calendar, Gmail, YouTube, and every other Google Internet service. By comparison, Windows operating system developed by Microsoft since the 1980s has some 50 million lines of code, a factor of 40 less than what Google has developed during two decades of existence.

The software stack for cloud computing has dramatically evolved in its quest to provide a unified higher-level view of the system, rather than a large collection of individual machines. Virtualization and containerization are ubiquitous abstractions that allow easier access to the increasingly larger and diverse population of cloud users. Distributed and semi-structured storage systems such as Google’s BigTable [94] or Amazon’s Dynamo [133] are widely used. The list of systems supporting higher-level abstractions includes Pig [190], FiumeJava [90], and Spark [532]. The effort to build one layer of abstraction removed from the underlaying hardware, a sort of operating system for Internet-scale jobs is underway. Systems such as Dryad [255], DryadLinq [530], Mesos [240], Borg [495], Omega [442], and Kubernets [80] attempt to bridge the gap between a clustered infrastructure and the assumptions made by applications about their environments. These systems manage a virtual computer aggregating the resources of a physical cluster with a very large number of independent servers.

In the early days of network-centric computing, it was postulated that Web searching is the “killer application” that will drive the software and hardware of large-scale systems for the next decades [52]. It turns out that the applications running on computer clouds are very diverse. *The broad spectrum of cloud applications adds to the challenges faced by cloud infrastructure. For example, controlling the tail latency of workloads consisting of a mix of time-critical and batch jobs is far from trivial [130]. Some critical system requirements are contradictory, e.g., multiplexing resources to increase efficiency and lowering the response time, while supporting performance and security isolation.*

The specific challenges differ for the cloud delivery models, but in all cases, the difficulties are created by the very nature of utility computing based on resource sharing and resource virtualization. Moreover, cloud computing requires a different trust model than the ubiquitous user-centric model we have been accustomed to for a very long time. *The most significant challenge is security; gaining the trust of a large user base is critical for the future of cloud computing. It is unrealistic to expect that a public cloud will provide a suitable environment for all applications.*

The SaaS model faces similar challenges because other online services are required to protect private information, such as financial or healthcare services. Since in this case, a user interacts with cloud services through a well-defined interface, in principle, it is less challenging for the provider of service

to close some of the attack channels. Still, such services are vulnerable to denial-of-service attacks, and the users are fearful of malicious insiders.

Data in storage is most vulnerable to attacks, so special attention should be devoted to the protection of storage servers. Data replication necessary to ensure continuity of service in the case of storage-system failure increases vulnerability. Data encryption may protect data in storage, but eventually, data must be decrypted for processing, and then it is exposed to attacks. Homomorphic encryption enables logic operations to be performed on encrypted data, but it is impractical at this time.

The IaaS is by far the most challenging model to defend against attacks; indeed, an IaaS user has considerably more degrees of freedom than allowed by the other two cloud delivery models. An additional source of concern is that the considerable resources of a cloud could serve as the host to initiate attacks against the networking and the computing infrastructure.

Virtualization is a critical design option for this model, but it exposes the system to additional sources of attacks. The trusted computing base (TCB) of a virtual environment includes not only the hardware and the hypervisor, but also the management operating system. As we shall see in Section 8.10, the entire state of a VM can be saved to a file to allow migration and recovery, both highly desirable operations; yet, this possibility challenges the strategies to bring the servers belonging to an organization to a desirable and stable state. Indeed, an infected VM can be inactive when the systems are cleaned up and an infected VM can wake up later and infect other systems. This is another example of the deep intertwining of desirable and undesirable effects of basic cloud computing technologies.

The next major challenge is related to cloud resource management. A systematic, rather than ad hoc, resource-management strategy requires the existence of controllers tasked to implement several classes of policies: admission control, capacity allocation, load balancing, energy optimization, and, last but not least, to provide Quality of Service (QoS) guarantees.

It seems reasonable to expect that such a complex system can only function based on self-management principles. But self-management and self-organization raise the bar for the implementation of logging and auditing procedures critical for the security and trust in a provider of cloud computing services. Under self-management, it becomes next to impossible to identify the reasons why a certain action that resulted in a security breach was taken.

The last major challenge we address is related to interoperability and standardization. Vendor lock-in, the fact that a user is tied to a particular cloud service provider is a major concern for cloud users. Standardization would support interoperability and, thus, alleviate some of the fears that a service critical for a large organization may not be available for an extended period of time. But imposing standards at a time when a technology is still evolving is not only challenging but can be counterproductive because it may stiffen innovation.

Further research in various areas of cloud computing is an important objective of the National Science Foundation (NSF). The NSF supports research-community access to two cloud facilities, CloudLab and Chameleon. *CloudLab* is a testbed allowing researchers to experiment with cloud architectures and new applications. Some 15 000 cores, at three sites in Utah, Wisconsin, and South Carolina, are available for such experiments. *Chameleon* is an OpenStack KVM experimental environment for large-scale cloud research.

2.13 Cloud computing as a disruptive technology

We present now a few statistics supporting the view that computer clouds are changing the enterprise computing landscape at a stunning speed. Gartner⁸ predicts that the public cloud service revenues will increase from \$196.7 billion in 2018 to \$354.6 billion in 2022. More than \$1.3 trillion in IT spending will be affected by the shift to the cloud. PaaS, SaaS, and IaaS cloud service providers will almost double their revenues, from \$26.4 to \$58.0 billion, from \$85.7 to \$151.1 billion, and from \$32.4 to \$74.1 billion, respectively.

Eighty-two percent of the enterprise workload resided on clouds and 40 zettabytes (40×10^{21} bytes) of data flowed through cloud servers and networks by the end of 2020. Ninety-four percent of enterprises use the cloud, and 66% of enterprises have dedicated cloud teams working on cloud cost optimization, selecting applications to be run on clouds and setting up policies for cloud users. Manufacturing, professional services, and banking plan to spend the most on cloud computing services, \$19.7 billion, \$18.1 billion, and \$16.7 billion, respectively.

Cost-cutting is the top reason why companies migrate to the cloud. Eighty percent of companies report operation improvements within the first few months of adopting cloud computing. Organizations with more than 1 000 employees are mainly looking for flexibility and reduced costs of operation, while smaller companies want to ensure their business continuity.

A 2021 study reports that: 78% of organizations participating in the study use hybrid clouds, i.e., public and private; 19% use only public and 2% use one private clouds; and 31% of respondents spend more than \$12 million annually on cloud computing. Companies with less than 1 000 employees (SMBs) are the driving force of the economy; 41% of SMBs favor the public cloud, and 94% of SMBs adopt the cloud due to upgraded security. Fifty percent of US government organizations are now using the cloud.

AWS has the largest cloud computing market share of 32%. Dropbox, Google Drive, and Microsoft OneDrive are the leading cloud storage providers with 47.3%, 26.9%, and 15.3% of the market, respectively. The multicloud is an aspirational goal for enterprises concerned about vendor lock-in. Many organizations are working hard to abstract their applications and allow them to be moved across clouds. Legacy vendors have created platforms that can plug into multiple clouds using VMware or Red Hat.

It is expected that in the immediate future artificial intelligence, analytics, the IoT, edge computing, and serverless and managed services will differentiate the top cloud service providers. All these areas require a highly skilled workforce. Ergo, cloud computing experts are in high demand; site reliability engineers, enterprise account executives, customer success managers, and solution architects are sought after, according to LinkedIn.

While enterprises migrate in large numbers to cloud computing to cut costs and reduce risk, the latter rationale is increasingly being challenged. On one hand, the number of potentially vulnerable and attractive targets increases, and, on the other hand, hackers are becoming more adept at exploiting vulnerabilities in cloud systems, see Section 8.15. Some of these vulnerabilities are exposed by the increasingly larger number of cloud users, some with little or no training and understanding of potential risks associated with cloud computing.

⁸ The sources for data presented in this section: (i) <https://www.gartner.com/en/newsroom/press-releases/2019-11-13-gartner-forecasts-worldwide-public-cloud-revenue-to-grow-17-percent-in-2020>; (ii) <https://resources.flexera.com/web/pdf/report-cm-state-of-the-cloud-2021.pdf>.

2.14 Exercises and problems

- Problem 1.** The list of desirable properties of a large-scale distributed system includes transparency of access, location, concurrency, replication, failure, migration, performance, and scaling. Analyze how each one of these properties applies to AWS.
- Problem 2.** Compare the three cloud computing delivery models, SaaS, PaaS, and IaaS, from the point of view of application developers and users. Discuss the security and the reliability of each one of them. Analyze the differences between the PaaS and the IaaS.
- Problem 3.** Compare the Oracle Cloud offerings (see <https://cloud.oracle.com>) with the cloud services provided by Amazon, Google, and Microsoft.
- Problem 4.** Read the IBM report [249] and discuss the workload preferences for private and public clouds and the reasons for the preferences.
- Problem 5.** Many organizations operate one or more computer clusters and contemplate the migration to private clouds. What are the arguments for and against such an effort?
- Problem 6.** Evaluate the SLA toolkit at <http://www.service-level-agreement.net/>. Is the interactive guide useful, and what does it miss? Does the SLA template include all clauses that are important in your view, and what is missing? Are the examples helpful?
- Problem 7.** Software licensing is a major problem in cloud computing. Discuss several ideas to prevent an administrator from hijacking the authorization to use a software license.
- Problem 8.** Annotation schemes are widely used by popular services such as Flickr photo-sharing service, which support the annotation of photos. Sketch the organization of a cloud service used for sharing medical x-ray, tomography, CAT-scans, and other medical images and discuss the main challenges for its implementation.
- Problem 9.** An organization debating whether to install a private cloud or to use a public cloud, e.g., the AWS, for its computational and storage needs, asks your advice. What information will you require to base your recommendation on, and how will you use each one of the following items: (a) the description of the algorithms and the type of the applications the organization will run; (b) the system software used by these applications; (c) the resources needed by each application; (d) the size of the user population; (e) the relative experience of the user population; and (f) the costs involved?
- Problem 10.** A university is debating the question in Problem 9. What will be your advice and why? Should software licensing be an important element of the decision?

Parallel processing and distributed computing

3

This chapter overviews concepts in parallel and distributed systems important for understanding the basic challenges of design and use of large-scale collections of autonomous and heterogeneous distributed systems such as computer clouds. Cloud computing is the result of knowledge and wisdom accumulated over some 60 years of computing and is intimately tied to parallel and distributed processing.

Parallel processing has mesmerized computational science and engineering communities since early days of the computing era, resulting in fascination with high-performance computer systems and, ultimately, with supercomputers. It is hard to expose parallelism in many scientific applications, but, the harder the problem, the more satisfying it was to develop parallel algorithms, implement them, wait for the next generation of processors running at a higher clock rate, and enjoy the impressive speedup. The enterprise computing world seemed more skeptical and less involved in the parallel processing movement.

Distributed computing studies distributed systems, i.e., systems whose components are running on networked computers, communicating and coordinating their actions by message passing or shared memory. Message passing has its roots in the operating system of the 1960s and the widespread use of distributed systems can be traced back to the invention of the Ethernet in the 1970s.

More than a half century after the dawn of the computing era, an eternity in the age of the silicon, disruptive multicore technology forced the community to realize the need to understand and exploit parallelism. Rather than wait for faster clock rates, we should better design algorithms and applications able to use all the cores of a modern processor.

Things changed again when cloud computing showed that there are new applications that can effortlessly exploit parallelism and, in the process, generate huge revenues. A new era in parallel and distributed systems began, the era of Big Data hiding nuggets of useful information and requiring massive amounts of computing resources. The new challenge is to obtain the results faster by effectively harnessing the power of millions of multicore processors and systems on a chip.

To process massive amounts of data efficiently, cloud applications distribute data to large numbers of servers. Cloud applications use a number of instances running concurrently. Many cloud applications are based on the *client-server* paradigm. A relatively simple software, a *thin-client*, is often running on the user's mobile device with limited resources, while the computationally intensive tasks are carried out on the cloud.

Transaction processing systems, including web-based services, represent a large class of applications hosted by computing clouds. Such applications run multiple instances of the service and require reliable and in-order delivery of messages.

Early on scientists and engineers understood that parallel processing requires specialized hardware and system software. It was also clear that the interconnection fabric was critical for the performance of parallel processing systems. Building high-performance computing systems proved to be a major

challenge. The list of companies aiming to support parallel processing and ending up as casualties of this effort is long and includes names such as: Ardent, Convex, Encore, Floating Point Systems, Inmos, Kendall Square Research, MasPar, nCube, Sequent, Tandem, Thinking Machines, and possibly others, now forgotten. The difficulties of developing new programming models and the effort to design programming environments for parallel applications added to the challenges faced by each of these companies.

Hardware parallelism is critical for the performance of a single-core, multicore processors, systems on a chip, multiprocessor systems, clusters, and warehouse-scale computers, the backbone of computer clouds. In this chapter, the first sections cover parallel and distributed system hardware stressing the quantitative rather than qualitative aspects of computer architecture. Basic architectural concepts of modern computer systems, optimizations of computer architecture including caching, out-of-order instruction execution, dynamic scheduling, branch predictions, and ARM architecture are discussed in Sections 3.1, 3.2, and 3.3. Sections 3.4, 3.5, 3.6, and 3.7 cover SIMD architectures, GPUs (Graphics Processing Units), TPUs (Tensor Processing Units), SOCs (Systems On a Chip), and edge cloud computing.

Section 3.8 covers data, task, and thread-level parallelism. Extracting parallelism depends on the application; Sections 3.9 and 3.10 analyze the speedup limits given by Amdahl's Law and Amdahl's Law for multicore processors. Section 3.11 reviews the evolution of the most powerful computing systems, from supercomputers to large-scale distributed systems.

Organization principles for distributed systems such as modularity and layering, presented in Sections 3.12 and 3.13, are applied to the design of peer-to-peer and large-scale systems discussed in Sections 3.14 and 3.15, respectively. Section 3.16 presents composability bounds and scalability, and Section 3.17 surveys fallacies in distributed computing. Finally, Section 3.18 discusses blockchain technologies and applications.

3.1 Computer architecture concepts

Architecture is the science of designing buildings, monuments, parks, cities, and even computers. Computer architecture describes the functionality, organization, and implementation of computer systems. A digital computer processes data and consists of three sub-systems with distinct functionality and performance: (i) CPU (Central Processing Unit)—transforms data; (ii) memory—stores data; and (iii) I/O (Input/Output)—supports communication with the outside world.

The three subsystems communicate with one another using computer busses, groups of wires transporting different information; for example, one set of wires transmits instructions and control information, a second transmits data address, and a third one moves data between communicating entities. Each unit is connected to the bus via an interface, which buffers data before and after any exchange, allowing its end points with different bandwidths to operate asynchronously. A computer bus is the most primitive type of communication device we encounter in computer clouds.

Latency hiding. This concise view of computer architecture is a first sign that *computing and communication are inseparable*, a theme reverberating throughout several chapters. The three subsystems of a processor have different *bandwidths*, i.e., numbers of operations per unit of time, and *latency*, i.e., the time between the start and the completion of an operation. CPU registers are the fastest, the main memory is slower by some two orders of magnitude, and I/O devices are orders of magnitude slower than

memory. The impact of this discrepancy on processor performance, due to physical and technological limitations of the three subsystems, can be hidden to some extent through architectural complications and data buffering.

Amazingly, the gap between the bandwidth and the latency of the three subsystems has widened significantly during the past three decades, increasing the pressure on computer architects. Processor bandwidth has improved by a factor of 10 000 – 25 000 versus a factor of 300 – 1 200 for memory and disks. The corresponding latency improvements are 30 – 80 times for processors versus 6 – 8 times for memory and disks [232].

Hiding the latency of slower subsystems is the holy grail of computer architecture. Latency hiding leads to architectural complications, eventually creating the hardware vulnerabilities discovered a few years ago, a topic discussed in depth in Section 3.2.

ISA—Instruction Set Architecture. ISA is an abstract model of a processor defining the set of instructions executed by the CPU. ISA reflects compiler writer's and programmer's view of the processor. CISC (Complex Instruction Set Computer) and RISC (Reduced Instruction Set Computer) are two competing CPU architectures.

CISC architecture is characterized by a large number of machine instructions, some operating with data in registers and others directly with data in the system's memory. CISC leads to a smaller footprint of the application code, but complicates CPU implementation and limits the efficiency of pipelining and other CPU optimization techniques. CISC architecture was used by IBM System 360 and 370 in the late 1960s. x86-64 is a 64-bit architecture developed by Intel as a successor of its 32-bit architecture, x86-32. x86-64 is implemented by Intel and AMD processors. x86-64 CISC instructions are not executed by the hardware, but compiled into microcode; then, the CPU executes the microcode.¹

RISC architecture, true to the name, limits the number of instructions. RISC is a *load-store architecture*. RISC ISA includes two instructions, one to fetch data into CPU registers and the other to store data from CPU registers in memory; all arithmetic and logic instructions operate only with data in registers. RISC architecture was introduced in the early 1980s. Power PC, MIPS, and SPAC are three ISA RISC architectures conceived independently by John Cocke at IBM Research, John Hennessy at Stanford, and David Patterson at UC Berkeley. The contribution of all three was recognized with Turing Awards, the computer science equivalent of a Nobel Prize, John Cocke in 1987 and the others in 2017. John Hennessy and David Patterson are the authors of a seminal text on computer architecture; they pioneered the transition from *qualitative* to *quantitative* methods in computer architecture emphasizing objective measurements. Qualitative methods provide only some insights and identify factors affecting processor performance.

Control flow versus data flow processor architecture. The dominant processor architecture is the *control flow* architecture pioneered by John von Neumann [79]. The implementation of the processor control flow is straightforward: The *program counter* (PC) determines the next instruction to be loaded into the *instruction register* (IR) and then executed. The execution is strictly sequential, until a branch is encountered.

There is an alternative to control flow architecture, the *data flow* architecture when an operation is executed as soon as the data required by the operation becomes available. Though only a few general

¹ x86 instructions are internally converted into simpler RISC-style micro-operations specific to a particular processor and stepping level.

purpose data flow systems are available today, this alternative computer architecture is widely used by network routers, digital signal processors, and other special-purpose systems. The fastest way to get the results of any computation is using a data flow model, i.e., carry out each operation as soon as the input data becomes available.

Still, the von Neumann architecture based on the control flow model, is implemented by the vast majority of computers today. Why? Lack of locality, inefficient use of cache, and ineffective pipelining are most likely some of the reasons why data flow general-purpose processors are not as popular as control flow processors. It should not be surprising that some of the systems discussed in Chapters 4, 9, and 11 apply the data flow model for task scheduling on large clusters. The power of this model for supporting optimal parallel execution is unquestionable, and we should probably expect soon the addition of general-purpose data flow systems to the cloud infrastructure.

Processor clock, CPI, and IPC. The basic subsystems of a CPU are a Control Unit (CU), an ALU (Arithmetic and Logic Unit), and a Register File (RF). All CPU activities are controlled by an internal clock. The CU implements a state machine and at each clock cycle dictates what latches should open along the control paths to different circuits and what actions should be carried out. R , the clock rate, is measured in GHz (Gigahertz) (a Hz is one cycle/second and one GHz is 10^9 Hz). CPU bandwidth is determined by IPC (Instructions per Clock Cycle) or CPI (Cycle per Instruction), $IPC = 1/CPI$. The larger the IPC, the more performant the processor architecture. T , the execution time of an application, is a function of: N , the number of instructions executed, CPI, and R:

$$T = N \times CPI \times R.$$

The only way to decrease the execution time of applications is to increase CPI, R, or both. Increasing the clock rate, R, was the preferred method to boost processor performance until 2003–2006, when the challenges of heat removal forced an abrupt halt to the frenetic pace of clock-rate increase; the clock rate increased from a high of 22 MHz in 1986 (VAX 8600) to 3.2 GHz in 2003 (Intel Xeon EE). The energy consumption of solid-state devices is proportional to the second or third power of the clock rate depending upon the technology, thus increasing the clock rate by two orders of magnitude led to a dramatic increase of energy dissipated and of heat generated by processors.

The implacable laws of physics show that energy dissipation and heat removal will eventually seal the fate of solid-state technology. The arguments are very simple: Transistors must be densely packed to minimize communication time; the speed of light cannot be exceeded. A sphere allows optimal transistor packing; then, the energy dissipated by the transistors is proportional to the volume of the sphere ($V = 4/3\pi r^3$); thus, the heat generated is proportional to the cube of the radius r of the sphere. The heat can only be removed through the surface of the sphere ($S = 4\pi r^2$), thus, the amount of heat removed is proportional to the square of the radius.

The only alternative to boost system performance is to increase the CPI of each CPU and to have multiple cores. This is feasible because the large number of transistors on a chip predicted by Moore's Law allows the implementation of various architectural optimizations and of multiple cores. The first step to increase the CPI was to exploit bit-level and instruction-level parallelism implemented since the early days of computing. Multicore processors were developed much later. In 2001, IBM introduced the world's first multicore processor, a chip with two 64-bit microprocessors comprising more than 170 million transistors. The total number of instructions executed by a multicore processor is the sum of instructions executed by individual cores.

Table 3.1 Basic superscalar processor pipeline. The pipeline has five stages: instruction fetch (IF), instruction decode (ID), instruction execution (EX), memory access (MEM), and write back (WB). In each clock cycle two instructions are executed; instructions i , $i + 2$, $i + 4$, $i + 6$ and $i + 8$ are executed by unit 1, and instructions $i + 1$, $i + 3$, $i + 5$, $i + 7$ and $i + 9$ are executed by unit 2. Once the pipeline is full, two instructions complete execution of every clock cycle.

	1	2	3	4	5	6	7	8	9
i	IF	ID	EX	MEM	WB				
$i + 1$	IF	ID	EX	MEM	WB				
$i + 2$		IF	ID	EX	MEM	WB			
$i + 3$		IF	ID	EX	MEM	WB			
$i + 4$			IF	ID	EX	MEM	WB		
$i + 5$			IF	ID	EX	MEM	WB		
$i + 6$				IF	ID	EX	MEM	WB	
$i + 7$				IF	ID	EX	MEM	WB	
$i + 8$					IF	ID	EX	MEM	WB
$i + 9$					IF	ID	EX	MEM	WB

Bit-level and instruction-level parallelism. Parallelism at different levels can be exploited by a von Neumann processor. The first two levels are:

1. *Bit-level parallelism.* A computer word is handled as a unit by the instruction set and the processor hardware. The number of bits in a word has increased gradually from 4-bit processors to 8-bit, 16-bit, and 32-bit processors. This has reduced the number of instructions required to process larger size operands and enabled a dramatic performance improvement. During this evolutionary process, the number of address bits have also increased from 32 bits to 64 bits in 2004 allowing instructions to reference a larger address space, from 2^{32} (about 4 GB), to 2^{64} (17 179 869 184 GB).
2. *Instruction-level parallelism (ILP).* Computers used multi-stage pipelines to speedup execution. *Once an n-stage pipeline is full, an instruction is completed at every clock cycle unless the pipeline is stalled.* This idea mimics assembly lines used in manufacturing as early as 1913 when Henry Ford introduced this innovation at the Highland Park assembly plant.

Pipelining means splitting an instruction into a sequence of steps that can be executed concurrently by multiple units on the chip. A basic pipeline of a RISC (Reduced Instruction Set Computing) architecture consists of five stages.²

A *superscalar processor* executes more than one instruction per clock cycle, see Table 3.1 [232]. A Complex Instruction Set Computer (CISC) architecture could have a much larger number of pipelines stages, e.g., an Intel Pentium 4 processor has a 35-stage pipeline.

The events in each pipeline stage involve multiple hardware components including: the PC (Program Counter), a register pointing to the next instruction to be executed, the IR (Instruction Register)

² The number of pipeline stages in different RISC processors varies. For example, ARM7 and earlier implementations of ARM processors have a three-stage pipeline, fetch, decode, and execute. Higher performance designs, such as the ARM9, have deeper pipelines: Cortex-A8 pipeline has 13 stages.

containing the current instruction, the register file, and the ALU. A generic descriptions of each stage of a five stage pipeline is:

1. IF—fetch the instruction and store it into IR, compute new program counter, store the incremented program counter in the PC register and into a pipeline register for later computing the branch target address if necessary.
2. ID—decode the instruction, fetch the data required by the instruction, and store it in the registers specified by the instruction in IR.
3. EX—perform an ALU operation or an address calculation and the condition code for the ALU operation; compute the target address if the instruction is a taken branch.
4. MEM—cycle the memory; write the PC if needed and pass along values needed in the final stage.
5. WB—update the register field from either the ALU output or the loaded value.

The execution flow is different for arithmetic and logic, load, store, and control and branch instructions. ALU instructions access two register operands; load and branch instructions access a register to obtain a base address; and store instructions access a register to obtain the register operand and another register for the base address.

Pipeline hazards occur when unchecked pipelining would produce incorrect results. Data, structural, and control hazards have to be handled carefully. *Data hazards* occur when the instructions in the pipeline are dependent upon one another. For example, a Read after Write (RAW) hazard occurs when an instruction operates with data in a register being modified by a previous instruction. A Write after Read (WAR) hazard occurs when an instruction modifies data in a register being used by a previous instruction. Finally, a Write after Write (WAW) hazard occurs when two instructions in a sequence attempt to modify data in the same register and a sequential execution order is violated.

Structural hazards occur when the circuits implementing different hardware functions are needed by two or more instructions at the same time. For example, a single memory unit is accessed during the instruction fetch stage when an instruction is retrieved from memory, and it is also accessed during the memory stage when data is written and/or read from memory. Structural hazards can be resolved by separating the component into orthogonal units (such as separate caches) or bubbling the pipeline. *Control hazard* are due to conditional branches. On some pipelined microarchitectures, a processor does not know the outcome of a branch when it needs to insert a new instruction into the pipeline during the fetch stage.

A *pipeline stall* is the delay in the execution of an instruction in the pipeline to resolve a hazard. Such stalls could drastically increase the CPI according to the expression:

$$\text{Pipeline CPI} = \text{Ideal Pipeline CPI} + \text{Structural Stalls} + \text{Data Hazard Stalls} + \text{Control Stalls}$$

Pipeline scheduling separates the dependent instruction from the source instruction by the pipeline latency of the source instruction. Its effect is to reduce the number of stalls.

Computer architecture taxonomy. In 1966, Michael Flynn proposed a taxonomy of computer architectures based on the number of *concurrent instructions* and the number of *data streams*: SISD (Single Instruction, Single Data); SIMD (Single Instruction, Multiple Data); and MIMD (Multiple Instructions, Multiple Data).

SISD processors have been around since ENIAC, the computer built at the University of Pennsylvania's Moore School of Electrical Engineering between 1943 and 1946 by J. Presper Eckert and John Mauchly.

Individual cores of a multicore processor are SISD (unless multithreaded) and support the execution of one or more threads at any given time. Multiple threads may run concurrently but only one is active at any given time. A *superscalar* processor executes more than one instruction per clock cycle. A single-core superscalar is still a SISD processor.

SIMD supports vector processing. When a SIMD instruction is issued, the operations on individual vector components are carried out concurrently. For example, to add two vectors $(a_0, a_1, \dots, a_{63})$ and $(b_0, b_1, \dots, b_{63})$, all 64 pairs of vector elements are added concurrently, all sums $(a_i + b_i)$, $0 \leq i \leq 63$ are available at the same time.

The first use of SIMD instructions was in vector supercomputers, such as the CDC Star-100 and Texas Instruments ASC in the early 1970s. Vector processing was especially popularized by Cray in the 1970s and 1980s by attached vector processors such as those produced by Floating Point Systems and by supercomputers such as Thinking Machines CM-1 and CM-2. The first widely deployed SIMD instruction set for gaming was Intel's MMX extensions to the *x86* architecture. IBM and Motorola then added AltiVec to the POWER architecture. There have been several extensions to the SIMD instruction sets for both architectures as we shall see in Section 3.4.

MIMD refers to a system with several processors that function asynchronously and independently; at any one time, different processors may be executing different instructions on different data. Multicore processors support true MIMD execution. Each core has its own register file, ALU, and floating-point execution units. Each core has its private L1 instruction and data caches, as well as L2 cache. All processor cores share the L3 cache.

Several processors can share a common memory, and we distinguish several types of multiprocessor systems: UMA, NUMA, and COMA, uniform, non-uniform, and cache only memory access, respectively. A MIMD system could have a distributed memory. Processors and memory modules communicate with one another using an interconnection network, such as a hypercube, a 2D torus, a 3D torus, an omega network, or another network topology. Today, most supercomputers are MIMD machines, and some use GPUs instead of traditional processors. Multicore processors with multiple processing units are now ubiquitous. As more powerful processors are needed, *hyper-threading* was introduced for Xeon and later Pentium 4 processors by Intel in 2002. Hyper-threading takes advantage of unused processor resources and presents itself to the operating system as a two-core processor.

MISD (Multiple Instructions, Single Data) is a fourth possible architecture, but it is very rarely used, mostly for fault tolerance.

Performance metrics and benchmarks. MIPS (Million Instructions per Second) quantify processor performance for integer arithmetic, while MFLOPS (Million Floating Point Operations per Second) quantify processor performance for floating-point arithmetic. The information provided by MIPS and MFLOPS is of limited usefulness because processor performance is application-specific. An application can be CPU-, memory-, I/O-, or all-intensive; the application will run optimally on a processor where the corresponding subsystem, CPU, memory, and I/O is optimized. This is the reason why AWS offers various types of instances, e.g., compute-optimized, memory-optimized, storage-optimized.

Applications running on the same processor may run at very different MIPS or MFLOPS rates depending on the type of resources intensively used. The only useful performance characterization of a processor is provided by *benchmarks*, suites of applications with similar characteristics. PARSEC benchmark suite developed at Princeton by Christian Bienia for his Ph.D. dissertation (<http://parsec.cs.princeton.edu>) is widely used. PARSEC compares favorably with other benchmarks, including SPEC CPU 2017, BioParallel, PhusicsBench, SPLAH-2, among others.

3.2 Grand architectural complications

In this section we assume a RISC architecture though there are only minor differences for CISC architecture. Two conditions guarantee computation completion in the shortest possible time: A—memory optimization: the instruction to be executed next should be in the IR register; after the instruction is decoded, the data referenced by instruction should be in the CPU registers before the instruction execution stage can begin; B—CPU optimization: once condition A is satisfied, we should keep the pipeline full at all times.

To address condition A, we have to consider memory latency, bandwidth, and cost. *Memory latency* is the time it takes to fetch one instruction, or a word or byte of data, measured in picoseconds (psec), nanoseconds (nsec), or milliseconds (msec); *memory bandwidth* is the amount of data delivered in one unit of time, measured in GB/sec or MB/sec. Ideally, the storage should have infinite capacity and bandwidth and zero latency and cost, but in reality, the storage has limited capacity, finite bandwidth, and non-zero latency. Memory cost is also a factor: The faster, the higher is the cost of physical memory. Computers need large memories because the code footprint, as well as the size of data, have been continually increasing. The cost of very large and very fast memory is prohibitive, and solutions balancing the memory performance and cost have been developed.

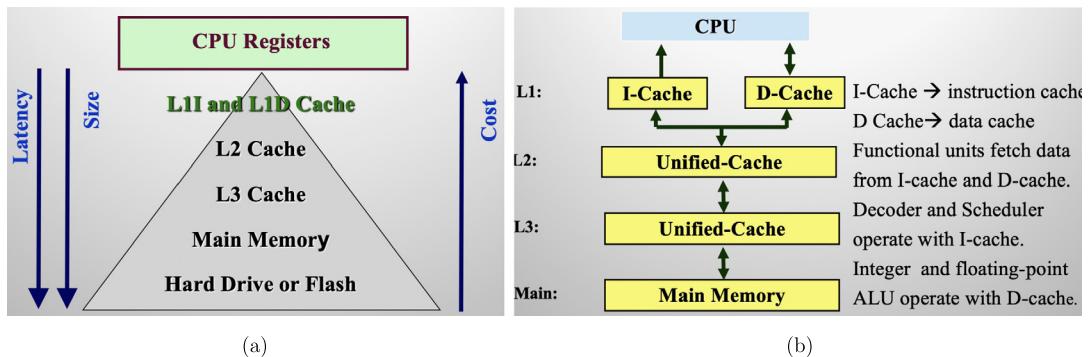
Keeping the pipeline full at all times has major implications:

- B1—avoid pipeline stalls;
- B2—deal with the fact that different arithmetic operations require a different number of clock cycles; for example, an integer addition may require two clock cycles, a multiplication 12 clock cycles, and a division 50–60 clock cycles;
- B3—avoid wasting clock cycles due to control flow instructions which alter the sequential execution pattern. A control flow instruction could be a condition generated by an arithmetic operation that modifies condition codes and may or may not force us to fetch a new instruction, or a subprogram call or an instruction that transfers control to a subprogram or elsewhere in the address space of the program.

A paramount complication: The architecture must ensure program correctness. This means that the architecture should: (i) *Preserve exception behavior*, any change in instruction execution order must not change the order in which exceptions are raised³; (ii) *preserve instruction flow*, the flow of data between instructions that produce results and consume them; and last but not least (iii) *preserve security*, i.e., it does not open side channels allowing an intruder to access data in the address space of the executable. Not even a data flow architecture could easily overcome the problems posed by B1–B3 and by the requirements for program correctness.

A. Memory hierarchy; data and instruction caches; cache organization and performance. The first challenge seems the easiest of the two, so we start with it. The goal is to decrease memory latency to match the CPU speed. The solution is to access blocks of code and data rather than individual items. Fortunately, code and data enjoy a very useful property, *locality*. There are two flavors of locality, *spatial* and *temporal*.

³ An *exception* is an unexpected event within the processor, e.g., an arithmetic exception or an addressing exception, while an *interrupt* is an unexpected outside event, e.g., an I/O event or a timing interrupt. Whenever an exception or an interrupt occurs, the hardware stops the execution of the running process and executes the code performing an action in response to the exception.

**FIGURE 3.1**

(a) Memory hierarchy; CPU registers have the shortest latency followed by several levels of cache. (b) Cache levels: the first level cache consists of L1I, instruction cache and L1D, and data cache; L2 and L3 level caches are unified, used for instruction and data.

Spatial locality means that, once a block of code is brought into a faster storage, there is a high probability that several, possibly many instructions in the block will be executed before another block of code is referenced. *Temporal locality* means that, once a block of code is brought into the faster storage, several, hopefully many, instructions in this block of code will be executed; this is a reasonable expectation because instructions are executed sequentially most of the time. Similar arguments hold for data locality: Both spatial and temporal data locality reflect the logical organization of data.

Handling blocks of code rather than individual instructions and blocks of data rather than words of data decreases the average latency per reference. Temporal locality has an important corollary, captured by the *working-set* concept. The working set of a program is a collection of blocks that have been recently referenced. A program with a small working set references a smaller number of data blocks in an interval of time Δt than one with a larger working set.

The first idea that comes to mind is to design a storage hierarchy consisting of small, very fast, and expensive memory working in concert with increasingly larger and less expensive memory. The fast memory, called *cache*, was introduced by Maurice Wilkes from the University of Cambridge, who received the Turing Award in 1967 for his lifelong contributions to computer architecture.

Fig. 3.1(a) displays the memory hierarchy and the bandwidth, latency, cost-per-bit, and the storage technology (SRAM—Static Random Access Memory, DRAM—Dynamic Random Access Memory, and magnetic) at each level of storage hierarchy. Memory technology, bandwidth, latency, and cost are:

CPU registers—multipoint SRAM, 200+ GB/sec, 300+ psec, \$ 0.01 per bit.

On-chip L1I or L1D cache—SRAM, 50+ GB/sec, 300+ psec, \$ 10^{-4} per bit.

On-chip L2 cache—SRAM, 10+ GB/sec, 2+ nsec, \$ 10^{-4} per bit.

Main memory—DRAM, 2+ GB/sec, 50+ nsec, $< \$ 2 \times 10^{-7}$ per bit.

Hard disk—magnetic, 10+ MB/sec, 10 msec, \$ 1×10^{-9} per bit.

L1 cache capacity is several hundred KB (384 KB for AMD Ryzen 5 5600X, 1 MB for Intel i9-9980XE); L2 cache capacity 1–10 MB, and L3 10–50 MB. Fig. 3.1(b) shows a further refinement of storage hierarchy with multiple cache levels. The fastest cache consists of on-chip instruction cache,

L1I, and on-chip data cache, L1D. The other cache levels, L2 and L3, are unified instruction and data caches. Recently, L4-level cache was added to the cache hierarchy. Different cache levels are accessed by ALU, instruction decoder, and scheduler.

Cache implementation requires a number of important decisions: What should the block size be, where to place a block brought in to cache, what to do when the cache is full and a new block must be fetched from memory, and how to handle a modified cache block?

Cache organization elements are: (i) *cache block*—the unit transferred between memory and cache; (ii) *cache row*—consists of three elements for each cache block: *Tag*: contains part of the address of the item fetched from memory, the *Data block*: contains data fetched from memory, and it flags either valid or dirty; and (iii) *Cache line*—identifies the cache set for an n-way set associative cache. The mapping between memory address and cache address is done by splitting the address generated by CPU into Tag, Index, and Word offset in the block. The Index portion of the CPU-generated address identifies the cache line, and the Tag of CPU-generated address is compared with the Tag bits of the cache line. The valid bit indicates if there is a cache hit or a cash miss. In case of an *n-way set associative cache*, each cache block of the line includes all three fields, valid bit, tag, and data.

There are several options for block placement in cache: (i) *fully associative cache*—a block can be placed in any cache location; (ii) *direct map cache*—a block can be placed only in one cache location determined by the starting address of the block in memory and the cache capacity in number of blocks. For example, if cache capacity is 16 blocks, a memory block starting at address 128 should be placed in cache block 8 (128 modulo 16); (iii) *set-associative cache*—the cache is organized in sets consisting of multiple blocks. This cache organization reduces the miss rate as multiple alternatives for block placement in the set exist, e.g., a four-way set associative cache offers four-options for block placement in the cache, i.e., the first, second, third, or fourth block of the set.

For example, consider a 32-bit architecture and a processor with a 64 Kbyte cache, 32 byte cache block and one block per cache set. Then the number of blocks is $64\ 000/32 = 2\ 000$ and the 32 bits of the address generated by the CPU are used by the cache management as follows: bits 0 – 4 give the offset of a word in a block of 32 bytes; bits 5 – 15 give index of one of the 2 000 cache blocks; and bits 16 – 31 are the block tag. Each cache line consists of one valid bit, the block tag, and 32 bytes of data. If another processor had the same cache and block size but a two-way set-associative cache there will be 1 000 cache lines, each line will have one valid bit, a 17 bit tag and 32 bytes of data for each one of the two blocks in the set.

When CPU generates an address in the address space of the process, the cache is accessed first. The address can be in a block already in cache, and we experience a *cache hit*, or a new block must be fetched from memory and a *cache miss* occurs. There are several types of cache misses: *compulsory* - first reference to a block; *capacity* - blocks discarded and later retrieved; and *conflict* - repeated references to multiple addresses from different blocks that map to the same location in the cache.

Cache capacity is limited: It is two to three orders of magnitude smaller, KB or MB as shown in the caption of Fig. 3.1, versus GB for storage capacity. Once the cache is full, the block to be evicted depends upon the placement strategy that dictates where the new block should reside. The block to be evicted can be dirty, i.e., modified, different from its copy left in memory, so the memory must be updated. A dirty bit will identify a modified cache block. There are two choices for handling the modified cache block: *Write through cache*—a modified cache block is immediately written to memory and *write back cache* when memory is updated when the block must be evicted from the cache.

The *cache miss rate* is the number of cache misses per unit of time. A cache miss incurs a *miss penalty* much larger than the time to access a block already in cache, the so-called *hit time*. The average memory access time (AMAT) is:

$$AMAT = HitTime + MissRate \times MissPenalty.$$

Several cache techniques are used to reduce AMAT:

1. Small and simple L1 caches—reduce the critical cache timing path consisting of the time to extract the Index and Tag from the CPU generated address to compare the tags of the blocks on the same cache line, and to select the block.
2. Fast hit times via way prediction—keep extra bits in cache to predict the “way,” i.e., the block within the set, at the next cache access.
3. Cache pipelining—use multiple cycles to access the cache. For example, in stage 1, read the tag and valid bit, in stage 2, combine the result of stage 1 and start read if hit, and in stage 3, finish data read and transfer data to the CPU. Cache pipelining improves bandwidth, but it has higher latency and increases the branch miss-prediction penalty.
4. Increase cache bandwidth with non-blocking caches—allow data cache to continue to supply cache hits during a miss; processors can hide L1 miss penalty but not an L2 miss penalty.
5. Independent banks; interleaving—the cache is divided into a number of *cache banks* that can be accessed concurrently.
6. Early restart and critical word first—as soon as the requested word of the block arrives, send it to the CPU, and let the CPU continue execution while filling the rest of the words in the block.
7. Merging write buffers to reduce miss penalty—write buffers allow CPU to continue while waiting to write to memory.

Several other cache-like units are used to speed up instruction execution. For example, ROB (Reorder Buffer) is a cache for out-of-order instruction execution, and BTB (Branch Translation Buffer) is a cache used by the branch prediction hardware, discussed next.

TLB (Translation Look aside Buffer) is used by the virtual memory management to speed up dynamic address translation, the translation of virtual addresses generated by the CPU to physical memory addresses. Virtual memory applies ideas very similar to caching and extends the physical memory of a system allowing multiple processes, each with an address space (the universe of addresses visible to the process) larger than the processor’s physical memory, to run concurrently. In this case, the units of data exchanged between the lower level of storage hierarchy, disks, and memory are called *pages*. Without virtual memory one would have to reorganize a program to fit in the physical memory of each processor, leading to an unthinkable debacle.

This concludes our analysis of challenge A for CPI optimization. Needless to say, this presentation omits many additional complications, e.g., the fact that the address generated by the CPU is a *virtual address* forcing the cache management to interact with virtual memory management.

B. CPU optimization. Modern von Neumann microprocessors attempt to emulate the data flow model. This means to execute an instruction as soon as the data it needs is available, while guaranteeing program-order execution. Multiple instructions issue, dynamic instructions scheduling, branch prediction, speculative execution, and multi-threading are some of the complications designed to speed up execution.

Recall that CPU performance is determined by the product of three terms: the number of program instructions, the CPI, and the clock rate. Performance optimization implies CPI minimization; thus, the maximization of three types of flows involved in instruction processing are: *instruction flow*—affected by branch instructions, *register data flow*—affected by ALU instructions, and *memory data flow*—affected by load and store instructions [232].

Sequential instruction execution is interrupted by control instructions. Handling unconditional and conditional branches generates empty instruction pipeline slots. In case of unconditional branches, the next instruction cannot be fetched until the target address of the branch is calculated. The number of delay slots needed for target address generation depends upon the addressing modes of the branch instruction. The three addressing modes are: (i) *PC-relative* addressing when the branch target address can be generated during the fetch stage; (ii) *register indirect* addressing when the branch instruction must traverse the decode stage to access the register; and (iii) *register indirect with an offset* addressing when the offset must be added after register access.

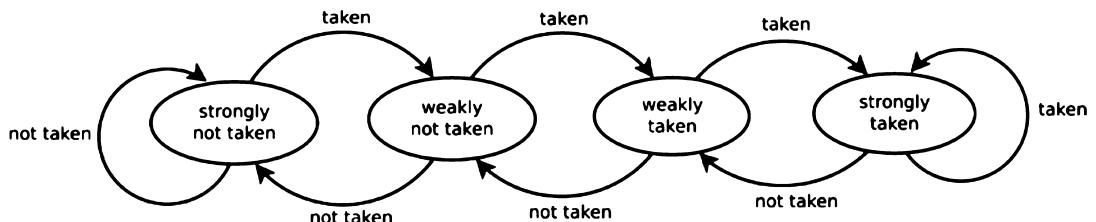
Conditional branches incur a larger penalty: the CPU must wait for the resolution of the branch condition. If the branch is taken, it must also wait until the target address is available; the delay involves several clock cycles for traversal of the decode, dispatch, and execute pipeline stages. The number of *delay slots* needed for condition evaluation of a conditional branch instruction is different when condition code registers are available or when the comparison of two general purpose registers generates the branch condition. The number of empty instruction pipeline slots is multiplied by the CPU width for a superscalar CPU.

Out-of-order instruction execution requires additional hardware, including: reservation stations, load-store buffers, a FIFO queue holding the instructions issued by functional units, a reorder buffer (ROB), and a common data bus (CDB). There are multiple reservation stations (RS) for each functional unit, e.g., floating point add/subtract, floating-point multiple/divide. A reservation station stores the instruction, buffered operand values (when available), and the ID of the reservation station providing the operand values. RS feed data to floating-point arithmetic units modified to accept ROB as a source.

ROB tracks the state of inflight instructions in the pipeline and provides the illusion of in-order program execution. After instructions are decoded and renamed, they are then dispatched to the ROB and the FIFO queue and marked as busy. When an instruction finishes execution, ROB is informed, and the status of the instruction becomes *not busy*. Once the “head” of the ROB is no longer busy, the instruction is *committed*, and its state is visible. If an exception occurs and the excepting instruction is at the head of the ROB, the pipeline is flushed, and no architectural changes that occurred after the excepting instruction are made visible. The ROB then redirects the PC to the appropriate exception handler. Register values and memory values are not written until an instruction commits.

The instruction execution steps and the actions carried in each steps are:

1. *Instruction issue*: (i) get next instruction from the FIFO queue; (ii) if a reservation station is available, issue the instruction to the reservation station including operand values, if available; and (iii) if operand values are not available, stall the instruction. When an operand value is unavailable, use a placeholder for the value. The placeholder is a tag indicating the reservation station to produce the value. The placeholder will be later replaced by the result produced by a functional which will broadcasts it on CDB.
2. *Execute instruction*: (i) when an operand becomes available, store it in reservation station(s) waiting for it; (ii) when all operands are ready, issue the instruction. To ensure correctness loads and store

**FIGURE 3.2**

Two-bit saturation counter FSM states and transitions.

are maintained in program-order through effective address and no instruction is allowed to initiate execution until all branches that precede it in program-order have completed.

3. Write instruction result: broadcast result on CDB for waiting reservation stations or for store buffers.

Branch prediction. The penalty for interrupting the sequential instruction execution by branches increases when the number of pipeline stages increases. The behavior of branch instruction is predictable, and accurate predictions lead to significant performance improvements. Two functions are required for branch prediction: *branch target speculation* - to predict the address of the next instruction and *branch condition speculation* - to predict if the branch will be taken or not. Both require additional hardware.

Branch target speculation uses a cache accessed during the instruction fetch to store the target address of the previous branch instruction, called BTB (Branch Target Buffer). BTB is indexed by the current PC, and it is accessed concurrently with the L1I-cache. BTB has two fields: (i) branch instruction address containing the address of the current branch instruction; and (ii) branch target address containing the speculative target address to be used as the new PC.

If the address of the current branch instruction has an entry in BTB, more precisely in field (i) of BTB, then field (ii) of the same BTB entry contains the predicted address when the branch is taken. There are two possibilities: (1) the current branch is taken; thus, the prediction was accurate and there is no penalty for accessing the BTB; (2) the current branch is not taken, the BTB prediction did not help, and there is a penalty of several clock cycles needed to compute the branch address.

BTB is a cache with limited capacity, thus, the address of the current branch instruction may not be in BTB. When the branch is not taken, there is no penalty; if the branch is actually taken, there is the penalty for computing the branch address, which depends on the addressing modes supported by the ISA.

Branch condition speculation can use a *biased strategy*. In this case, the branch target offset, the difference between the address of the current instruction (ci), and the address of the target instruction (ti) are used. When $ci < ti$, it is likely to have a loop closing branch, and the prediction is that the branch is taken; if $ci > ti$, the prediction is that the branch is not taken. This is not a very effective predictor.

History-based prediction uses a finite state machine (FSM) to decide if T (taken) or N (not taken). The FSM for two-bit predictor has four states, each state labeled by the last two events (TT, NT, TN, NN) and the predicted state either T or N. The prediction is changed when the prediction is wrong two consecutive times. The *two-bit saturation counter* is a better alternative to the two-bit predictor. In this case a single iteration will not change the predicted direction. For example, if a branch has been taken many times in succession, the FSM in Fig. 3.2 will be in the *strongly taken* state; if next time this

branch is not taken, the new state will be *weakly taken*, and, if the next branch is taken, it will switch back to the *strongly taken* state. Only if the branch is not taken two or more times consecutively will the prediction change to not taken. This hysteresis effect can significantly boost prediction, e.g., to 85%, for SPECint92.

Yeh and Patt algorithm for branch condition speculation looks at previous behavior when this same sequence has occurred, using historical data to generate a prediction. For each branch it maintains a storage register with the state of the last k branches, 1 if taken and 0 if not taken and a pattern table of 2^k entries each implementing a two-bit saturating counter that tracks the results of previous iterations that occurred when the history register was in a given state. When a branch is encountered, the contents of the history register are used as an index into the pattern table, selecting the entry that corresponds to the recent history of that branch. After the branch is resolved, the result (taken or not taken) is shifted into the history register and used to update the appropriate entry in the pattern table.

Correlating predictors use global history. To record what happened at the last m branches, we need a shift register with m bits, one for each branch; each bit is set to 1 if the branch was taken and to 0 if the branch was not taken. Multiple two-bit predictors for each branch are used, one predictor for each possible combination of outcomes of preceding n branches.

The concepts and ideas discussed in this section project the view of a very complex mechanism with many moving parts that have to work in concert, and they do. Billions and billions of processors used today run flawlessly on all continents and in space on satellites and space probes. Unfortunately, the interaction of cash and processor optimization discussed in this section has led to the surprising hardware vulnerabilities discussed next.

Critical vulnerabilities in modern processors. Two hardware vulnerabilities, Meltdown and Spectre, affect systems with Intel x86, IBM Power architectures, and some ARM microprocessors. *Meltdown vulnerability* exploits a race condition between memory access and privilege checking during instruction processing and affects x86-based systems. The mechanism used to create a side-channel attack is:

1. The CPU is instructed to read data at an address A that it is forbidden to access. The read is speculatively executed, and address A is not saved in the visible CPU state, but it would become visible only after access check confirms that access to A is permitted.
2. The CPU continues speculative execution and is instructed to read the contents of an address calculated using a base register B and A as displacement. The access to this address is permitted and causes caching of Base+A location.
3. The instruction in step 1 attempts to complete and save its result in visible CPU state architecture. The instruction retirement logic detects that access is forbidden, and the results produced by this and all subsequent instructions are discarded. However, the effect of caching data at the address given Base+A is not undone.
4. The attacker measures how fast elements of Base[x] array can be accessed. If the data at the address Base+A is cached, and all other elements Base+x are not; only instruction referencing Base+A execute faster. The measurement can detect the timing difference and determine the value of A obtained by the rejected instruction of step 1, thus allowing the attacker to access the contents of the forbidden memory address.

Spectre enables a whole class of potential vulnerabilities by exploiting branch prediction. The speculative execution resulting from a branch miss-prediction enables an attacker to trick a program into

accessing arbitrary locations in the program's address space and potentially obtain sensitive data. The attacker trains the branch prediction logic to reliably hit or miss cache, accurately times the difference between cache hits and cache misses, and uses this information to open a covert channel.

The existing code in the address space of the victim process is searched for places where speculation touches upon otherwise inaccessible data. Then the attacker manipulates the processor into a state where speculative execution has to touch that data. Next, it times the side effect of the processor being faster and determines if its by-now-prepared prefetch machinery indeed did load a cache line. This attack can also be remotely conducted using the following sequence of actions: flush cache, mis-train the branch predictor, and time the read operations.

Software patches to limit the impact of these hardware vulnerabilities have been released making even harder the exploitation of the two vulnerabilities. Only very sophisticated attackers could exploit these vulnerabilities, but in the age of state-sponsored cyber terrorism such attacks could affect the critical infrastructure of a country.

The conclusion of this section is obvious: Advances in computer architecture have increased the complexity of processor design, sometimes with the unexpected and dangerous side effects just discussed. Solid-state technology will unquestionably reach its limits and will most likely be replaced by quantum technology, a super-disruptive set of concepts and ideas that requires a very different way of thinking. The wealth of ingenious solutions for CPU optimizations will then be forgotten.

3.3 ARM architecture

There is no better way to conclude the discussion of general-purpose processors and illustrate the dramatic evolution of computer architecture than to overview the ARM architecture. ARM—Advance RISC Machine—is a state-of-the-art microarchitecture for a broad range of applications and workloads. ARM microprocessors are used not only for mobile devices such as smartphones, tables, and laptops, but also by some cloud instances. ARM architecture has some notable features:

- a. ARM and Thumb are two different instruction sets supported by ARM cores with a “T” in their name. For instance, ARM7 TDMI supports Thumb mode. ARM instructions are 32 bits wide, and Thumb instructions are 16 wide. Thumb mode allows for code to be smaller and potentially faster when the target has slow memory.
- b. Uniform 16×32 -bit register file including program counter, stack pointer, and link register. Fixed instruction width of 32 bits eases decoding and pipelining at the cost of decreased code density; the 16-bit Thumb instruction increases code density.
- c. Mostly single clock-cycle execution.
- d. Powerful indexed addressing modes.
- e. Support unaligned accesses for halfword and single-word load/store instructions with some limitations, such as no guaranteed atomicity.
- f. Fast leaf function calls use a link register. When calling the leaf function, the return address does not have to be pushed to the stack since it's stored in the link register. However, there are situations where the leaf function must save data to the stack.
- g. Arithmetic instructions alter condition codes *only when desired*.
- h. Conditional execution of most instructions reduces branch overhead and compensates for the lack of a branch predictor in early chips.

- i. 32-bit barrel shifter that can be used with most arithmetic instructions and address calculations without performance penalty. A barrel shifter is used to shift and rotate n-bits typically within a single clock cycle.
- j. Almost every ARM instruction has a conditional execution feature, the predication, implemented with a four-bit condition code selector (the predicate). To allow for unconditional execution, one of the four-bit codes causes the instruction to be always executed. Most other CPU architectures only have condition codes on branch instructions.
- k. Integer add, subtract, and multiply arithmetic operations; ARMv7-M and ARMv7E-M versions of the architecture also support divide operations, while others, including ARMv7-A, only optionally support it.
- l. Simple and fast, two-priority-level interrupt subsystem has switched register banks.

There are several CPU modes, depending on the implemented architecture features. At any moment in time, the CPU can be in only one mode, but the mode can be switched in response to external events or programmatically. These modes are:

1. User—the only non-privileged mode.
2. FIQ—privileged mode entered whenever the processor accepts a fast interrupt request.
3. IRQ—privileged mode entered whenever the processor accepts an interrupt.
4. Supervisor (SVC)—privileged mode entered when the CPU is reset or when an SVC instruction is executed.
5. Abort—privileged mode entered whenever a prefetch abort or data abort exception occurs.
6. Undefined—privileged mode entered whenever an undefined instruction exception occurs.
7. System (ARMv4 and above)—the only privileged mode not entered by an exception. It can only be entered by executing an instruction that explicitly writes to the mode bits of the Current Program Status Register (CPSR) from another privileged mode (not from user mode).
8. Monitor (ARMv6 and ARMv7 Security Extensions, ARMv8 EL3)—introduced to support Trust-Zone extension in ARM cores.
9. Hyp (ARMv7 Virtualization Extensions, ARMv8 EL2)—hypervisor mode supporting Popek and Goldberg virtualization requirements for non-secure CPU operation.
10. Thread (ARMv6-M, ARMv7-M, ARMv8-M)—can be specified as either privileged or unprivileged. Whether the Main Stack Pointer (MSP) or Process Stack Pointer (PSP) is used can also be specified in CONTROL register with privileged access. This mode is designed for user tasks in RTOS (Real Time OS) environment, but it's typically used in bare-metal for super-loop.
11. Handler (ARMv6-M, ARMv7-M, ARMv8-M)—mode dedicated for exception handling (except RESET handled in Thread mode).

The *state of a processor/core* is given by the contents of the register file. This state is saved at the time of a context switch and restored when the interrupted process or thread is scheduled to run again. The overhead of a context switch depends upon the number of registers and this overhead is low because the ARM register file consists of a small number of registers. Register functions are: registers R0 through R7 are the same across all CPU modes and are never banked; registers R8 through R12 are the same across all CPU modes except FIQ mode. FIQ mode has its own distinct R8 through R12 registers. R13 and R14 are banked across all privileged CPU modes except system mode, that is, each mode that can be entered because of an exception has its own R13 and R14. These registers generally

contain the stack pointer and the return address from function calls, respectively. R13 is also referred to as SP, the Stack Pointer, R14 is also referred to as LR, the Link Register and R15 is also referred to as PC, the Program Counter.

The Current Program Status Register (CPSR) contains critical CPU state information, the CPU execution mode, whether the CPU can be interrupted or not, the condition code set by arithmetic and logic instruction, and so on. The 32 bits of CPSR are: processor mode bits, M (0 to 4); Thumb state bit, T (bit 5); FIQ disable bit, F (bit 6); IRQ disable bit, I (bit 7); imprecise data abort disable bit, A (bit 8); data endianness bit, E (bit 9); if–then state bits, T (bits 10 to 15 and 25, 26); greater-than-or-equal-to bits, GE (bits 16 to 19); do not modify bits, DNM (bits 20 to 23); Java state bit, J (bit 24); sticky overflow bit, Q (bit 27); overflow bit, V (bit 28); carry/borrow/extend bit, C (bit 29); zero bit, Z (bit 30); and negative/less than bit, N (bit 31).

Register banking refers to providing multiple copies of a register at the same address. FIQ stands for Fast Interrupt reQuest, and it is basically a higher priority interrupt. FIQ will always have precedence over regular interrupts; regular interrupts won't mask or interrupt an FIQ, while an FIQ will mask or interrupt any IRQ (Interrupt reQuest).

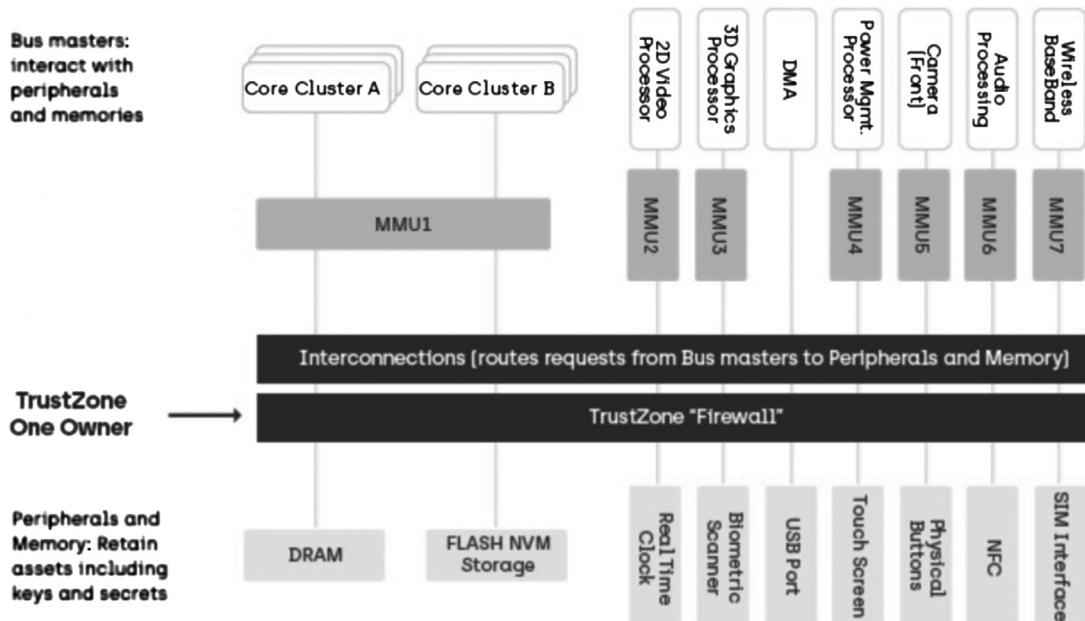
One cannot conclude even an overview of ARM architecture without mentioning the TrustZone, a technology to build a Trusted Execution Environment (TTE) in a cost-effective manner without complicating the development of different components of a System on a Chip (SoC), see also Section 3.7. Traditional isolation for a system with only one Operating System (OS) managing system resources and functionality is MMU-based. The Memory Management Unit (MMU) splits memory into regions isolated from each other, but this solution does not work well for mobile device with several cores.

TrustZone provides a hardware solution for separating a rich operating system (OS), supporting a massive set of features and consisting of a massive amount of code, from a considerably smaller but secure OS hosting the TTE. A firewall separates the TTE from a normal environment running general code. This solution allows independent bus masters, possibly running different OS on different cores, to support the isolation required by security as shown in Fig. 3.3 (from <https://www.trustonic.com/technical-articles/what-is-trustzone>). Moreover, the different operating systems do not need to be modified for the implementation of a security subsystem; at the same time, TrustZone encourages applications to separate security from general purpose aspects.

ARM implementation is hardwired without microcode. ARM architecture is cheap to produce and has low energy consumption and astounding performance. It is not surprising that more than 150 billion ARM microprocessors have been shipped since 2019, roughly 20 for every inhabitant of planet Earth [257]!! Once dominant, Intel's x86-64 architecture now lags behind ARM in popularity, as shown by Table 3.2.

ARM Ltd. founded in 1990 is located outside Cambridge, UK, and is a joint venture of Acorn Computing, Apple, and VLSI Technology. ARM licenses the architecture to companies who design and produce the microprocessors after paying the ARM licenses; in FY 2018, ARM received US \$1.8 billion from royalties [257]. The 2020 acquisition of ARM by NVIDIA for US \$40 billion is likely to have a significant impact on the computer-game industry and expanding AI computing in computer clouds and beyond.

This outline aims to offer convincing arguments for the elegance, power, and energy consumption frugality rather than an in-depth presentation of an architecture designed to be flexible, supporting virtualization, real-time execution, and security.

**FIGURE 3.3**

ARM hardware solution for security. The TrustZone Firewall allows independent bus masters to support isolation.

Table 3.2 Shipments of ARM and x86-64 units [408].

Year	2010	2011	2012	2013	2014	2015	2016
ARM (in billions)	6.1	7.9	8.7	10.4	12.0	14.8	17.1
x86-64 (in millions)	8.9	9.52	9.67	9.89	10.09	11.09	11.1

3.4 SIMD architectures

SIMD architectures have significant advantages over the other systems described by Flynn's taxonomy. Some of these advantages are:

- Exploit a significant level of data-parallelism. Enterprise applications in data mining and multimedia applications, as well as the applications in computational science and engineering using linear algebra, benefit the most.
- Allow mobile device to exploit parallelism for media-oriented image and sound processing using SIMD extensions of traditional ISA.
- Are more energy efficient than MIMD architecture. Only one instruction is fetched for multiple data operations, rather than fetching one instruction per operation.
- Have a higher potential speedup than MIMD architectures. SIMD potential speedup could be twice as large as that of MIMD.
- Allow developers to continue thinking sequentially.

Three flavors of the SIMD architecture are encountered in modern processor design: vector architecture; SIMD extensions for mobile systems and multimedia applications; and Graphics Processing Units.

Vector architectures use vector registers holding 64, 128, 256, or more vector elements. Vector functional units carry out arithmetic and logic operations using data from vector registers as input and disperse the results back to memory. Vector load-store units are pipelined, hide memory latency, and leverage memory bandwidth. The memory system spreads access to multiple *memory banks*, which can be addressed independently.

Vector length registers allow handling of vectors when the number of elements is not a multiple of the physical vector registers size, e.g., a vector with 100 elements when the vector register can only contain 64 vector elements. *Vector mask registers* disable/select vector elements and are used by conditional statements that select specific vector elements.

Non-adjacent vector elements of a multidimensional array can be loaded into a vector register by specifying the *stride*, the distance between elements to be gathered in one register. Scatter-gather operations support processing of sparse vectors. A *gather* operation takes an *index vector* and fetches the vector elements at the addresses given by adding a base address to the offsets given by the index vector; as a result, a dense vector is loaded in a vector register. A *scatter* operation does the inverse, it scatters the elements of a vector register to addresses given by the index vector and the base address.

Chaining allows vector operations to start as soon as individual elements of vector source operands become available and operate on *convoys*, sets of vector instructions that can potentially be executed together. Multiple *lanes* process several vector elements per clock cycle. Each lane contains a subset of the vector register file and one execution pipeline from each functional unit.

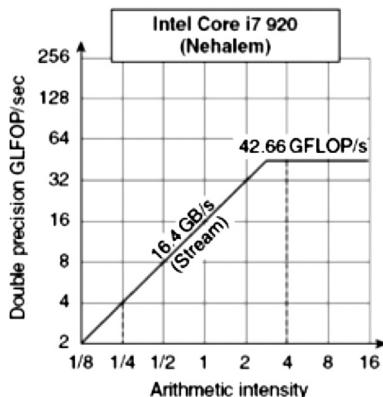
SIMD extensions for multimedia applications. The name of this class of SIMD architectures reflects the basic architectural philosophy—augmenting an existing instruction set of a scalar processor with a set of vector instructions. SIMD extensions have obvious advantages over vector architecture:

1. Low cost to add circuitry to an existing ALU.
2. Little extra state is added, thus the extensions have little impact on context-switching.
3. Need little extra memory bandwidth.
4. Do not pose additional complications to the virtual memory management for cross-page access and page-fault handling.

Multimedia applications often run on mobile devices and operate on narrower data types than the native word size. For example, graphics applications use three 8-bit for colors and one 8-bit for transparency, audio applications use 8, 16, or 24-bit samples. To accommodate narrower data types, carry chains have to be disconnected. For example a 256-bit adder can be partitioned to perform simultaneously 32, 16, 8 or 4 additions on 8, 16, 32, or 64 bit, respectively. The instructions *opcode* now encode the data type and neither sophisticated addressing modes supported by vector architectures, such as stride-base addressing or scatter-gather, nor mask registers, are supported.

Intel extended its *x86-64* instruction set architecture, and in 1996 introduced MMX (Multi-Media Extensions) which supports eight 8-bit, or four 16-bit integer operations. MMX was followed by multiple generations of streaming SIMD extensions (SSE) in 1999 and ending with SSE4 in 2007. The SSEs operate on eight 8-bit integers, four 32-bit or two 64-bit either integer or floating-point operations.

AVX (Advanced Vector Extensions) introduced by Intel in 2010 operates on four 64-bit either integer or floating-point operations. Several members of the AVX family of Intel processors are: Sandy

**FIGURE 3.4**

Rooftline performance model for Intel i7 920. When $ArI < 3$, the memory bandwidth of 16.4 GB/sec is the bottleneck. The processor delivers 42.66 Gflops, and this limits the performance of applications with arithmetic intensity larger than three.

Bridge, Ivy Bridge, Haswell, Broadwell, Skylake, and its follower, the Baby Lake, announced in August 2016. AMD offers several processor families with multimedia extensions, including the Steamroller.

Floating-point performance models for SIMD architecture. The gap between the processor and the memory speed, though bridged by different level of caches, is still a major factor affecting the performance of many applications. Applications displaying low spatial and temporal locality are particularly affected by the gap. The effects of this gap are also most noticeable for SIMD architecture and floating-point operations. *Arithmetic intensity* (ArI), defined as the number of floating-point operations per byte of data read, is used to characterize application scalability and to quantify the performance of SIMD systems.

Arithmetic intensity of applications involving dense matrices is high, and this means that dense matrix operations scale with problem size, while sparse matrix applications have a low arithmetic intensity and do not scale well with the problem size. Applications involving spectral methods and FFT (Fast Fourier Transform) have an average arithmetic intensity.

The *rooftline* model captures the fact that the performance of an application is limited by its arithmetic intensity and the memory bandwidth. A graph depicting the floating-point performance function of the arithmetic intensity is shown in Fig. 3.4. Memory bandwidth limits the performance at low arithmetic intensity, and this effect is captured by the sloped line of the graph. As arithmetic intensity increases, the floating-point performance of the processor is the limiting factor that is captured by the straight line of the graph.

3.5 Graphics processing units

The desire to support real-time graphics with vectors of two, three, or four dimensions led to the development of Graphics Processing Units, which are very efficient at manipulating computer graph-

ics. GPUs produced by NVIDIA, and AMD/ATI are also used in embedded systems, mobile phones, personal computers, workstations, and game consoles. GPU processing is based on a heterogeneous execution model with a CPU acting as the *host* connected with a GPU called the *device*.

The highly parallel structures of GPUs are based on SIMD execution and support parallel processing of large blocks of data. A GPU has multiple *multithreaded SIMD* processors. The current generation of GPUs, e.g., Fermi of NVIDIA, have 7–15 multithreaded SIMD processors. Compared with vector processors, each multithreaded SIMD processor has several wide and shallow SISD *lanes*. For example, an NVIDIA GPU has 32 768 registers divided among the 16 physical SIMD lanes; each lane has 2 048 registers.

A typical execution includes several steps: (i) CPU copies the input data from the main memory to the GPU memory; (ii) CPU instructs the GPU to start processing using the executable in the GPU memory; (iii) GPU uses multiple cores to execute the parallel code; and (iv) when done, the GPU copies the result back to the main memory.

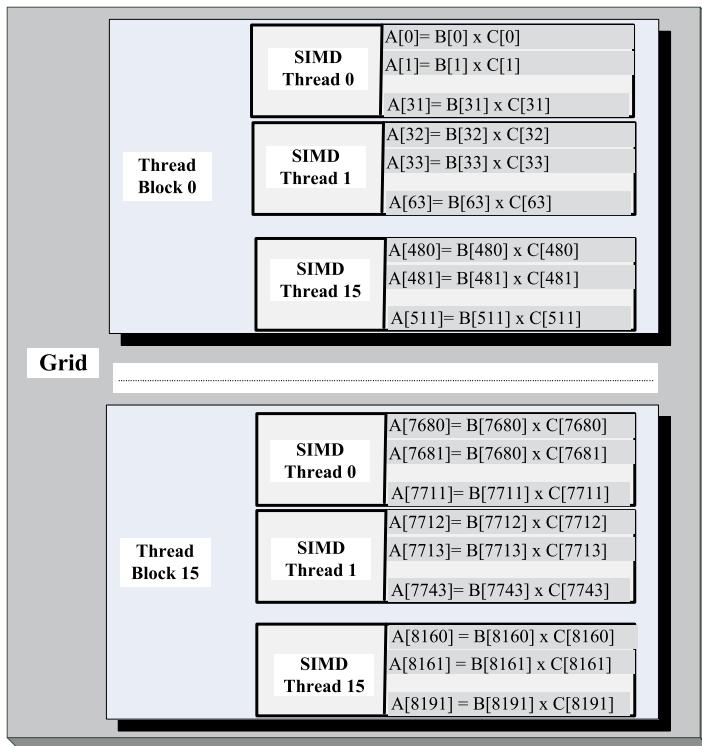
The GPU programming model is Single-Instruction-Multiple-Thread (SIMT). GPUs are often programmed in CUDA, a C-like programming language developed by NVIDIA, the pioneer of GPU-accelerated computing. All forms of GPU parallelism are unified as CUDA threads. In the SIMT model a *thread* is associated with each data element. Thousands of CUDA threads can run concurrently. Threads are organized in *blocks*, groups of 512 vector elements, and multiple blocks form a grid. Fig. 3.5 shows the grid representing a vector *A* with 8 192 components; there are 16 blocks, each block has 16 SIMD threads; each thread operates on 32 elements of the vector.

A two-level scheduling mechanisms assigns threads to the multiple *lanes* of a multithreaded SIMD processor. A *thread block scheduler* assigns thread blocks to multithreaded SIMD processors, and then a *tread scheduler* assigns threads to SIMD *lanes*. Fig. 3.6 illustrates scheduling for the grid in Fig. 3.5. The thread blocks must be able to run independently. NVIDIA GPU memory has the following organization:

- a. Each SIMD lane has an off-chip private memory for stack frame, spilling registers, private variables, and GPU data in L1 and L2 caches.
- b. Each multithread SIMD processor has on-chip local memory shared by all its lanes, thus, by the threads within the block scheduled on the processor.
- c. GPU memory. The host can read from and write to this off-chip memory.

In 2018, Nvidia and AMD controlled nearly 100% of the GPU market. Nvidia launched RTX 20 series in 2018 with ray-tracing cores to improve performance on lighting effects. One of the most powerful GPUs produced by NVIDIA is GEFORCE RTX 3090 with a 1.7 GHz clock, 24 GB of memory with a 384-bit memory interface width and 10 496 cores. In 2019 AMD launched GCN (Graphics Core Next) ISA and the RX 6000 series.

It should not be surprising that cloud service providers now offer GPU instances. For example, AWS EC2 P4 instances feature up to eight NVIDIA A100 Tensor Core GPUs with 600 GB/s peer to peer GPU, and communication with NVIDIA NVSwitch and 400-Gbps instance networking with support for Elastic Fabric Adapter (EFA) and NVIDIA GPUDirect RDMA (remote direct memory access).

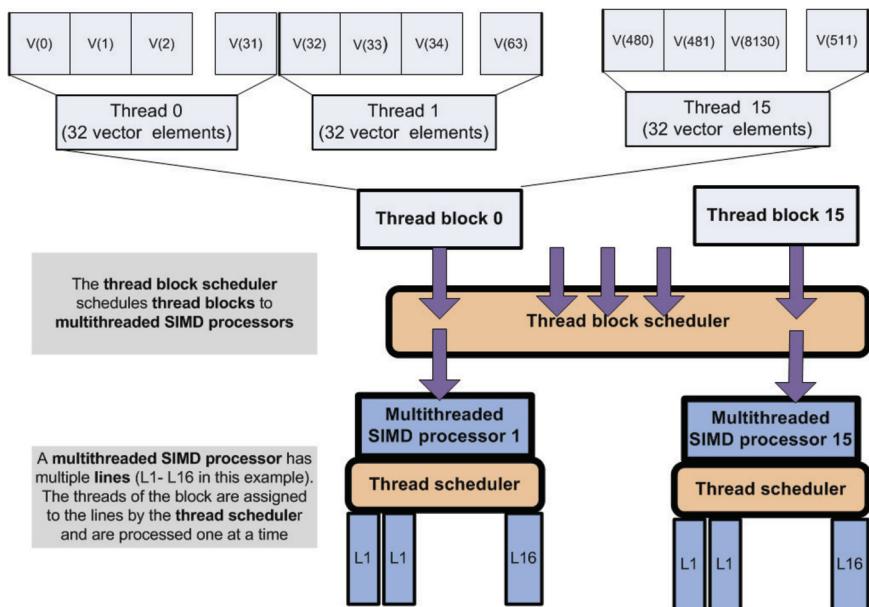
**FIGURE 3.5**

Grid, blocks, and threads. Grids with 8 192 components for vectors A , B , C . There are 16 blocks with 512 vector components each. Each block has 6 threads, and each thread operates on 32 vector elements.

3.6 Tensor processing units

For more than half a century, Moore's Law enabled major architectural enhancements including deep pipelines, branch prediction, out-of-order instruction execution, speculative prefetching, multithreading, deep memory hierarchies, and wide SIMD units. In the first decade of the new millennium, it became apparent that the end of an era was in sight and the age of Domain Specific Architectures (DSA) began.

With the end of Moore's Law in sight, Dennard's Law has already limited processor clock rate. Moore's Law, formulated in 1965 by one of Intel's founders, Gordon Moore, states that the number of transistors per chip at constant cost doubles every two years. This empirical, rather than physical law, has often been confused as stating that processor performance doubles every two years because transistor size and speed are related, the smaller the transistor the higher its switching speed. In 1974, Dennard observed that voltage and current should be proportional to the linear dimensions of a transistor. Power is proportional to transistor area and, as transistors get smaller, power density increases. Since around 2006, Dennard's Law has limited processor clock rate to 4 GHz or less.

**FIGURE 3.6**

The execution for the grid in Fig. 3.5. The *thread block scheduler* assigns thread blocks to multithreaded SIMD processors. A *thread scheduler* running on each multithreaded SIMD processor assigns threads to the SIMD *lanes* of the processors.

The DSA objective is to extract performance from software oblivious to system architecture. In 2015, Google deployed the first custom Application Specific Integrated Circuit (ASIC) for its cloud computing infrastructure, a Tensor Processing Unit. TPUs were designed for Machine Learning (ML) applications, in particular, for inference stages of Deep Neural Networks (DNNs). TPUs were followed by Microsoft Catapult, a Field-Programable Gate Array (FPGS), Nervana's ASIC called CREST, and Google's own Pixel Visual Core.

TPU design aimed to achieve at least one-order-of-magnitude performance improvement versus a GPU. TPUs are single-threaded processors, based on a deterministic execution model, a good match to the 99th-percentile response-time requirement of a typical DNN inference application [262]. TPUs are designed to be plugged into existing servers as co-processors. The host CPU sends TPU instructions over the PCIe bus into an instruction buffer. The internal blocks are typically connected together by 256-byte-wide (2048-bits) paths. TPU organization is shown in Fig. 3.7; TPA ISA includes the following instructions:

Read_Host_Memory—reads data from the CPU host memory into Unified Buffer.

Write_Host_Memory—writes data from the Unified Buffer into CPU host memory.

Read_Weights—reads weights from Weight Memory into Weight FIFO as input to Matrix Unit.

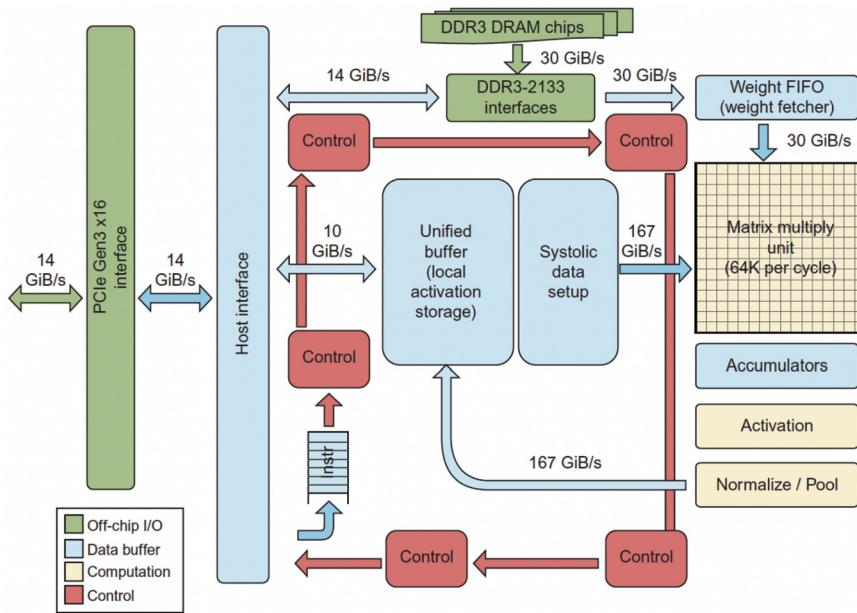


FIGURE 3.7

TPU organization and its connections to the CPU via a host interface and a 14 GiB/sec PCIe gen3 x 16 bus and with the DDR3 Dram memory via a 30 GiB/sec memory bus. The Matrix Multiply Unit gives the TPU its enormous computational power. The activation unit performs nonlinear functions, such as sigmoid on data supplied by accumulators, and transmits them to the unified buffer [262].

MatrixMultiply/Convolve—causes the Matrix Multiply Unit to perform a matrix-matrix multiply, a vector-matrix multiply, an element-wise matrix multiply, an element-wise vector multiply, or a convolution from the Unified Buffer into the Accumulators.

Activate—performs nonlinear function of the artificial neuron, with options for ReLU, Sigmoid, tanh, and so on. Its inputs are the Accumulators, and its output is the Unified Buffer.

Application code running on the TPU is typically written in TensorFlow and is compiled into an API that can run on GPUs or TPUs. A third-generation TPU, the Edge TPU, much smaller and using far less power compared to Cloud TPUs, was released in 2018. The Edge TPU is capable of four trillion operations per second while using 2W and runs a version of Tensor Flow called TensorFlow Lite.

A 2017 paper [262] reports on performance analysis of a datacenter TPU with a 65 536 8-bit MAC (multiply-accumulate unit) matrix multiply unit with a peak throughput of 92 TeraOps/second (TOPS) and a 28 MiB software-managed on-chip memory. It concludes that “TPU leverages the order-of-magnitude reduction in energy and area of 8-bit integer systolic matrix multipliers over 32-bit floating-point data paths of a K80 GPU to pack 25 times as many MACs (65 536 8-bit vs. 2 496 32-bit) and 3.5 times the on-chip memory (28 MiB vs. 8 MiB) while using less than half the power of the K80 in a relatively small die.”

3.7 Systems on a chip

A system on a chip (SoC) is an integrated circuit hosting on the same substrate/chip multiples cores, including a mix of traditional CPUs, alongside GPUs, TPUs, other types of functional units, e.g., DSPs (Digital Signal Processors), memory, input/output ports, secondary storage, and sometimes WiFi and cellular modems. SoC organization improves performance and reduces power consumption, as well as semiconductor die area, compared with motherboard-based architecture with equivalent functionality but discrete components.

All this is possible due to cutting-edge technologies such as the five-nanometer process, which can populate the chip with billions of transistors. The trend to processor customization reflects the reality that Moore's Law is slowing down and power dissipation limits the clock rates [408]. The side effect of this organization is increased replacement cost. Once any component of an SoC becomes inoperative, the entire SoC must be replaced.

SoC processor cores typically use RISC instruction-set architectures and in particular ARM due to its advantages discussed in Section 3.3. SoC memory consists of a memory hierarchy and cache hierarchy including read-only memory (ROM), random-access memory (RAM), Electrically Erasable Programmable ROM (EEPROM), and flash memory. SRAM is used to implement processor registers and cores' L1 caches, whereas DRAM is used for low level(s) of memory hierarchy discussed in Section 3.2. SoC units communicate with one another using data bus architectures, often based on ARM's royalty-free Advanced Microcontroller Bus Architecture (AMBA) standard and, more recently, sparse intercommunication networks known as networks-on-a-chip (NoC).

SoCs are optimized to maximize power efficiency quantified by performance per watt. Applications, such as edge computing, distributed processing, and ambient intelligence, require high-performance computing SoCs with high-power efficiency, while power is limited for SoCs used for mobile device such as smartphones, watches, tablets, laptops, or GPS. The five-fold discrepancy between battery density improvements and Moore's law implies there is a widening gap between compute requirements and energy availability. Increasing power density leads to thermal challenges exceeding the inherent passive cooling capability of a mobile device. All SoCs minimize energy dissipation because the heat from high energy consumption of one unit can damage other circuit components. Mobile devices operate under strict thermal and energy budgets.

Apple A14 Bionic is a 64-bit ARMv8.5a SoC with 64-bit six-core CPU with two high-performance cores and four energy-efficient GPU cores for faster graphics. It is manufactured by TSMC (Taiwan Semiconductor Manufacturing) using their first-generation 5-nm fabrication process. The transistor density is 134 million transistors per m², for a total of 11.8 billion transistors. This SoC includes a custom silicon block called the Neural Engine, hardware tailor-made to perform operations used in machine learning and AI. A14 has also a Neural Engine with 16 cores delivering 11 trillion operations per second.

Another processor, M1, inspired by A14, is the first ARM-based SoC designed as CPU for Apple computers. M1 has four high-performance and four energy-efficient cores, The high-performance cores have 192 KB of L1 instruction cache and 128 KB of L1 data cache and share a 12 MB L2 cache. The energy-efficient cores have a 128 KB L1 instruction cache, 64 KB L1 data cache, and shared 4 MB L2 cache and consume one tenth of the energy consumed by high-performance cores. M1 has a GPU with eight cores; each one contains eight *executions units* (Apple terminology for the multithreaded SIMD processors in Fig. 3.6), which in turn contain eight ALUs per execution unit. M1's GPU can execute nearly 25 000 threads simultaneously and has a maximum performance of 2.6 Tflops.

The computing power of today SoCs, even of those for mobile devices, is on par with, or higher than, that of supercomputers of two decades ago and can carry out complex computational tasks locally with lower power consumption and without the need to transfer data to the cloud. Computer clouds are in a symbiotic relationship not only with billions of mobile devices used to access data stored on clouds but also with systems outside the cloud infrastructure and carrying out computational tasks together with clouds as edge computing, also discussed in Section 12.12. Edge computing is a distributed computing framework that brings enterprise applications closer to data sources, such as IoT devices or local edge servers. This proximity to data at its source can deliver strong benefits: faster insights, improved response times, and better bandwidth availability.

3.8 Data, thread-level, and task-level parallelism

As demonstrated by nature, the ability to work as a group and carry out many tasks in parallel represents a very efficient way to reach a common goal. Thus, we should not be surprised that the thought that individual computer systems should work in parallel for solving complex applications was formulated early on, at the dawn of the computer age. Parallel processing allows us to solve large problems by splitting them into smaller ones and solving them concurrently. Parallel processing was considered for many years the holy grail for solving data-intensive problems encountered in many areas of science, engineering, and enterprise computing and required major advances in several areas, including algorithms, programming languages and environments, performance monitoring, computer architecture, interconnection networks, and, last but not least, solid-state technologies. Often discovering parallelism is quite challenging, and the development of parallel algorithms requires a considerable effort.

Fine-grained versus coarse-grained parallelism. We distinguish *fine-grained* from *coarse-grained* parallelism. In the former case only relatively small blocks of the code can be executed in parallel without the need to communicate or synchronize with other threads or processes, whereas in the latter case large blocks of code can be executed concurrently.

Numerical computations involving linear algebra operations exhibit fine-grained parallelism. For example, many numerical analysis problems, such as solving large systems of linear equations, or solving systems of Partial Differential Equations (PDEs) require algorithms based on domain decomposition methods. The speedup of applications displaying fine-grained parallelism is considerably lower than those of coarse-grained applications. Indeed, even systems with a fast interconnect processor speed are orders of magnitude higher than communication speed.

Concurrent processes or threads communicate using message-passing or shared-memory. Shared-memory is the defining attribute of distributed shared-memory multiprocessor systems. Shared-memory is also used by multicore processors where each core has private L1 instruction and data caches, as well as an L2 cache, and all cores share the L3 cache. Shared-memory is not scalable, thus, it is seldom used in supercomputers and large clusters. Message passing is used exclusively in large-scale distributed systems, and our discussion is restricted to this communication paradigm. Debugging a message-passing application is considerably easier than debugging a shared-memory one.

Shared-memory is extensively used by system software. The system stack is an example of shared-memory used to save the state of a process or thread at the time of a context switch. The kernel of an operating system uses control structures, such as processor and core tables for multiprocessor and multicore system management, process and thread tables for process and thread management, page tables

for virtual memory management, among others. Multiple application threads running on a multicore processor often communicate via the shared-memory of the system.

Data-level parallelism. This is an extreme form of coarse-grained parallelism. It is based on partitioning data into large chunks/blocks/segments and running concurrently either multiple programs or copies of the same program, each on a different data block. In the latter case the paradigm is called *Same Program Multiple Data* (SPMD). There are the so-called *embarrassingly parallel* problems where little or no effort is needed to extract parallelism and to run a number of concurrent tasks with little or no communication among them.

Assume that we wish to search for the occurrence of an object in a set of n images or of a string of characters in n records. Such a search can be conducted in parallel. In all these instances the time required to carry out the computational task using N servers is reduced by a factor of N , and the speedup is almost linear to the number of servers used. This type of data-parallel applications is at the heart of enterprise computing on computer clouds and will be discussed in depth in Chapter 11. The MapReduce programming model will be presented in Section 11.5, followed by the discussion of Hadoop and Yarn in Section 11.7.

Decomposition of a large problem into a set of smaller problems that can be solved concurrently is sometimes trivial and can be implemented in hardware, a topic discussed in Section 3.1. For example, assume that we wish to manipulate the image of a three-dimensional object represented as a 3D lattice of $n \times n \times n$ points. To rotate the image, we apply the same transformation to each one of the n^3 points. Such a transformation can be done by a *geometric engine*, a processor designed to carry out the transformation of a subset of n^3 points concurrently. Graphics Processing Units (GPUs) discussed in Section 3.5 initially designed as graphics engines are now widely used by data-intensive applications.

Thread-level and task-level parallelism. The term thread-level parallelism is overloaded. In the computer architecture literature it is used for data-parallel execution using a GPU. In this case a *thread* is a subset of vector elements processed by one of the lanes of a multithreaded processor, see Section 3.5. *Hyper-threading* is used to describe multiple execution threads possibly running concurrently, but on a single core, see Section 3.1. Threads are also used in a multicore processor to run multiple processes concurrently. Database applications are memory-intensive and I/O-intensive; and multiple *threads* are used to hide the latency of memory and I/O access.

In the context of scheduling, a job consists of multiple tasks scheduled either independently or co-scheduled when they need to communicate with one another. Often fine-grained execution units are given control of resources for a relatively short time to guarantee a low-latency response time.

Cloud computing is an appealing environment of running applications attempting to identify optimal model parameters that best fit experimental data. Such applications involve computationally intensive tasks. Multiple instances running concurrently test the fitness of different sets of parameters, and the results are then compared to determine the optimal set of model parameters.

There are also numerical simulations of complex systems that require an optimal design of a physical system. Multiple design alternatives are compared, and the optimal one is selected according to several optimization criteria. Consider for example the design of a circuit using Field Programmable Gate Arrays (FPGAs). An FPGA is an integrated circuit designed to be configured by the customer using a hardware description language (HDL), similar to that used for an application-specific integrated circuit (ASIC).

As multiple choices for the placement of components and for interconnecting them exist, the designer could run concurrently N versions of the design choices and choose the one with the best

performance, e.g., minimum power consumption. Alternative optimization objectives could be to reduce cross-talk among the wires or to minimize the overall noise. Because each alternative configuration requires hours, or maybe days, of computing, running them concurrently reduces the design time considerably.

3.9 Speedup, Amdahl's law, and scaled speedup

Parallel hardware and software systems allow us to solve problems demanding more resources than those provided by a single system and, at the same time, to reduce the time required to obtain a solution. There are applications, the so-called *embarrassingly parallel applications*, when the *parallel execution time*, the time when running on N cores is $1/N$ of the *sequential execution time*, the time when running on a single core. There are other applications when the parallel execution time is only slightly smaller than the sequential execution.

This begs questions such as: Is there a way to quantify the amount of parallelism that can be extracted from an application?; what is the largest number of cores that can be used effectively to run an application?; how does one measure the parallelization effectiveness?; and does the amount of data processed by the application matter? These questions were addressed in 1960s by Gene Amdahl and more recently by Gustafson and are discussed in this and the next section. The effectiveness of parallelization is measured by the speedup. In the general case the *speedup* of the parallel computation is defined as

$$S(N) = \frac{T(1)}{T(N)}, \quad (3.1)$$

with $T(1)$ the execution time of the sequential computation and $T(N)$ the execution time when N parallel computations are carried out.

Amdahl's Law. Gene Myron Amdahl⁴ was a theoretical physicist turned computer architect best known for Amdahl's Law. In a 1967 seminal paper [21] Amdahl argued that the fraction of a computation that is not parallelizable is significant enough to favor single-processor systems. He reasoned that large-scale computing capabilities can be achieved by enhancing the performance of single processors, rather than building multiprocessor systems.

Though this thesis was disproved, Amdahl's Law is a fundamental result used to predict the theoretical maximum speedup for a program using multiple processors. This law states that *the portion of the computation that cannot be parallelized determines the overall speedup*. If α is the fraction of running time a sequential program spends on non-parallelizable segments of the computation, then the maximum speedup achievable S is

$$S = \frac{1}{\alpha}. \quad (3.2)$$

⁴ Gene Amdahl contributed significantly to the development of several IBM systems, including System/360, and in the 1970s, started his own company, Amdahl Corporation, to produce high-performance systems.

To prove this result, call σ the sequential time and π the parallel time and start from the definitions of $T(1)$, $T(N)$, and $\alpha = \sigma/(\sigma + \pi)$:

$$T(1) = \sigma + \pi, \quad T(N) = \sigma + \frac{\pi}{N}, \quad \text{and} \quad \alpha = \frac{\sigma}{\sigma + \pi}. \quad (3.3)$$

Then

$$S = \frac{T(1)}{T(N)} = \frac{\sigma + \pi}{\sigma + \pi/N} = \frac{1 + \pi/\sigma}{1 + (\pi/\sigma) \times (1/N)}. \quad (3.4)$$

But

$$\pi/\sigma = \frac{1 - \alpha}{\alpha}. \quad (3.5)$$

Thus, for large N ,

$$S = \frac{1 + (1 - \alpha)/\alpha}{1 + (1 - \alpha)/(N\alpha)} = \frac{1}{\alpha + (1 - \alpha)/N} \approx \frac{1}{\alpha}. \quad (3.6)$$

An alternative formulation of Amdahl's Law is that, if a fraction f of a computation is enhanced by a speedup S , then the overall speedup is

$$S_{overall}(f, S) = \frac{f}{(1 - f) + \frac{f}{S}} \quad \text{or} \quad S_{overall}(f, S) = \frac{1}{\frac{1}{f} + \frac{1}{S} - 1}. \quad (3.7)$$

Scaled speedup. Amdahl's Law applies to a *fixed problem size*; in this case, the amount of work assigned to each one of the parallel processes decreases when the number of processes increases, and this affects the efficiency of the parallel execution. When the problem size is allowed to change, Gustafson's Law gives the *scaled speedup* with N parallel processes as

$$S(N) = N - \alpha(N - 1). \quad (3.8)$$

As before, we call σ the sequential time; now π is the *fixed parallel time per process*; α is given by Eq. (3.3). The sequential execution time, $T(1)$, and the parallel execution time with N parallel processes, $T(N)$, are

$$T(1) = \sigma + N\pi \quad \text{and} \quad T(N) = \sigma + \pi. \quad (3.9)$$

Then, the scaled speedup is

$$S(N) = \frac{T(1)}{T(N)} = \frac{\sigma + N\pi}{\sigma + \pi} = \frac{\sigma}{\sigma + \pi} + \frac{N\pi}{\sigma + \pi} = \alpha + N(1 - \alpha) = N - \alpha(N - 1). \quad (3.10)$$

Amdahl's Law expressed by Eq. (3.2) and the *scaled speedup* given by Eq. (3.8) assume that all processes are assigned the same amount of work. The scaled speedup assumes that the amount of work assigned to each process is the same regardless of the problem size. Then, to maintain the same execution time, the number of parallel processes must increase with the problem size. The scaled speedup captures the essence of efficiency, namely, that the limitations of the sequential part of a code can be balanced by increasing the problem size.

3.10 Multicore processor speedup

We now live in the age of multicore processors brought about by the limitations imposed on solid-state devices by the laws of physics. Increased power dissipation due to faster clock rates makes the heat removal more challenging, and this implies that in the future, we could expect only a modest increase of the clock rate. Multicore processors use the billions of transistors on a chip to deliver significantly higher computing power and process more data every second. Yet, the ability to get data in and out of the chip is limited by the number of pins. Increasing the number of cores on a chip faces its own physical limitations.

There are alternative designs of multicore processors, and the next question is to investigate chip configurations most useful for applications exhibiting a limited parallelism. The cores can be identical or different from one another, or there could be a few powerful cores or a larger number of less powerful cores. Theoretically, the cores could be configured automatically or be immutable.

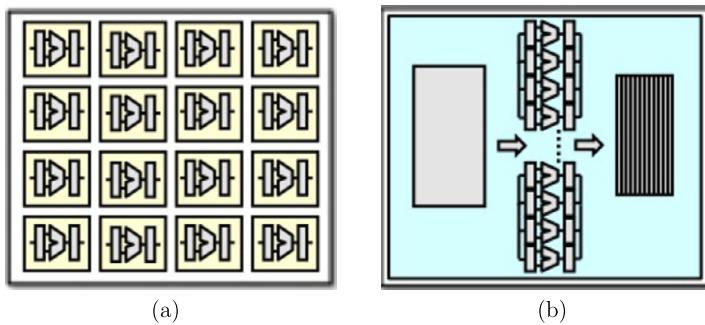
More cores will lead to a large speedup of highly parallel applications; a powerful core will favor highly sequential applications. If feasible, a chameleonic system will adapt to the level of parallelism, though changing the core configuration will incur overhead and challenge application developers. Even considering re-configuration overhead, the speedup of automatic core configuration will be superior to either symmetric or asymmetric core design.

The design space of multicore processors should be driven by cost–performance considerations. A design will be cost-effective if the speedup achieved will exceed the *cost up* defined as the multicore processor cost divided by the single-core processor cost. The cost of a multicore processor depends on the number of cores and the complexity, ergo, the power of individual cores.

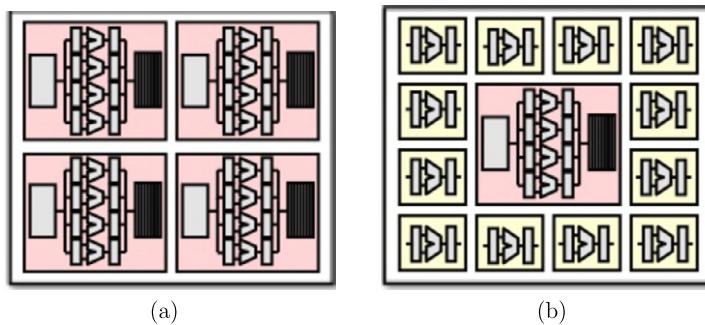
The *Basic Core Equivalent* (BCE) concept was introduced to quantify the resources of individual cores. A *symmetric core* processor may have n BCEs with r resources each. Alternatively, $n \times r$ resources may be distributed unevenly in an *asymmetric core* processor.

A quantitative analysis of design choices based on an extension of Amdahl’s Law to multicore processors is presented in [238]. We expect this analysis to confirm the obvious, that is: *The larger the parallelizable fraction f of an application, the larger the speedup*. This analysis is based on a number of simplifying assumptions:

- i. A number of factors such as the chip area, the power dissipation, or combinations of these two with other factors limit the number of cores to n BCE. These limitations consider only on-chip resources. Off-chip resources such as shared caches, memory controllers, or interconnection networks are assumed to be approximately identical for the alternative designs.
- ii. The performance of a single BCE core is the unity. When the chip is limited to n BCEs, all cores are identical and the performance of each core is r BCE, then the total number of cores on a chip is $\lceil n/r \rceil$. Fig. 3.8 shows a symmetric 16-BCE chip with two configurations: one with sixteen 1-BCE cores and the other with one 16-BCE core.
- iii. The sequential performance of r BCEs is denoted as $perf(r)$. When $perf(r) > r$, the speedup of both sequential and parallel execution increases; therefore, the designers should increase as much as feasible the core resources and implicitly the individual core performance. On the other hand, when $perf(r) < r$, increasing individual core performance increases the performance for sequential execution but lowers that of the parallel execution. Therefore, the following analysis is focused on the case when $perf(r) < r$. A good model is $perf(r) = \sqrt{r}$ when the performance doubles, triples, and quadruples for 4, 9, 16 cores, respectively.

**FIGURE 3.8**

16-BCE chip. Symmetric core processor with two different configurations: (a) sixteen 1-BCE cores; (b) one 16-BCE core.

**FIGURE 3.9**

16-BCE chip. (a) Symmetric core processor with four 4-BCE cores; (b) Asymmetric core processor with 1 4-BCE core and 12 1-BCE cores.

The first case analyzed assumes a symmetric multicore chip. There are n/r cores on the chip; for example, when $n = 32$, the chip could have 16 cores of 2 BCE each, 8 cores of 4 BCE each, and so on. Call f the parallelizable fraction of a computation. One core will run the $(1 - f)$ sequential component of the computation, and n/r cores will run the parallel component of the computation, f , with a performance $\text{perf}(r)$. According to Eq. (3.7) the speedup will be

$$\text{Speedup}_{\text{symcore}}(f, n, r) = \frac{1}{\frac{1-f}{\text{perf}(r)} + \frac{f \cdot r}{n \cdot \text{perf}(r)}}. \quad (3.11)$$

In an asymmetric multicore processor, more powerful cores will coexist with less powerful ones. Fig. 3.9 illustrates the differences between symmetric and asymmetric cores; the asymmetric core processor has a 4-BCE core and 12 1-BCE cores. The sequential performance will benefit from the more powerful core running four times faster, while the parallel performance is $\text{perf}(r)$ from the 4-BCE

core and *one* each from the remaining $(n - r)$, in our case, 12 1-BCE cores. The speedup with one powerful core and $(n - r)$ 1-BCE cores is then

$$\text{Speedup}_{\text{asymcore}}(f, n, r) = \frac{1}{\frac{1-f}{\text{perf}(r)} + \frac{f}{\text{perf}(r)+n-r}}. \quad (3.12)$$

A dynamic multicore chip could configure its n BCE depending on the fraction f of a particular application. If the sequential component of the application, $(1 - f)$ is large, then configure the chip as one r -BCE core, while in parallel mode use all base cores in parallel. In this case

$$\text{Speedup}_{\text{dyncore}}(f, n, r) = \frac{1}{\frac{1-f}{\text{perf}(r)} + \frac{f}{n}}. \quad (3.13)$$

What non-obvious conclusions can be drawn from this analysis? First, the speedup of asymmetric multicore processors is always larger and, in some cases, could be significantly larger than the speedup of symmetric core processors. For example, the largest speedup when $f = 0.975$ and $n = 1024$ is achieved for a configuration with one 345-BCE core and 679 1-BCE cores. Second, increasing the power of individual cores is beneficial even for symmetric core processors. For example, when $f = 0.975$ and $n = 256$, the maximum speedup occurs for seven 1-BCE cores.

Not to be forgotten should be the fact that task scheduling on asymmetric and dynamic multicore processors will be fairly difficult. There are also other factors affecting the performance ignored by the simple model discussed in [238].

3.11 From supercomputers to distributed systems

Computational science and engineering applications require high-performance computing systems able to exploit fine-grained parallelism. The stringent requirements of such applications have motivated many of the architectural advancements we see in modern processors. For example, the supercomputers of the mid 1960s, CDC 6400, 6500, and 7600 manufactured by Control Data Corporation (CDC), pioneered pipelining and multiple functional units for addition, multiplication, and other arithmetic and logic operations. In 1969, IBM unsuccessfully entered the scientific computing market with IBM 360 model 91, which had novel features including an implementation of Tomasulo's algorithm for out-of-order instruction execution.

Modern supercomputers derive their power from architecture and parallelism rather than faster processors running at higher clock rates. The supercomputers of today consist of a very large number of processors and cores communicating through very fast and expensive *custom interconnects*. In mid 2012 the most powerful supercomputer was an IBM Sequoia-BlueGene/Q Linux-based system powered by Power BQC 16-core processors running at 1.6 GHz. The system, installed at Lawrence Livermore National Laboratory, which had a total of 1 572 864 cores and 1 572.864 TB of memory, achieved a sustainable speed of 16.32 PFlops, and consumed 7.89 MW of power. Later in 2012 a Cray XK7 system, Titan, installed at the Oak Ridge National Laboratory (ORNL) with 560 640 processors, including 261 632 Nvidia K20x accelerator cores, achieved a speed of 17.59 PFlops on the Linpack benchmark. In 2016, the most powerful supercomputer was Sunway TaihuLight at National Supercomputer Center

in Wixi, China, with 10 649 600 cores with a peak bandwidth of 125.436 PFlops. The system needed 15.371 MW of power. Its Linpack performance is 93.0146 PFlops, and it has 1 310.720 TB of memory.

As of June 2020, the Fugaku supercomputer at the RIKEN Center for Computational Science, Japan, powered by A64FX 2.2 GHz, an ARM architecture microprocessor designed by Fujitsu topped the list. It has 7 299 072 cores, a Linpack performance of 415.530 TFlop/sec, and needs 28.3345 MW. It is followed by the Summit at DOE/SC/Oak Ridge National Laboratory in the US using IBM Power System AC922, IBM POWER9 22C 3.07 GHz, with NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband. It has 2 414 592 cores, a Linpack performance of 148.600 TFlops/sec and needs 10 MW. Several of the most powerful systems listed in [478] are powered by Nvidia 2050 GPU. Some of the top 10 supercomputers use the Infiniband interconnect discussed in Chapter 6.

The next natural step to increase the computing power was possible due to advances in communication networks when low-latency and high-bandwidth Wide Area Networks (WANs) allowed individual systems, many of them multiprocessors, to be geographically separated. Large-scale distributed systems were first used for scientific and engineering applications and took advantage of the advancements in system software, programming models, tools, and algorithms developed for parallel processing.

A distributed system is a collection of autonomous and heterogeneous systems able to communicate effectively with each other. The question is: How could such a collection be organized to cooperate and compute efficiently? In spite of intensive research efforts spanning many years, an optimal organization of a large-scale system has eluded us.

The inability to conceive an optimal general-purpose system reflects the realization that a system cannot be looked upon in isolation, rather it should be analyzed in the context of the environment it is expected to operate in. The more complex and diverse this environment, the less likely it is to display an asymptotically optimal performance for all, or virtually all, applications. The organization and the system management may be optimal for a class of applications, yet have a sub-optimal performance for others.

Distributed systems have existed for several decades. For example, distributed file systems and network file systems have been used for user convenience and for improving reliability and functionality of file systems for many years, see Section 7.4. Modern operating systems allow a user to *mount* a remote file system and access it the same way a local file system is accessed, but with a performance penalty due to larger communication costs. *Remote Procedure Calls* (RPCs) allow a procedure on one system to invoke a procedure running in another address space, possibly on a remote system.

A *distributed system* is a collection of computers connected through a network and a distribution software called *middleware*, which enables computers to coordinate their activities and to share the resources of the system; the users perceive the system as a single, integrated computing facility. The middleware should support this set of desirable properties of the distributed system:

1. *Access transparency*; local and remote information objects should be accessed using identical operations.
2. *Location transparency*; information objects should be accessed without knowledge of their location.
3. *Concurrency transparency*; processes running concurrently should shared information objects without interference among them.
4. *Replication transparency*; multiple instances of information objects should be used to increase reliability without the knowledge of users or applications.
5. *Failure transparency*; the faults should be concealed.

6. *Migration transparency*; the information objects in the system should be moved without affecting the operation performed on them.
7. *Performance transparency*; the system should be reconfigured based on the load and quality of service requirements.
8. *Scaling transparency*; the system and the applications should scale without a change in the system structure and without affecting the applications.

A distributed system has several characteristics: Its components are autonomous, meaning scheduling and other resource management and security policies are implemented by each system. There are multiple points of control and multiple points of failure in a distributed system, and some sources may not be accessible at all times. Distributed systems can be scaled-up by adding more resources and can be designed to maintain availability even at low levels of hardware/software/network reliability. The availability is a system metric aiming to ensure an agreed-upon level of operational performance for an extended period of time.

3.12 Modularity. Soft modularity versus enforced modularity

Modularity is a basic concept in the design of man-made systems; a system is made out of components, or modules, with well-defined functions. Modularity supports the separation of concerns, encourages specialization, improves maintainability, reduces costs, and decreases the development time of a system.

Modularity has been used extensively since the Industrial Revolution to build every imaginable product, from weaving looms to steam engines, from watches to automobiles, from electronic devices to airplanes. Individual modules are often made of sub-assemblies. Modularity can reduce cost for the manufacturer and for the consumers. The same module may be used by a manufacturer in multiple products; to repair a defective product, a consumer only replaces the module causing the malfunction rather than the entire product. Modularity encourages specialization as individual modules can be developed by experts with deep understanding of a particular field. It also supports innovation since it allows a module to be replaced with a better one, without affecting the rest of the system.

It is no surprise that the hardware, as well as the software systems are composed of modules interacting with one another through well-defined interfaces. Software development for distributed systems is more challenging than for sequential systems, and these challenges are amplified by the scale of the system and the diversity of applications.

Modularity, layering, and hierarchy are some of the means to cope with the complexity of a distributed application software. Software modularity, the separation of a function into independent, interchangeable modules, requires well-defined interfaces specifying the elements provided and supplied to a module [389]. A modular software design is driven by several principles outlined in [139], namely:

- a. *Information hiding*; the user of a module does not need to know anything about the internal mechanism of the module to make effective use of it.
- b. *Invariant behavior*; the functional behavior of a module must be independent of the site or context from which it is invoked.

- c. *Data generality*; the interface to a module must be capable of passing any data object an application may require.
- d. *Secure arguments*; the interface to a module must not allow side-effects on arguments supplied to the interface.
- e. *Recursive construction*; a program constructed from modules must be usable as a component in building larger programs or modules.
- f. *System resource management*; resource management for program modules must be performed by the computer system and not by individual program modules.

Some of these principles are implicitly supported by the enforced modularity. A system should prevent modules to make private resource-allocation decisions and should support a global address space. The modularity concept is dissected, and its applications are reviewed in the next section.

The progress made in system design is notable not in the least due to a number of principles guiding the design of parallel and distributed systems. One of these principles is specialization; this means that a number of functions are identified and an adequate number of system components are configured to provide these functions. For example, data storage is an intrinsic function, and storage servers are a ubiquitous presence in most systems. This brings us to the modularity concept.

Modularity allows us to create a complex software system from a set of components built and tested independently. A requirement for modularity is to clearly define the interfaces between modules and enable the modules to work together. The steps involved in the transfer of the flow of control between the caller and the callee are: (i) the caller saves on the stack its state including registers, arguments, and return address; (ii) the callee loads the arguments from the stack, carries out the calculations, and then transfers control back to the caller; and (iii) the caller adjusts the stack, restores its registers, and continues its processing.

Soft modularity. We distinguish *soft modularity* from *enforced modularity*. The former implies dividing a program into modules that call each other and communicate using shared-memory or follow the procedure call convention.

Soft modularity hides the details of the implementation of a module and has many advantages: Once the interfaces of the modules are defined, the modules can be developed independently; a module can be replaced with a more elaborate, or with a more efficient one, as long as its interfaces with the other modules are not changed. The modules can be written using different programming languages and can be tested independently.

Soft modularity presents a number of challenges. It increases the difficulty of debugging, for example, a call to a module with an infinite loop it will never return. There could be naming conflicts and wrong context specifications. The caller and the callee are in the same address space and may misuse the stack, e.g., the callee may use registers that the caller has not saved on the stack, and so on.

Strongly-typed languages may enforce soft modularity by ensuring type safety at compile time or at run time, it may reject operations or function class that disregard the data types, or it may not allow class instances to have their class altered. Soft modularity may be affected by errors in the run-time system, errors in the compiler, or by the fact that different modules are written in different programming languages.

Enforced modularity. The ubiquitous client–server paradigm is based on enforced modularity; this means that the modules are forced to *interact only by sending and receiving messages*. This paradigm

leads to a more robust design: The clients and the servers are independent modules and may fail separately.

Moreover, the servers are stateless, i.e., they do not have to maintain state information. A server may fail and then come back up without the clients being affected, or even noticing the failure of the server. The system is more robust because it does not allow errors to propagate. Enforced modularity makes an attack less likely because it is difficult for an intruder to guess the format of the messages or the sequence numbers of segments, when messages are transported by TCP.

Last, but not least, resources can be managed more efficiently. For example, a server typically consists of an ensemble of systems, a *front-end* system which dispatches the requests to multiple *back-end* systems which process the requests. Such an architecture exploits the elasticity of a computer cloud infrastructure: The larger the request rate, the larger is the number of back-end systems activated.

The client–server paradigm. This paradigm allows systems with different processor architecture, e.g., 32-bit or 64-bit, with different operating systems, e.g., multiple versions of operating systems, such as Linux, Mac OS, or Microsoft Windows, libraries and other system software, to cooperate. The client–server paradigm increases flexibility and choice; the same service could be available from multiple providers, a server may use services provided by other servers, a client may use multiple servers, and so on.

Heterogeneity of systems based on the client–server paradigm is less of a blessing, e.g., the problems it creates outweigh its appeal. Heterogeneity adds to the complexity of the interactions between a client and a server as it may require conversion from one data format to another, e.g., from little-endian to big-endian or vice-versa, or conversion to a canonical data representation. There is also uncertainty in terms of response time because some servers may be more performant than others or may have a lower workload.

A major difference between the basic models of grid and cloud computing is that the former does not impose any restrictions regarding heterogeneity of the computing platforms, whereas homogeneity used to be a basic tenet of computer clouds' infrastructure. Originally, a computer cloud was a collection of homogeneous systems, systems with the same architecture and running under the same or very similar system software. We have already seen in Chapter 2 that nowadays, computer clouds exhibit some heterogeneity. For example, AWS instances use now x86 64-bit and ARM 64-bit processors; accelerated computing instances have attached co-processors with NVIDIA cores and TPUs.

The clients and the servers communicate through a network that can be congested. Transferring large volumes of data through the network can be time-consuming; this is a major concern for data-intensive applications in cloud computing. Communication through the network further lengthens the response time. Security becomes a major concern as the traffic between a client and a server can be intercepted.

Remote Procedure Call (RPC). RPCs were introduced in the early 1970s by Bruce Nelson and used for the first time at PARC (Palo Alto Research Park), a place credited with many innovative ideas in distributed systems including the development of the Ethernet, the GUI interfaces, bitmap displays, and the Alto system. The Network File System (NFS) introduced in 1984 was based on Sun's RPC. Many programming languages support RPCs; for example, Java Remote Method Invocation (Java RMI) provides a functionality similar to the one of UNIX RPC methods, and XML-RPC uses XML to encode HTML-based calls.

RPC is often used for implementation of client–server systems' interactions. To use an RPC, a process may use special services *PORTMAP* or *RPCBIND* available at port 111 to register and for

service lookup. RPC messages must be well-structured; they identify the RPC and are addressed to an RPC demon listening at an RPC port. *XDP* is a machine-independent representation standard for RPC. The RPC standard is described in RFC 1831.

RPCs reduce the *fate sharing* between caller and the callee. RPC's take longer than local calls due to communication delays. Several RPC semantics are used to overcome potential communication problems:

- i. *At least once*: a message is resent several times, and an answer is expected. The server may end up executing a request more than once, but an answer may never be received. This semantics is suitable for operation free of side effects.
- ii. *At most once*: a message is acted upon at most once. The sender sets up a timeout for receiving the response. When the timeout expires, an error code is delivered to the caller. This semantics requires the sender to keep a history of the time stamps of all messages because messages may arrive out-of-order. This semantics is suitable for operations which have side effects.
- iii. *Exactly once*: it implements the *at most once* semantics and requests an acknowledgment from the server.

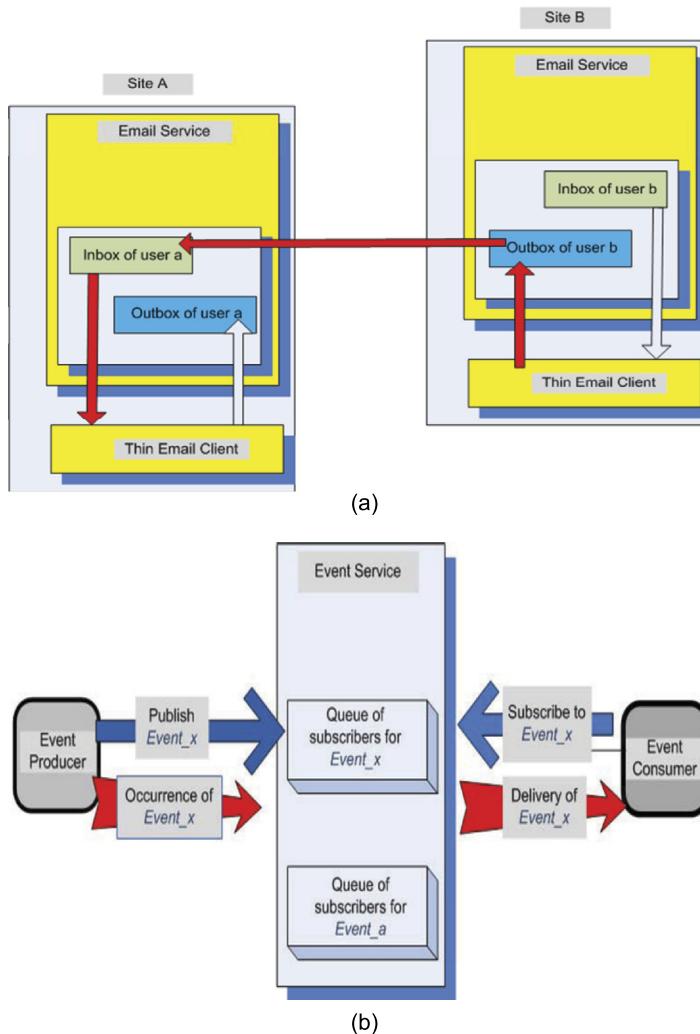
Applications of the client–server paradigm. The large spectrum of applications attests to the role played by the client–server paradigm in the modern computing landscape. Examples of popular applications of the client–server paradigm are numerous and include: the World Wide Web, the Domain Name System (DNS), the X-windows, electronic mail, see Fig. 3.10(a), event services, see Fig. 3.10(b), and so on.

The World Wide Web illustrates the power of the client–server paradigm and its effects on society. The web enables users to access *resources* such as text, images, digital music, and any imaginable type of information previously stored in a digital format. A *web page* is created using a description language called HTML (Hypertext Description Language). Information in each web page is encoded and formatted according to some standard, e.g., GIF or JPEG for images, MPEG for videos, MP3 or MP4 for audio, and so on.

The web is based upon a “pull” paradigm; the resources are stored at the server’s site and the client pulls them from the server. Some web pages are created “on the fly”; others are fetched from the disk. The client, called a *web browser*, and the server communicate using an application-level protocol called HTTP (HyperText Transfer Protocol) built on top of the TCP transport protocol.

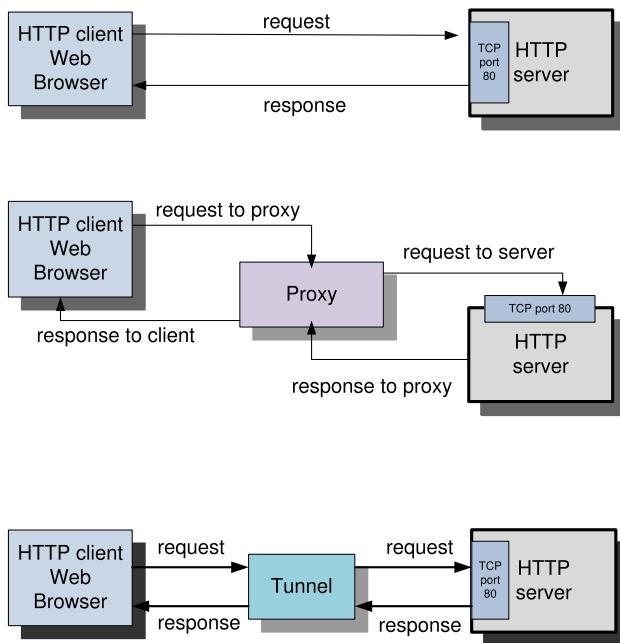
The web server also called an *HTTP server* listens for connections from clients at port 80. The sequence of events when a client browser sends an HTTP request to a server to retrieve some information and the server constructs the page on the fly and then the browser sends another HTTP request for an image stored on the disk is discussed next. First, a TCP connection between the client and the server is established using a *three-way handshake*. The client provides an arbitrary initial sequence number in a special segment with the *SYN* control bit on; then, the server acknowledges the segment and adds its own arbitrarily chosen initial sequence number; finally, the client sends its own acknowledgment *ACK*, as well as the HTTP request, and the connection is established.

The time elapsed from the initial request until the server’s acknowledgment reaches the client is called the RTT (Round-Trip Time). Once the TCP connection is established, the HTTP server takes its time to construct the page to respond to the first request; to satisfy the second request, the HTTP server must retrieve an image from the disk. The *response time* includes the RTT, the server residence time, and the data transmission time.

**FIGURE 3.10**

(a) Email service; sender and receiver communicate asynchronously using inboxes and outboxes. Mail daemons run at each site. (b) An event service supports coordination in a distributed system environment. The service is based on the *publish–subscribe* paradigm; an event producer publishes events, and an event consumer subscribes to events. The server maintains queues for each event and delivers notifications to clients when an event occurs.

The *response time*, defined as the time from the instance the first bit of the request is sent until the last bit of the response is received, consists of several components: the RTT, the *server residence time*, the time it takes the server to construct the response, and the data transmission time. RTT depends on the network latency, the time it takes a packet to cross the network from the sender to the receiver.

**FIGURE 3.11**

A client can communicate directly with the server, can communicate through a proxy, or may use tunneling to cross the network.

The data transmission time is determined by the network bandwidth. In turn, the server residence time depends on the server load.

Often, the client and the server do not communicate directly, but through a proxy server as shown in Fig. 3.11. Proxy servers could provide multiple functions; for example, they may filter client requests and decide whether or not to forward the request based on some filtering rules. A proxy server may redirect a request to a server in close proximity of the client or to a less-loaded server. A proxy can also act as a cache and provide a local copy of a resource, rather than forward the request to the server.

Another type of client–server communication is *HTTP-tunneling*, used most often as a means of communication from network locations with restricted connectivity. Tunneling means encapsulation of a network protocol; in our case HTTP acts as a wrapper for the communication channel between the client and the server, see Fig. 3.11.

3.13 Layering and hierarchy

Layering and hierarchy have been present in social systems since ancient times. For example, the Spartan Constitution, called Politeia, describes a Dorian society based on a rigidly layered social system and

a strong military. Nowadays, in a modern society, we are surrounded by organizations structured hierarchically. We have to recognize that layering and hierarchical organization have their own problems, could negatively affect the society, impose a rigid structure and affect social interactions, increase the overhead of activities, and prevent the system from acting promptly when such actions are necessary.

Layering demands modularity because each layer fulfills a well-defined function. The communication patterns in the case of layering are more restrictive; a layer is expected to communicate only with the adjacent layers. This restriction, the limitation of communication patterns, clearly reduces the complexity of the system and makes it easier to understand its behavior.

Modularity, layering, and hierarchy are critical for computer and communication systems. Large programs have been split into modules, each with a well-defined functionality since the early days of computing. Modules with related functionalities have then been grouped together into numerical, graphics, statistical, and many other libraries. Layering helps us dealing with complicated problems when we have to separate concerns that prevent us from making optimal design decisions. To do so, we define layers that address each concern and design clear interfaces between the layers.

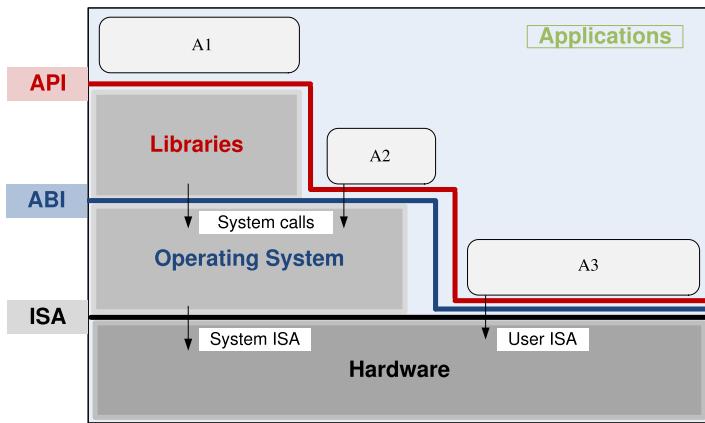
Probably the best example is layering of communication protocols. Early on, the need of accommodating a variety of physical communication channels that carry electromagnetic, optical, or acoustic signals, thus, the need for a *physical* layer, was recognized. The next concern is how to transport bits, not signals, between two systems directly connected to one another by a communication channel, i.e., the need for a *data link* layer.

Communication requires networks with multiple intermediate nodes. When bits have to traverse a chain of intermediate nodes from a source to the destination, the concern is how to forward the bits from one intermediate node to the next, so the *network* layer was introduced. Then, it was recognized that the source and the recipient of information are in fact outside the network and only want the data to reach the destination unaltered. Therefore, the *transport* layer was deemed necessary. Finally, the data sent and received has a meaning only in the context of an application, thus, the need for the *application* layer.

Strictly enforced layering can prevent optimizations. For example, cross-layer communication in networking was proposed to enable wireless applications to take advantage of information available at the Media Access Control (MAC) sub-layer of the data link layer. This example shows that layering gives us insight into where to place the basic mechanisms for error control, flow control, and congestion control of the network protocol stack.

An interesting question is whether a layered cloud architecture could be designed that has practical implications for the future development of computing clouds. One could argue that it may be too early for such an endeavor, that we need time to fully understand how to better organize a cloud infrastructure, and that we need to gather data to support the advantages of one approach over another. On the other hand, there are systems where it is difficult to envision a layered organization because of the complexity of the interaction between the individual modules. Consider for example an operating system that has a set of well-defined functional components:

1. The processor management subsystem, responsible for processor virtualization, scheduling, interrupt handling, and execution of privileged operations and system calls.
2. The virtual memory management subsystem, responsible for translating virtual addresses to physical addresses.
3. The multi-level memory management subsystem, responsible for transferring storage blocks between different memory levels, most commonly between primary and secondary storage.

**FIGURE 3.12**

Layering and interfaces between layers in a computer system. The software components including applications, libraries, and operating systems interact with the hardware via several interfaces: the Application Program Interface (API), the Application Binary Interface (ABI), and the Instruction Set Architecture (ISA). An application uses library functions (A1), makes system calls (A2), and executes machine instructions (A3).

4. The I/O subsystem, responsible for transferring data between the primary memory and the I/O devices.
5. The networking subsystem responsible for network communication.

The processor management interacts with all other subsystems, and there are also multiple interactions between the other subsystems; therefore, it seems unlikely that a layered organization would be feasible in this case.

Layers and interfaces between layers. A common approach to managing system complexity is to identify a set of *layers* with well-defined *interfaces* among them. The interfaces separate different levels of abstraction. Layering minimizes the interactions among the subsystems and simplifies the description of the subsystems. Each subsystem is abstracted through its interfaces with the other subsystems, thus, we are able to design, implement, and modify the individual subsystems independently.

The ISA (Instruction Set Architecture) defines the set of instructions of a processor; for example, the Intel architecture is represented by the *x86-32* and *x86-64* instruction sets for systems supporting 32-bit addressing and 64-bit addressing, respectively. The hardware supports two execution modes, a *privileged*, or *kernel* mode, and a *user* mode.

The instruction set consists of two sets of instructions, *privileged* instructions that can only be executed in kernel mode and the *non-privileged* instructions that can be executed in user mode. There are also *sensitive instructions* that can be executed in kernel and in user mode, but behave differently, see Section 5.4.

Computer systems are fairly complex and their operation is best understood when we consider a model similar with the one in Fig. 3.12, which shows the interfaces between the software components and the hardware [450]. The hardware consists of one or more multicore processors, a system intercon-

nect (e.g., one or more busses), a memory translation unit, the main memory, and I/O devices, including one or more networking interfaces.

Applications written mostly in High-Level Languages (HLL) often call library modules and are compiled into *object code*. Privileged operations, such as I/O requests, cannot be executed in user mode; instead, application and library modules issue *system calls*, and the operating system determines if the privileged operations required by the application do not violate system security or integrity and, if so, executes them on behalf of the user. Binaries resulting from translation of HLL programs target a specific hardware architecture.

The first interface, at the boundary of hardware and software, is the *Instruction Set Architecture* (ISA). The next interface is the *Application Binary Interface* (ABI), which allows the ensemble consisting of the application and the library modules to access the hardware. ABI does not include privileged system instructions; rather, it invokes system calls.

Finally, the *Application Program Interface* (API) defines the set of instructions the hardware was designed to execute and gives the application access to the ISA. API includes HLL library calls, which often invoke system calls. Recall that a *process* is the abstraction for the code of an application at execution time; a *thread* is a light-weight process. ABI is the projection of the computer system seen by the process, and API is the projection of the same system from the perspective of the HLL program.

Binaries created by a compiler for a specific ISA and a specific operating systems are not portable. Such code cannot run on a computer with a different ISA, or on the computer with the same ISA but a different OS. But it is possible to compile an HLL program for a VM environment. In this case, portable code is produced and distributed and then converted by binary translators to the ISA of the host system. A *dynamic binary translation* converts blocks of guest instructions from the portable code to the host instructions and leads to a significant performance improvement because such blocks are cached and reused.

3.14 Peer-to-peer systems

The distributed systems discussed in this chapter allow access to resources in a tightly controlled environment. System administrators enforce security rules and control the allocation of physical, rather than virtual resources. In all models of network-centric computing prior to utility computing, a user maintained direct control of the software and the data residing on remote systems.

This user-centric model, in place since the early 1960s, was challenged in the 1990s by the peer-to-peer (P2P) model. P2P systems share some ideas with computer clouds. The new distributed computing model promoted the idea of low-cost access to storage and CPU cycles provided by participant systems. In this case, the resources are located in different administrative domains. The P2P systems are self-organizing and decentralized, while the servers in a cloud are in a single administrative domain and have a central management.

P2P systems exploit the network infrastructure to provide access to distributed computing resources. Decentralized applications developed in the 1980s, such as SMTP (Simple Mail Transfer Protocol), a protocol for Email distribution, and NNTP (Network News Transfer Protocol), an application protocol for dissemination of news articles, are early examples of P2P systems. Systems developed in the late 1990s, such as the music-sharing system Napster, gave participants access to storage distributed over the

network. The first volunteer-based scientific computing, SETI@home, used free cycles of participating systems to carry out compute-intensive tasks.

P2P model represents a significant departure from the client–server model, the cornerstone of distributed applications for several decades. P2P systems have several desirable properties [420]:

- i. Require a minimally dedicated infrastructure; participating systems provide resources.
- ii. Are highly decentralized.
- iii. Are scalable; individual nodes are not required to be aware of the global state.
- iv. Are resilient to faults and attacks because few of their elements are critical for the delivery of service and the abundance of resources can support a high degree of replication.
- v. Individual nodes do not require excessive network bandwidth.
- vi. The dynamic unstructured system architecture shields P2P systems from censorship.

The undesirable properties of peer-to-peer systems are also notable. Decentralization raises the question if P2P systems can be managed effectively and provide the security required by various applications. The fact that they are shielded from censorship makes them a fertile ground for illegal activities, including the distribution of copyrighted content.

In spite of its problems, the new paradigm was embraced by applications other than file sharing. Since 1999, new P2P applications, such as the ubiquitous Skype, a voice over IP telephony service,⁵ data-streaming applications such as Cool Streaming [537] and BBC’s online video service, content distribution networks such as CoDeeN [503], and volunteer computing applications based on the BOINC (Berkeley Open Infrastructure for Networking Computing) platform [26], have proved their appeal to users.

Skype reported in 2008 that 276 million registered users have used more than 100 billion minutes for voice and video calls. The site www.boinc.berkeley.edu reports that at the end of June 2012, volunteer computing involved more than 275 000 individuals and more than 430 000 computers providing a monthly average of almost 6.3×10^9 MFlops. It is also reported that the P2P traffic accounts for a very large fraction of the Internet traffic, with estimates ranging from 40% to more than 70%.

Many groups from industry and academia have rushed to develop and test new ideas taking advantage of the fact that P2P applications do not require a dedicated infrastructure. Applications such as Chord [457] and Credence [502] address issues critical for the effective operation of decentralized systems. Chord is a distributed lookup protocol to identify the node where a particular data item is stored. The routing tables are distributed and, while other algorithms for locating an object require the nodes to be aware of most of the nodes of the network, Chord maps a key related to an object to a node of the network using routing information about a few nodes only.

Credence is an object reputation and ranking scheme for large-scale P2P file sharing systems. Reputation is of paramount importance for systems that often include many unreliable and malicious nodes. In the decentralized algorithm used by *Credence*, each client uses local information to evaluate the reputation of other nodes and shares its own assessment with its neighbors. The credibility of a node depends only on the votes it casts.

⁵ Skype allows close to 700 million registered users from many countries around the globe to communicate using a proprietary voice-over-IP protocol. The system developed in 2003 by Niklas Zennström and Julius Fris was acquired by Microsoft in 2011 and nowadays is a hybrid P2P and client–server system.

Each node computes the reputation of another node based solely on the degree of matching with its own votes and relies on like-minded peers. Overcite [461] is a P2P application to aggregate documents based on a three-tier design. The web front-ends accept queries and display the results, while servers crawl through the web to generate indexes and to perform keyword searches; the web back-ends store documents, metadata, and coordination state on the participating systems.

The rapid acceptance of the new paradigm has triggered the development of a new communication protocol allowing hosts at the network periphery to cope with the limited network bandwidth available to them. BitTorrent is a peer-to-peer file sharing protocol enabling a node to download/upload large files from/to several hosts simultaneously.

The P2P systems differ in their architecture. Some do not have any centralized infrastructure, while others have a dedicated controller, but this controller is not involved in resource-intensive operations. For example, Skype has a central site to maintain the user accounts; the users sign in and pay for specific activities on this site. The controller for a BOINC platform maintains membership and is involved in task distribution to participating systems. The nodes with abundant resources in systems without any centralized infrastructure often act as *supernodes* and maintain information useful to increasing the system efficiency, e.g., indexes of the available content.

Regardless of the architecture, P2P systems are built around an *overlay network*, a virtual network superimposed over the real network. Each node maintains a table of *overlay links* connecting it with other nodes of this virtual network, each node being identified by its IP addresses. Two types of overlay networks, *unstructured* and *structured*, are used by P2P systems. Random walks starting from a few bootstrap nodes are usually used by systems desiring to join an unstructured overlay.

Each node of a structured overlay has a unique key that determines its position in the structure; the keys are selected to guarantee a uniform distribution in a very large name space. Structured overlay networks use *key-based routing* (KBR); given a starting node v_0 and a key k , the function $KBR(v_0, k)$ returns the path in the graph from v_0 to the vertex with key k . Epidemic algorithms are often used by unstructured overlays to disseminate the network topology.

3.15 Large-scale systems

The developments in computer architecture, storage technology, networking, and software during the last several decades of the 20th century, coupled with the need to access and process information, led to several large-scale distributed system developments:

The web and the semantic web expected to support composition of services (not necessarily computational services) available on the web. The web is dominated by unstructured or semi-structured data, while the semantic web advocates inclusion of semantic content in web pages.

The Grid, initiated in the early 1990s by National Laboratories and universities primarily for applications in science and engineering.

The need to share data from high energy physics experiments motivated Sir Tim Berners-Lee, who worked at CERN at Geneva in the late 1980s, to put together the two major components of the World Wide Web: HTML (Hypertext Markup Language) for data description and HTTP (Hypertext Transfer Protocol) for data transfer. The web opened a new era in data sharing and ultimately led to the concept of network-centric content.

The *Semantic Web* is an effort to enable lay people to find, share, and combine information available on the web more easily. The name was coined by Berners-Lee to describe “a web of data that can be processed directly and indirectly by machines.” It is a framework for data sharing among applications based on the Resource Description Framework (RDF). In this vision, information can be readily interpreted by machines, so machines can perform the tedious work involved in finding, combining, and acting upon information on the web.

The semantic web is “largely unrealized” according to Berners-Lee. Several technologies are necessary to provide a formal description of concepts, terms, and relationships within a given knowledge domain; they include the Resource Description Framework (RDF), a variety of data interchange formats, and notations such as RDF Schema (RDFS) and the Web Ontology Language (OWL).

Gradually, the need to make computing more affordable and to liberate the users from the concerns regarding system and software maintenance reinforced the idea of concentrating computing resources in data centers. Initially, these centers were specialized, each running a limited palette of software systems, as well as applications developed by the users of these systems. In the early 1980s major research organizations, such as the National Laboratories and large companies, had powerful computing centers supporting large user populations scattered throughout wide geographic areas. Then the idea to link such centers in an infrastructure resembling the power grid was born; the model known as network-centric computing was taking shape.

A *computing grid* is a distributed system consisting of a large number of loosely coupled, heterogeneous, and geographically dispersed systems in various administrative domains. The term *computing grid* is a metaphor for accessing computer power with similar ease as we access power provided by the electric grid. Software libraries known as *middleware* have been furiously developed since the early 1990s to facilitate access to grid services.

The vision of the grid movement was to give a user the illusion of a very large virtual supercomputer. The autonomy of the individual systems and the fact that these systems were connected by wide-area networks with latency higher than the latency of the interconnection network of a supercomputer posed serious challenges to this vision. Nevertheless, several “Grand Challenge” problems, such as protein folding, financial modeling, earthquake simulation, and climate/weather modeling, run successfully on specialized grids. The Enabling Grids for E-science project is arguably the largest computing grid; along with the LHC Computing Grid (LCG), the E-science project aims to support the experiments using the Large Hadron Collider (LHC) at CERN which generates several gigabytes of data per second, or 10 PB (petabytes) per year.

In retrospect, two basic assumptions about the infrastructure prevented the grid movement from having the impact its supporters were hoping for. The first is the heterogeneity of the individual systems interconnected by the grid. The second is that systems in different administrative domain are expected to cooperate seamlessly. Indeed, the heterogeneity of the hardware and of the system software poses significant challenges for application development and for application mobility.

At the same time, critical areas of system management including scheduling, optimization of resource allocation, load balancing, and fault-tolerance are extremely difficult in a heterogeneous system. The fact that resources are in different administrative domains further complicates many, already difficult, problems related to security and resource management. While very popular in the science and the engineering communities, the grid movement did not address the major concerns of enterprise computing community and did not make a noticeable impact on the IT industry.

Cloud computing is largely viewed as the next big step in the development and deployment of an increasing number of distributed applications. The companies promoting cloud computing seem to have learned the most important lessons from the grid movement. Computer clouds are typically homogeneous. An entire cloud shares the same security, resource management, cost, and other policies, and last, but not least, it targets enterprise computing. These are some of the reasons why several agencies of the US government, including Health and Human Services, the Center for Disease Control (CDC), NASA, the Navy's Next Generation Enterprise Network (NGEN), and Defense Information Systems Agency (DISA), have launched cloud computing initiatives and conduct actual system developments intended to improve the efficiency and effectiveness of their information processing needs.

3.16 Composability bounds and scalability (R)

Nature creates complex systems from simple components. For example, a vast variety of proteins are linear chains assembled from 20 amino acids, the building blocks of proteins found in human body. These amino acids are naturally incorporated into polypeptides and are encoded by the genetic code. Imitating nature, manmade systems are assembled from sub-assemblies; in turn, a sub-assembly is made from several modules, each module could consist of sub-modules, and so on. Composability has natural bounds imposed by the laws of physics as we have seen when discussing heat dissipation of solid-state devices. As the number of components increases, the complexity of a system also increases.

The limits of composability can be reached when the physical size of individual components changes. A recent paper with the suggestive title “When every atom counts” [350] shows that even the most modern solid-state fabrication facilities cannot produce chips with consistent properties. The percentage of defective or substandard chips has been constantly increasing as components have become smaller and smaller.

The lack of consistency in the manufacturing process of solid-state devices is attributed to the increasingly smaller size of the physical components of a chip. This problem is identified by the International Technology Roadmap for Semiconductors as “a red brick wall,” a problem without a clear solution, a wall that could prevent further progress. Chip consistency is no longer feasible because the transistors and the “wires” on a chip are so small that random differences in the placement of an atom can have a devastating effect, e.g., it can increase the power consumption by an order of magnitude and slowdown the chip by as much as 30%.

As features become smaller and smaller, the range of the *threshold voltage*, the voltage needed to turn a transistor on and off, has been widening; since now many transistors have this threshold voltage at or near zero, they cannot operate as switches. While the range for 28-nm technology was approximately between +0.01 and +0.4 V, the range for 20-nm technology is between -0.05 and +0.45 V, and the range becomes even wider, from -0.18 to +0.55 for 14-nm technology.

There are physical bounds for the composition of analog systems: Noise accumulation, heat dissipation, cross-talk, the interference of signals on multiple communication channels, and several other factors limit the number of components of an analog system. Digital systems have more distant bounds, but composability is still limited by physical laws.

There are virtually no bounds on the composition of digital computing and communication systems controlled by software. The Internet is a network of networks and a prime example of composability

with distant bounds. Computer clouds are another example; a cloud is composed of a very large number of servers and interconnects, each server is made up of multiple processors, and each processor has multiple cores. Software is the ingredient that pushes the composability bounds and liberates computer and communication system from the limits imposed by physical laws.

In the physical world the laws valid at one scale break down at a different scale, e.g., the laws of classical mechanics are replaced at the atomic and subatomic scales by quantum mechanics. Thus, we should not be surprised that scale really matters in the design of computing and communication systems. Indeed, architectures, algorithms, and policies that work well for systems with a small number of components very seldom scale up.

For example, many computer clusters have a front-end that acts as the nerve center of the system, manages communication with the outside world, monitors the entire system, and supports system administration and software maintenance. A computer cloud has multiple such nerve centers, and new algorithms to support collaboration among these centers must be developed. Scheduling algorithms that work well within the confines of a single system cannot be extended to collections of autonomous systems when each system manages local resources; in this case, as in the previous example, entities must collaborate with one another, and this requires communication and consensus.

Another manifestation of this phenomenon is in the vulnerabilities of large-scale distributed systems. The implementation of Google's Bigtable revealed that many distributed protocols designed to protect against network partitions and fail-stop are unable to cope with failures due to scale [94]. Memory and network corruption, extended and asymmetric network partitions, systems that fail to respond, and large clock skews occur with increasing frequency in a large-scale system, and they interact with one another in a manner that greatly affects the overall system availability.

Scaling has other dimensions than just the number of components, and space plays an important role. The communication latency is small when the component systems are clustered together within a small area, and this allows us to implement efficient algorithms for global decision making, e.g., consensus algorithms. When, for the reasons discussed in Section 1.4, the data centers of a cloud provider are distributed over a large geographic area, transactional database systems are of little use for most online transaction-oriented systems, and a new type of data store has to be introduced in the computational ecosystem.

Societal scaling means that a service is used by a very large segment of the population and/or is a critical element of the infrastructure. There is no better example to illustrate how societal scaling affects the system complexity than communication supported by the Internet. The infrastructure supporting the service must be highly available. A consequence of redundancy and of the measures to maintain consistency is increased system complexity.

At the same time, the popularity of the service demands simple and intuitive means to access the infrastructure. Again, the system complexity increases due to the need to hide the intricate mechanisms from a lay person with little understanding of the technology. The vulnerability of wireless systems has increased due to the desire to design wireless devices that: (a) operate efficiently in terms of power consumption; (b) present the user with a simple interface and few choices; and (c) satisfy a host of other popular functions. This is happening at a time when not many smartphone and tablet users understand the security risks of wireless communication.

3.17 Distributed computing fallacies and the CAP theorem

Distributed processing became a reality after Berkeley sockets were released in 1983 as a programming interface in 4.2 BSD Unix operating system. A *socket* is an abstract representation of the local endpoint of a network communication path. A socket is like a file descriptor in the Unix philosophy; it provides a common interface for input and output streams of data. Sockets became the standard interface for applications running in the Internet. All operating systems implement a version of Berkeley socket interface.

The group at UC Berkeley, credited with the developments of sockets, founded Sun Microsystems in 1982. Sun is credited with a number of remarkable contributions to distributed computing, including the Java programming language, the Solaris operating system, and the Network File System (NFS), as well as SPARC microprocessors. At Sun, Bill Joy and Dave Lyon formulated a list of flawed assumptions about distributed computing. James Gosling, best known as the founder and lead designer of the Java programming language, codified four of these assumptions as “The Fallacies of Networked Computing.”

According to the Merriam–Webster dictionary, a fallacy is a false or mistaken idea or an often plausible argument using false or invalid inferences. In the early days of distributed processing, and even later, many developers of distributed applications made a number of false assumptions. The list of fallacies was then extended: Peter Deutsch added another three, and James Gosling contributed the most recent one, see <http://java.sys-con.com/read/38665.htm>. Ignoring the problems caused by each one of these fallacies has severe consequences:

1. *The network is reliable.* This assumption causes applications to stall or infinitely wait for an answer; memory and other resources may fail to retry stalled operations or require a manual restart in case of network outage.
2. *The latency is zero and no packets are lost.* This causes application- and transport-layer developers to assume unbounded traffic, greatly increasing the number of dropped packets and wasting bandwidth.
3. *The bandwidth is infinite.* Assuming infinite bandwidth and repeated packet retransmissions can cause bottlenecks.
4. *The network is secure.* Ignoring malicious users is a capital sin.
5. *The topology does not change.* Topology changes affect both bandwidth and latency.
6. *There is one administrator.* There may be conflicting policies preventing packets to reach their destination, and these may limit user access to some resources.
7. *The transport cost is zero.* The costs of building and maintaining networks cannot be ignored.
8. *The network is homogeneous.* The Internet is a network of networks with various different latencies and bandwidths.

Grid computing movement initiated in late 1980s captured the interest of the distributed research community for more than two decades, in spite of ignoring the consequences of the eight fallacies. This movement aimed to create a “super virtual computer” composed of many loosely coupled computers acting together to perform large computational tasks. Grid computers tended to be heterogeneous with systems in different administrative domains; moreover, the computers were mostly connected by networks with limited communication bandwidth and high latency. The Grid movement faded away for these reasons and because it did not attract applications other than those in computational science and engineering, it was eventually replaced by cloud computing.

Networking has dramatically changed since 1991. The effects of some of these fallacies have been amplified; in fact, malicious users are now able to threaten the critical infrastructure of the society. The consequences of other fallacies have been attenuated; the communication technology has evolved, bandwidth of terabytes per second are not uncommon, error rates are low, and today's routers have substantial resources to support various types of traffic. Application developers have now access to cloud computing services that mostly shield them from the effects of the eight fallacies.

The CAP theorem due to Eric Brewer from UC Berkeley states that it is impossible for a distributed data store to simultaneously provide more than two out of the three guarantees: (i) consistency—every read receives either the most recent write or an error; (ii) availability—every request receives a (non-error) response, without the guarantee that it contains the most recent write; and (iii) partition tolerance—the system continues to operate despite network errors when an arbitrary number of messages are dropped or delayed [198].

Tradeoffs are common, for example, consistency can be traded for availability. For example, when running of a cluster, a subset of nodes being consistent is sufficient for practical consistency. Write Quorum requires that the number of nodes writing successfully should be more than half the number of nodes involved in replication.

In 2010, Daniel Abadi from Yale University argued that “Ignoring the consistency/latency trade-off of replicated systems is a major oversight [in CAP], as it is present at all times during system operation, whereas CAP is only relevant in the arguably rare case of a network partition” [2]. The PACELC theorem due to Abadi is an extension of CAP stating that “in case of network partitioning (P) in a distributed computer system, one has to choose between availability (A) and consistency (C) (as per the CAP theorem), but else (E), even when the system is running normally in the absence of partitions, one has to choose between latency (L) and consistency (C).” Section 7.9 includes an in-depth discussion of PACELC implications for cloud storage systems.

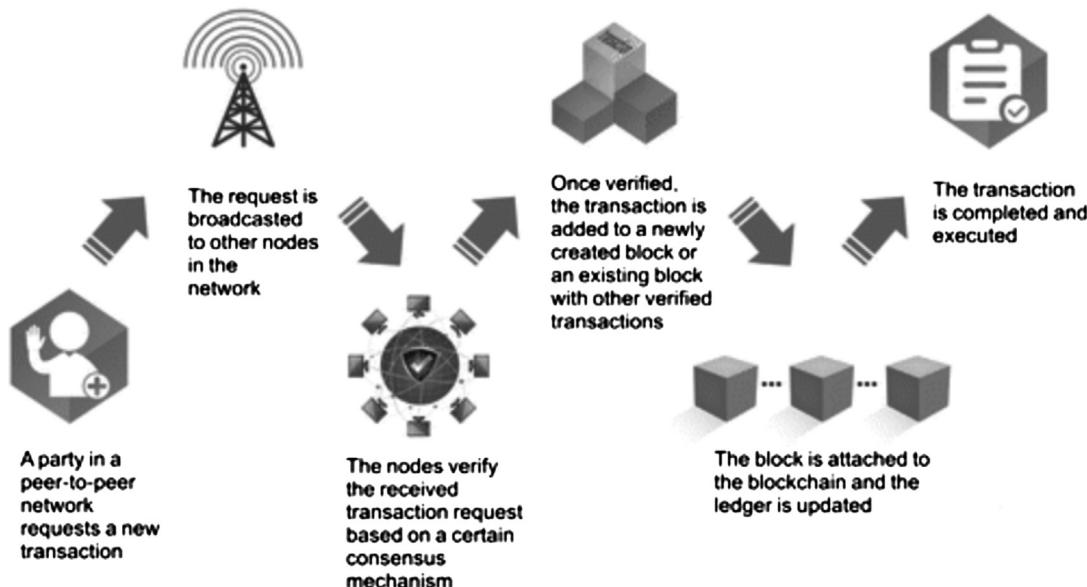
Note that the blockchain technology gives the impression that CAP is invalid since it sacrifices consistency for availability and partition tolerance, but availability and partition tolerance is achieved through validation among the nodes over time.

3.18 Blockchain technology and applications

A *blockchain* is a decentralized, distributed, digital ledger consisting of records called *blocks* recording transactions so that no block can be altered retroactively without the alteration of all subsequent blocks; it was invented in 2008 by an individual of unknown identity. A blockchain is managed by a P2P network collectively adhering to a protocol for inter-node communication and for new block validation.

A block contains batches of valid transactions hashed and encoded into a Merkle tree, see also Section 7.13. Every node of the system has a copy of the blockchain, and data quality is maintained by massive database replication and computational trust. Fig. 3.13 illustrates the flow of blockchain information. Cryptocurrencies such as Bitcoin made blockchain technology ubiquitous.

Bitcoin mining is the process by which new Bitcoins are entered into circulation and it is also a critical component of the maintenance and development of the blockchain ledger. Powerful ASICs (Application Specific Integrated Circuits) are used for bitcoin mining; for example, WhatsMiner M30S++ delivers 112 TH/s (Tera Hashing operations per second). Aside from the short-term Bitcoin payoff, being a coin miner gives “voting” power when changes in the Bitcoin network protocol are proposed.

**FIGURE 3.13**

Blockchain flow of information [318].

While blockchain is widely used by cryptocurrencies, there are many other potential applications of this technology with immense economic and social implications, [318]. In recent years, we have seen an explosion of interest in blockchain that “could someday underlie everything from how we vote to who we connect with online to what we buy” (Wall Street Journal 2018, p. B4). For example, GSA (General Service Administration) with contracts of some \$55 billions per year is working with a technology company to develop a new procurement blockchain for its vendors network. Blockchain could be used in detecting counterfeits by associating unique identifiers to individual items and storing records associated with transactions involving the item that cannot be forged or altered. A non-fungible token (NFT) is a unit of data stored on a blockchain that certifies a digital asset to be unique. NFTs are now widely used to sign digital art and other collectibles.

Blockchain exhibits fundamental properties unreachable by other technologies. *Decentralized consensus* is a hard problem; achieving consensus in a decentralized network requires careful design of the algorithms. Blockchain enables data integrity because the complete copy of identical information can be held by anyone who has access to it and wants to keep it. *Machine-based automation* gives blockchain the freedom to bypass human actors’ unpredictability and inability to process massive amounts of information.

The rise of blockchain transactions has led to increased environmental criticism. The proof-of-work blockchain process is computationally intensive, requiring high energy consumption contributing to global warming. According to the University of Cambridge’s Bitcoin electricity consumption index, Bitcoin miners are expected to consume roughly 130 Terawatt-hours of energy (TWh), roughly 0.6% of global electricity consumption in 2022 [125].

3.19 History notes and further readings

Ex nihilo nihil, the philosophical dictum first appearing in Aristotle's Physics⁶ translated as "nothing comes from nothing," applies to every human endeavor and means to search for the origins of everything to grasp the intellectual challenges faced along the way by the great minds of the past. This motivates us to take a brief excursion into the history of computing.

Two theoretical developments in 1930s were critical for the development of modern computers. The first was the publication of Alan Turing's 1936 paper [481]. The paper provided a definition of a universal computer, called a Turing Machine, which executes a program stored on tape; the paper also proved that there were problems, such as the halting problem, that could not be solved by any sequential process. The second major development was the publication in 1937 of Claude Shannon's master's thesis at MIT "A Symbolic Analysis of Relay and Switching Circuits" in which he showed that any Boolean logic expression can be implemented using logic gates. This should remind us how much the digital era owes to George Boole, a humble professor at University College Cork, in Ireland!

The first Turing complete⁷ computing device was Z3, an electro-mechanical device built by Konrad Zuse in Germany in May 1941; Z3 used a binary floating-point representation of numbers and was program-controlled by a film-stock. The first programmable electronic computer ENIAC, built at the Moore School of Electrical Engineering at the University of Pennsylvania by a team led by John Prosper Eckart and John Mauchly, became operational in July 1946 [341]; ENIAC, unlike Z3, used a decimal number system and was program-controlled by patch cables and switches.

John von Neumann, the famous mathematician and theoretical physicist, contributed fundamental ideas for modern computers [79,497,498]. He was one of the most brilliant minds of the 20th century, with an uncanny ability to map fuzzy ideas and garbled thoughts to crystal clear and scientifically sound concepts. John von Neumann drew the insight for the stored-program computer from Alan Turing's work⁸ and from his visit to the University of Pennsylvania; he thought that ENIAC was an engineering marvel, but was less impressed with the awkward manner to "program" it by manually connecting cables and setting switches. He introduced the so-called "von Neumann architecture" in a report published in the 1940s; to this day, he is faulted by some because he failed to mention in this report the sources of his insight.

John von Neumann led the development at the Institute of Advanced Studies at Princeton of MANIAC, an acronym for "mathematical and numerical integrator and computer." MANIAC was closer to the modern computers than any of its predecessors; it was used for sophisticated calculations required by the development of the hydrogen bomb nicknamed "Ivy Mike" secretly detonated on November 1, 1952, over an island that no longer exists in the South Pacific. In a recent book [157], the historian of science George Dyson writes: "The history of digital computing can be divided into an Old Testament whose prophets, led by Leibnitz, supplied the logic, and a New Testament whose prophets led by von Neumann built the machines. Alan Turing arrived between them."

⁶ This dictum comes from ancient Greek cosmology and states that everything has its origin in something.

⁷ A Turing complete computer is equivalent to a universal Turing machine except for memory limitations.

⁸ Alan Turing came to the Institute of Advanced Studies at Princeton in 1936 and got his Ph.D. there in 1938; von Neumann offered him a position at the Institute, but, as the dark clouds signaling the approaching war were gathering over Europe, Turing decided to go back to England.

In 1951, Sir Maurice Vincent Wilkes developed the concept of microprogramming first implemented on the EDSAC 2 computer. Wilkes is also credited with the idea of symbolic labels, macros and subroutine libraries, and caching.

Third-generation computers were built during the period 1964–1971; they made extensive use of integrated circuits (ICs) and ran under the control of an operating systems. MULTIX (Multiplexed Information and Computing Service) was an early time-sharing operating system developed in 1963 by MIT, GE, and Bell Labs for the GE 645 mainframe. MULTIX had a lasting impact on the design and implementation of computer systems, had numerous novel features, and implemented a fair number of interesting concepts such as: a hierarchical file system, access control lists for file information sharing, dynamic linking, and online reconfiguration [115,116].

In his address “A Career in Computer System Architecture”, MIT Professor Jack Dennis wrote: “In 1960 Professor John McCarthy, now at Stanford University and known for his contributions to artificial intelligence, led the *Long Range Computer Study Group* (LRCSG) which proposed objectives for MIT’s future computer systems. I had the privilege of participating in the work of the LRCSG, which led to Project MAC and the Multics computer and operating system, under the organizational leadership of Prof. Robert Fano and the technical guidance of Prof. Fernando Corbató. At this time, Prof. Fano had a vision of the Computer Utility—the concept of the computer system as a repository for the knowledge of a community—data and procedures in a form that could be readily shared—a repository that could be built upon to create ever more powerful procedures, services, and active knowledge from those already in place. Prof. Corbató’s goal was to provide the kind of central computer installation and operating system that could make this vision a reality. With funding from DARPA, the Defense Advanced Research Projects Agency, the result was Multics...in the 1970s I found it easy to get government funding. The agencies were willing to fund pretty wild ideas, and I was supported to do research on *data flow* architecture, first by NSF and later by the DOE” (<http://csg.csail.mit.edu/Users/dennis/essay.htm>).

The development of the UNIX system was a consequence of the withdrawal of Bell Labs from the MULTIX project in 1968. UNIX was developed in 1969 for a DEC PDP minicomputer by a group led by Kenneth Thompson and Dennis Ritchie [417]. According to [416], “the most important job of UNIX is to provide a file-system;” the same reference discusses another concept introduced by the system: “For most users, communication with UNIX is carried on with the aid of a program called the Shell. The Shell is a command line interpreter: it reads lines typed by the user and interprets them as requests to execute other programs.”

The first microprocessor, Intel 4004, was announced in 1971; it performed binary-coded decimal (BCD) arithmetic using 4-bit words. Intel 4004 was followed in by Intel 8080, the first 8-bit microprocessor, and by its competitor, Motorola 6800, released in 1974. The first 16-bit multi-chip microprocessor, IMP-16, was announced in 1973 by National Semiconductors. The 32-bit microprocessors appeared in 1979; it widely used Motorola MC68000, had 32-bit registers, and supported 24-bit addressing. Intel’s 80286 was introduced in 1982. The 64-bit processor era was inaugurated by AMD64, an architecture called x86-64, backward compatible with Intel x86 architecture. Dual-core processors appeared in 2005; multicore processors are ubiquitous in today’s servers, PCs, tablets, and even smartphones.

The development of distributed systems was only possible after major advances in communication technology. In early 1960s Leonard Kleinrock from UCLA developed theoretical foundations for packet

switching networks and in early 1970s for hierarchical routing in packet switching networks. Kleinrock published the first paper on packet switching theory in 1961 and the first book in 1964.

The Advanced Research Projects Agency (ARPA), created in 1958, funded research at multiple universities and businesses sites and the ARPANET project to connect them all within a network was initiated. The Internet is a global network based on the Internet Protocol Suite (TCP/IP); its origins can be traced back to 1965 when Ivan Sutherland, the Head of the Information Processing Technology Office (IPTO) at ARPA, encouraged Lawrence Roberts, who had worked previously at MITs Lincoln laboratories, to become the Chief Scientist at IPTO and to initiate a networking project based on packet switching rather than circuit switching.

In August 1968, DARPA (Defence Advanced Research Projects Agency), the successor of ARPA, released a request for quotation (RFQ) for the development of packet switches called Interface Message Processors (IMPs). A group from Bolt Beranek and Newman (BBN) won the contract. Several researchers including Robert Kahn from BBN, Lawrence Roberts from DARPA, Howard Frank from Network Analysis Corporation, and Leonard Kleinrock from UCLA and their teams played a major role in the overall ARPANET architectural design. The idea of open-architecture networking was first introduced by Kahn in 1972, and his collaboration with Vincent Cerf from Stanford led to the design of TCP/IP. Three groups, one at Stanford, one at BBN, and one at UCLA, won the DARPA contract to implement TCP/IP.

In 1969, BBN installed the first IMP at UCLA. The first two ARPANET nodes interconnected were the Network Measurement Center at the UCLA and SRI International in Menlo Park, California. Two more nodes were added at UC Santa Barbara and the University of Utah. By the end of 1971, there were 15 sites interconnected by ARPANET.

Ethernet technology, developed by Bob Metcalfe at Xerox PARC in 1973 and other local area network technologies, such as token passing rings, allowed PCs and workstations to be connected to the Internet in the 1980s. The Domain Name System (DNS) was invented by Paul Mockapetris of USC/ISI. The DNS permitted a scalable distributed mechanism for resolving hierarchical host names into an Internet address.

UC Berkeley with support from DARPA rewrote the TCP/IP code developed at BBN and incorporated it into the Unix BSD system. In 1985, Dennis Jennings started the NSFNET program at NSF to support the general research and academic communities.

The first distributed computing programs were a pair of worm programs called Creeper and Reaper. In the 1970s, Creeper was using the idle CPU cycles of processors in the ARPANET to copy itself onto the next system and then delete itself from the previous one. Then it was modified to remain on all previous computers. Reaper deleted all copies of the Creeper.

Further readings. Several texts are highly recommended. “Computer Architecture: A Quantitative Approach” [232] by John Hennessy and David Patterson is the authoritative reference for computer architecture. The text “Principles of Computer Systems Design” co-authored by Jerome Saltzer and Frans Kaashoek [430] covers basic concepts in computer system design. “Computer Networks: A Top-Down Approach Featuring the Internet” by James Kurose and Keith Ross is a good introduction to networking.

Amdahl’s paper [21] is a classic, and [416] and [417] are the references for UNIX. [238] covers multicore processors. Load balancing in distributed system is analyzed in [158]. A comprehensive survey of peer-to-peer systems was published in 2010 [420]. *Chord* [457] and *Credence* [502] are important references in the area of peer-to-peer systems.

3.20 Exercises and problems

- Problem 1.** Do you believe that the homogeneity of a large-scale distributed systems is an advantage? Discuss the reasons for your answer. What aspects of hardware homogeneity are the most relevant in your view and why? What aspects of software homogeneity do you believe are the most relevant and why?
- Problem 2.** Peer-to-peer systems and clouds share a few goals, but not the means to accomplish them. Compare the two classes of systems in terms of architecture, resource management, scope, and security.
- Problem 3.** Explain briefly how the *publish–subscribe* paradigm works, and discuss its application to services such as bulletin boards, mailing lists, etc. Outline the design of an event service based on this paradigm, see Fig. 3.10(b). Can you identify a cloud service that emulates an event service?
- Problem 4.** Tuple spaces can be thought of as an implementation of a distributed shared-memory. Tuple spaces have been developed for many programming languages, including Java, Lisp, Python, Prolog, Smalltalk, and Tcl (Tool Command Language). Explain briefly how tuple spaces work. How secure and scalable are the tuple spaces, e.g., JavaSpaces, you are familiar with?
- Problem 5.** Arithmetic intensity is defined as the number of floating-point operations divided by the number of bytes in the main memory accessed for running a program. (1) Give examples of computations with low, medium, and high arithmetic intensity. (2) Derive a formula relating the attainable performance to the peak performance of a processor, the peak memory bandwidth, and the arithmetic intensity. (3) How is this formula related to the roofline model?
- Problem 6.** Read Section 3.5 of [232] regarding dynamic instruction scheduling and answer the following questions: (1) What is the role of the reservation stations used by Tomasulo's approach for dynamic instruction scheduling? (2) What are instruction execution steps for dynamic instruction scheduling? (3) Draw a diagram of a system with two reservation stations.
- Problem 7.** There are several approaches to hardware multithreading, fine-grain, coarse-grain, and simultaneous (SMT). What are the benefits and the disadvantages of each one of them, and if and where are they used.
- Problem 8.** Amazon offers EC2 instances with GPU coprocessors for data-intensive applications. (1) Is it beneficial to have the GPUs attached as coprocessors, or is it a disadvantage? (2) Is the thread scheduling inside a GPU done by hardware or controlled by the application? (3) Is it beneficial to add more CUDA threads to an application?
- Problem 9.** Reference [52] analyzes the power consumption of computing, storage, and networking infrastructure in a cloud data center. Find the percentage of the total power consumed by the CPUs, the DRAM, the disks, the networking, and the other consumers, and then suggest the most effective means to reduce the power consumption.
- Problem 10.** Given the 32-bit processor with 64 Kbyte cache, 32 byte block and one block per cache set from Section 3.2 draw a figure showing the cache lines and the 32 bit address generated by the CPU. Show how the bits used for the offset of a word in a block, the index, and the tag are used to identify the block in cache. Repeat the task for a processor with the same cache and block size, but with a two-way set associative cache.

Cloud hardware and software

4

This chapter presents the hardware and the software stack for cloud computing. In their quest to provide reliable, low-cost services, cloud service providers (CSPs) exploit the latest computing, communication, and software technologies to offer a highly available, easy-to-use, and efficient cloud computing infrastructure.

The cloud hardware infrastructure is built with inexpensive, off-the-shelf components to deliver cheap computing cycles. The millions of servers operating in today's cloud data centers deliver the computing power necessary to solve problems that in the past could only be solved by large supercomputers assembled from expensive, one-of-a-kind components.

The cloud system software is complex. There is not a single *killer application* for cloud computing thus, investing in large-scale computing systems can only be justified if the cloud infrastructure can effectively accommodate a mix of applications. In addition to scalability challenges, modern cluster management systems address the problems posed by a mix of workloads. Typical cloud workloads include not only coarse-grained, batch applications, but also fine-grained, long-running applications with strict timing constraints. Only strict performance isolation and sophisticated scheduling can eliminate the undesirable effects of long-tail distribution of latency-sensitive jobs response time.

Resource virtualization is widely used to hide the complexity of a physical system and facilitate system access. Virtual machines (VMs) and containers are key components of the cloud infrastructure as discussed in depth in Chapter 5. A VM abstracts the hardware and *exploits virtualization by multiplexing* allowing multiple virtual systems to share a physical system. A container exploits *virtualization by aggregation* and bridges the gap between a clustered infrastructure and assumptions made by applications about their environments. Containers abstract an OS and include applications or tasks, as well as all their dependencies. Containers are portable, independent objects that can be easily manipulated by the software layer managing a large virtual computer. This virtual computer exposes to users the vast resources of a physical cluster with a very large number of independent processors.

The management of large-scale systems poses significant challenges and has triggered a flurry of developments in hardware and software systems. Several milestones in the evolution of ideas in cluster architecture, along with algorithms and policies for resource sharing and effective implementation of the mechanisms to enforce these policies, are analyzed in this chapter. Our analysis of the cloud software stack is complemented by the discussion of software systems closely related to applications presented in Chapter 11 and by topics related to Big Data applications discussed in Chapter 12.

Section 4.1 examines challenges faced by cloud infrastructure and the benefits of virtualization and containerization. The next two Sections, 4.2 and 4.3, analyze Warehouse Scale Computers (WSCs) and their performance.

Then the focus switches to software and presents VMs and hypervisors in Section 4.4 and frameworks, such as Dryad, Mesos, and Borg, in Sections 4.5, 4.6, and 4.7, respectively. Dryad is an execution

engine for coarse-grained data parallel applications, Mesos is used for fine-grained cluster resource sharing, and Borg is a cluster management system. Section 4.8 presents the evolution of cluster management systems and Section 4.9 covers Omega, a system based on state sharing. Section 4.10 analyzes Quasar, a system supporting QoS-aware cluster management. Resource isolation discussed in Section 4.11 is followed by an analysis of in-memory cluster computing with Spark and Tachyon in Section 4.12. Docker containers and Kubernetes are covered in Sections 4.13 and 4.14.

4.1 Cloud infrastructure challenges

Computing systems have evolved from single processors to multiprocessors, to multicore multiprocessors, and to clusters. The next step in this evolution, the Warehouse-scale computers (WSCs) with hundreds of thousands of processors are no longer a fiction, but serve millions of users, and are analyzed in computer architecture textbooks [52,232].

These systems are controlled by increasingly complex software stacks. Software helps integrate a very large number of system components and contributes to the challenge of ensuring efficient and reliable operation. The scale of the cloud infrastructure, combined with the relatively low mean-time to failure of the off-the-shelf components used to assemble a WSC, make the task of ensuring reliable services quite challenging.

At the same time, long-running cloud services require a very high degree of availability. For example, a 99% availability translates to 22 hours of downtime per quarter. Only a fair level of hardware redundancy combined with software support for error detection and recovery can ensure such a level of availability [232].

Virtualization. The goal of virtualization is to support portability, improve efficiency, increase reliability, and shield the user from infrastructure complexity. For example, threads are virtual processors, abstractions that allow a processor sharing among different activities thus, increase utilization and effectiveness. RAIDs are abstractions of storage devices designed to increase reliability and performance.

Processor virtualization, running multiple independent instances of one or more Operating Systems (OS), pioneered by IBM in early 1970s, was revived for computer clouds. Running multiple VMs on the same server enables applications to better share the server resources and lead to higher processor utilization. The instantaneous demands for resources by the applications running concurrently are likely to be different and complement each other, so the idle time of the server is reduced.

Processor virtualization is beneficial for both users and cloud service providers (CSPs). Cloud users appreciate virtualization because it allows better isolation of applications from one another than the traditional process-sharing model. Another advantage is that an application developer can chose to develop the application in a familiar environment and under the OS of choice. CSPs enjoy larger profits due to the lower cost for providing cloud services. Virtualization also provides more freedom for the system resource management because VMs can be easily migrated. VM migration proceeds as follows: the VM is stopped, its state is saved as a file, and the file is transported to another server where the VM is restarted. Nowadays, several hypervisors including KVM and Vmware ESXi fully support live migration, thus eliminating the need to stop the VM.

On the other hand, virtualization contributes to increased complexity of system software and has undesirable side effects on application performance and security. Processor sharing is now controlled by a new layer of software, the *hypervisor*, also called a Virtual Machine Monitor. It is often argued

that a hypervisor is a more compact software with only a few hundred thousand lines of code versus the million lines of code of a typical OS, thus it is less likely to be faulty.

Unfortunately, though the footprint of the hypervisor is small, the server must run a management OS. For example, Xen, the hypervisor used by AWS until 2017, initially starts Dom0, a privileged domain that starts and manages *DomU* unprivileged domains. *Dom0* runs the Xen management *toolstack*, is able to access the hardware directly, and provides Xen with virtual disks and network access for guests.

Virtual Machines run applications inside a guest OS which runs on virtual hardware under the control of a hypervisor. One can build immutable VM images, but VMs are heavyweight and non-portable.

Containers. Containers are based on *operating-system-level virtualization rather than hardware virtualization*; they isolate an application running inside a container from another application running in a different container and isolate them both from the physical system where they run. Moreover, containers are portable, and the resources used by a container can be limited. They are more transparent than VMs, thus easier to monitor and manage. Containers have several other benefits including:

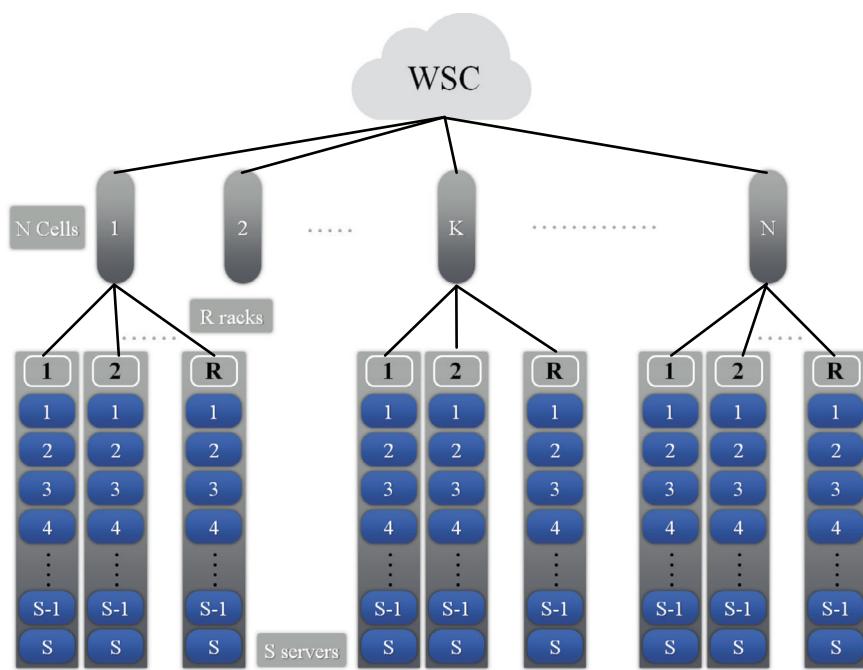
1. Ease creation and deployment of applications.
2. Applications are decoupled from the infrastructure; application container images are created at build time rather than deployment time.
3. Support portability; containers run independently of the environment.
4. Benefit from an application-centric management.
5. Have an optimal philosophy for application deployment; applications are broken into smaller, independent pieces and can be managed dynamically.
6. Support higher resource utilization.
7. Lead to predictable application performance.

Containers were initially designed to support the isolation of the *root* file system. The concept can be traced back to the *chroot* system call implemented in 1979 in Unix to change the root directory for the running process issuing the call and for its children, and to prohibit access to files outside the directory tree. Later, BSD and Linux adopted the concept, and in 2000, FreeBSD expanded it and introduced the *jail* command. The environment created with *chroot* was used to create and host a new virtualized copy of the software system.

Container technology has emerged as an ideal solution for application developers who no longer need to be aware of cluster organization and management details. Container technology is now ubiquitous and has a profound impact on cloud computing. Docker's containers gained widespread acceptance for ease of use, while Google's Kubernetes are performance-oriented.

Cluster management systems have evolved, and each system has benefited from the experience gathered from the previous generation. Mesos, a system developed at UC Berkeley, is now widely used by more than 50 organizations and has also morphed into a variety of systems, such as Aurora used by Twitter, Marathon offered by Mesosphere,¹ and Jarvis used by Apple. Borg, Omega, and Kubernetes are the milestones in Google's cluster management development effort discussed in this chapter.

¹ Mesosphere is a startup selling the Datacenter Operating System, a distributed OS, based on Apache Mesos.

**FIGURE 4.1**

Organization of a WSC with N cells, R racks, and S servers per rack.

4.2 Cloud hardware; warehouse-scale computer (WSC)

Cloud computing had an impact on large-scale systems architecture. WSCs [52,232] form the backbone of the cloud infrastructure of Google, Amazon, and other CSPs. WSCs are hierarchically organized systems with 50 000–100 000 processors capable of exploiting request-level and data-level parallelism.

At the heart of a WSC is a *hierarchy of networks* which connect the system components, servers, racks, and cells/arrays together, as in Fig. 4.1. Typically, a *rack* consists of 48 servers interconnected by a 48 port, 10 Gbps Ethernet (GE) switch. In addition to the 48 ports, the GE switch has two to eight uplink ports connecting a rack to a cell. Thus, the level of *oversubscription*, the ratio of internal to external ports, is between $48/8 = 6$ and $48/2 = 24$. This has serious implications for the performance of an application; two processes running on servers in the same rack have a much larger bandwidth and lower latency than the same processes running on servers in different racks.

The next component is a *cell*, sometimes called an *array*, consisting of a number of racks. The racks in a cell are connected by an *array switch*, a rather expensive communication hardware with a cost two orders of magnitude higher than that of a rack switch. The cost is justified as the bandwidth of a switch with n ports is of order n^2 . To support a 10 times larger bandwidth for 10 times as many ports, the cost increases by a factor of 10^2 . An array switch can support up to 30 racks.

Table 4.1 The memory hierarchy of a WSC with the latency given in microseconds, the bandwidth in MB/sec, and the capacity in GB [52].

Location type	DRAM			Disk		
	Latency	Bandwidth	Capacity	Latency	Bandwidth	Capacity
Local	0.1	20 000	16	10 000	200	2 000
Rack	100	100	1 040	11 000	100	160 000
Cell	3 000	10	31 200	12 000	10	4 800 000

WSCs support both interactive and batch workloads. The communication latency and the bandwidth within a server, a rack, and a cell are different, thus, the execution time and the costs for running an application is affected by the volume of data, the placement of data, and by the proximity of instances. For example, the latency, the bandwidth, and the capacity of the memory hierarchy of a WSC with 80 servers/rack and 30 racks/cell is shown in Table 4.1 based on the data from [52].

DRAM latency increases by more than three orders of magnitude, while the bandwidth decreases by a similar factor from server, to rack, and to cell. The latency and the bandwidth of the disks follow the same trend, but the variation is less dramatic. To put this into perspective, the memory-to-memory transfer of 1 000 MB takes 50 msec within a server, 10 seconds within the rack, and 100 seconds within a cell, while disk transfers take 5, 10, and 100 seconds, respectively.

WSCs are expected to supply cheap computing cycles, but are expensive; the cost of a WSC is of the order of \$150 million, but the cost-performance is what makes them appealing. The capital expenditures for a WSC includes the costs for servers, for the interconnect and for the facility. A case study reported in [232] shows a capital expenditure of \$167 510 000 including \$6 670 000 for 45 978 servers, \$12 810 000 for an interconnect with 1 150 rack switches, 22 cell switches, 2 layer 3 switches, and 2 border routers. In addition to the initial investment, the operation cost of the cloud infrastructure including the cost of energy is significant. In this case study the facility is expected to use 8 MW.

We now take a closer look at the WSC servers and ask what type of processors are best suited as server components. Multicore processors are ideal components of WSC servers because they support not only data-level parallelism for search and analysis of very large data sets but also request-level parallelism for systems expected to support a very large number of transactions per second. *Data-parallel* and *request-parallel* applications are the two major components of the workloads experienced by cloud service providers such as Google.

There are two basic types of multicore processors often called *browny* and *wimpy* cores [242]. The single-core performance of a browny core is impressive, but so is its power dissipation. The wimpy cores are less powerful but consume less power. Power consumption is a major concern for cloud as we see in Section 1.2; for solid-state technologies the power dissipation is approximately $\mathcal{O}(f^2)$ with f being the clock frequency.

An application task needs to spawn a larger number of threads when running on wimpy rather than browny cores, and this has two major implications: First, it complicates the software development process because it requires an explicit parallelization of the application, thus increases the application development cost. The second, equally important implication is that running a larger number of threads increases the response time. Multi-phase algorithms use *barrier-synchronization* requiring all threads have to finish before the next phase of the computation. This means that all threads have to wait for the slowest one.

The cost of systems using wimpy core may increase, e.g., the cost for DRAM will increase as the kernel and system processes consume more aggregate memory. Data structures used by applications might need to be loaded into memory on multiple wimpy-core machines instead of a single brawny-core machine with negative effects on performance. Lastly, managing a larger number of threads will increase the scheduling overhead and diminish performance.

Hölzle [242] concludes “Once a chip’s single-core performance lags by more than a factor of two or so behind the higher end of current-generation commodity processors, making a business case for switching to the wimpy system becomes increasingly difficult because application programmers will see it as a significant performance regression: their single-threaded request handlers are no longer fast enough to meet latency targets.”

4.3 WSC performance

The central questions discussed now are: how to extract the maximum performance from a warehouse-scale computer; what are the main sources of WSC inefficiency; and how these inefficiencies could be avoided. Even slight WSC performance improvements translate into large cost savings for the CSPs and noticeably better service for cloud users.

WSCs workload is very diverse; there are no typical, or “killer,” applications that would drive the design decisions and, at the same time, guarantee optimal performance for such workloads. It is not feasible to experiment with systems at this scale or to simulate them effectively under realistic workloads. The only alternative is to profile realistic workloads and analyze carefully the data collected during production runs, but this is only possible if low-overhead monitoring tools that minimize intrusion on the workloads are available. Monitoring tools minimize intrusion by random sampling and by maintaining counters of relevant events, rather than detailed event records.

Google-Wide-Proiling (GWP) is a low-overhead monitoring tool used to gather the data through random sampling. GWP randomly selects a set of servers to profile every day, uses mostly Perf² to monitor their activity for relatively short periods of time, collects the *callstacks*³ of the samples, aggregates the data, and then stores it in a database [413].

Data collected at Google with GWP over a period of 36 months is presented in [265] and discussed in this section. Only data for C++ codes were analyzed because C++ codes dominate the CPU cycle consumption, though the majority of codes are written in Java, Python, and Go. The data was collected from some 20 000 servers built with Intel Ivy Bridge processors.⁴

The analysis is restricted to 12 application binaries with distinct execution profiles: batch versus latency sensitive and low-level versus high-level services. The applications are: Gmail and Gmail-fe, the back- and front-end Gmail application; BigTable, a storage system discussed Section 7.11; *disk*, low-level distributed storage driver; *indexing1* and *indexing2* of the indexing pipeline; *search1*, *2*, *3*, application for searching leaf nodes; *ads*, an application targeting ads based on web-page contents;

² Perf is a profiler tool for Linux 2.6+ systems; it abstracts CPU hardware differences in Linux performance measurements.

³ A *callstack*, also called execution stack, program stack, control stack, or run-time stack, is a data structure that stores information about the active subprograms invoked during the execution of a program.

⁴ The Ivy Bridge-E family is made in three different versions with the postfix -E, -EN, and - EP, with up to six, ten, and twelve cores per chip, respectively.

video, a transcoding and feature extraction application; and *flights-search*, the application used to search and price flights.

In spite of the workload diversity, there are common procedures used by a vast majority of applications running on WSCs. Data-intensive applications run multiple tasks distributed across several servers, and these tasks communicate frequently with one another. Cluster management software is also distributed and daemons running on every node of the cluster communicate with one or more schedulers making system-wide decisions. It is not unexpected that common procedures accounting for a significant percentage of CPU cycles are dedicated to communication. This is the case of RPCs (Remote Procedure Calls), as well as serialization, deserialization, and compression of buffers used by communication protocols.

A typical communication pattern involves the following steps: (a) serialize the data to the protocol buffer; (b) execute an RPC and pass the buffer to the callee; and (c) caller deserializes the buffers received in response to the RPC. Data collected over a period of 11 months shows that these common procedures translate into an unavoidable “WSC architecture tax” and consume 22–27% of the CPU cycles. About one-third of RPCs are used by cluster management software to balance the load, encrypt the data, and detect failures. The remaining RPCs move data between application tasks and system procedures. Data movement is also done using library functions, such as *memmove* and *memcpy*, with descriptive names⁵ which account for 4–5% of the “tax.”

Compression, decompression, hashing, and memory allocation and reallocation procedures are also common and account for more than one fourth of this “tax.” Applications spend about one fifth of their CPU cycles in the scheduler and the other components of the kernel. An optimization effort focused on these common procedures will undoubtedly lead to a better WSC utilization.

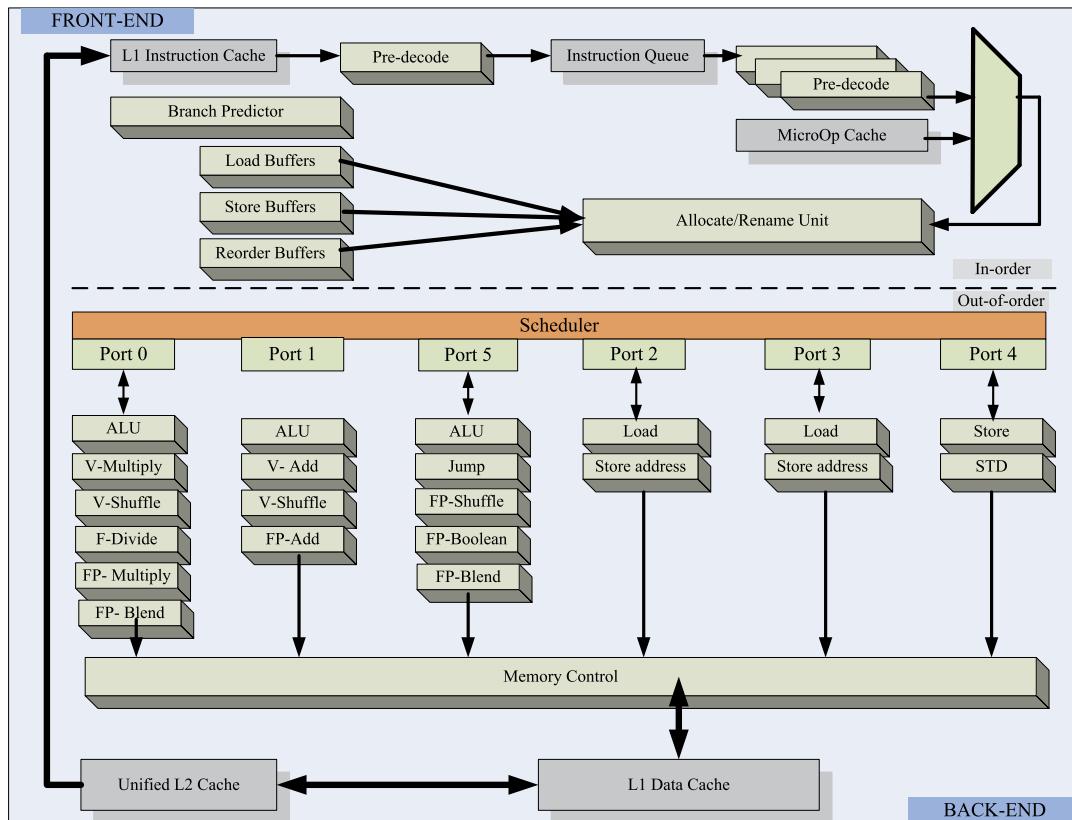
Most cloud applications have a substantial memory footprint, binaries of hundreds of MB are not uncommon, and some do not exhibit either spatial or temporal locality. Moreover, the memory footprint of applications shows a significant rate of increase, about 30% per year. Also, the instruction-cache footprints grow at a rate of some 2% per year. As more Big Data applications run on computer clouds, neither the footprint nor the locality of these applications are likely to limit the pressure on cache and memory management. These functions represent a second important target for the performance optimization effort.

Memory latency rather than memory bandwidth affects a processor ability to deliver a higher level of performance through Instruction Level Parallelism (ILP). The performance of such processors is significantly affected by stall cycles due to cache misses. It is reported that data cache misses are responsible for 50–60% of the stall cycles and, together with instruction cache misses, contribute to a lower IPC (Instructions Per Clock cycles).

There are patterns of software that hinder execution optimization through pipelining, hardware threading, out-of-order execution, and other architectural features designed to increase the level of ILP. For example, linked data structures cause indirect addressing that can defeat hardware prefetching and build bursts of pipeline idleness when no other instructions are ready to execute.

Understanding cache misses and stalls requires a microarchitecture-level analysis. A generic organization of the microarchitecture of a modern core is illustrated in Fig. 4.2. The core front-end processes

⁵ In Linux, *memmove* and *memcpy* copy n bytes from one memory area to another; the areas may overlap for the former but do not overlap for the latter.

**FIGURE 4.2**

Schematic illustration of a modern processor's core microarchitecture. Shown are the front-end, the back-end, L1 instruction and data caches, the unified L2 cache, and the microoperation cache. Branch prediction unit, load/store, and reorder buffers are components of the front-end. The instruction scheduler manages the dynamic instruction execution. The five ports of the instruction scheduler dispatch microinstructions to ALU and to load and store units. Vector (V) and floating-point operations (FP) are dispatched to ALU units.

instructions in order, while the instruction scheduler of the back-end is responsible for dynamic instruction scheduling and feeds instructions to multiple execution units including Arithmetic and Logic Units (ALU)s and Load/Store units.

The microarchitecture-level analysis is based on a top-down methodology [527]. According to <https://software.intel.com/en-us/top-down-microarchitecture-analysis-method>: “The Top-Down characterization is a hierarchical organization of event-based metrics that identifies the dominant performance bottlenecks in an application. Its aim is to show, on average, how well the CPU’s pipeline(s) were being utilized while running an application.”

This methodology identifies the *micro-op* (μ op) queue as the separator between the front-end and the back-end components of a microprocessor core. The μ op pipeline slots are then classified as *Retiring*, *Front-end bound*, *Bad speculation*, or *Back-end bound*, with only the first one doing useful work. *Front-end bound* includes overheads associated with fetching, instruction caches, decoding, and some other shorter-penalties and *Back-end bound* includes overheads due to the data-cache hierarchy and the lack of ILP; *Bad speculation* is self-explanatory.

In a typical SPEC CPU2006 benchmark, the front-end wasted execution slots are typically two–three times lower than those reported for the Google workload which account for 15–30% of all wasted slots. A reason for this behavior is that SPEC applications⁶ do not exhibit the combination of low retirement rates⁷ and high front-end boundedness of WSCs.

Data shows that the core back-end dominates the overhead and limit the ILP. The back-end and the front-end stalls limit the number of cores active during an execution cycle. To be more precise, only one or two cores of a six-core Ivy Bridge processor are active in 72% of execution cycles, while three cores are active during the balance of 28% of cycles.

The observation that memory latency is more important than memory bandwidth is a consequence of the low memory-bandwidth utilization at an average of 31% and a maximum of 68% with a heavy tail distribution. In turn, the low memory utilization is due in part to low CPU utilization. A surprising result reported in [265] is that the median CPU bandwidth utilization is 10%, while [52] reports a median CPU utilization in a much higher range, 40–70%. A low CPU utilization is also reported for the CloudSuite [170].

Several conclusions regarding optimal processor architecture can be reached from the top-down data analysis. Data analysis shows that cloud workloads display access patterns involving bursts of computations intermixed with bursts of stall cycles. Processors supporting a higher level of *simultaneous multithreading* (SMT) are better equipped to hide the latency by overlapping stall cycles than two-wide SMP processors. SMT is an architectural feature allowing instructions from more than one thread to be executed in any given pipeline stage at a time. SMT requires the ability to fetch instructions from multiple threads in a cycle; it also requires a larger register file to hold data from multiple threads.

The large working sets of the codes are responsible for the high rate of instruction cache misses. L2 caches show that MPKI (misses per kilo instructions) are particularly high. Larger caches would alleviate this problem but at the cost of higher cache latency. Separate cache policies that give priority to instructions over data or separate L2 caches for instructions and data could help in this regard.

⁶ The following applications in the SPEC CPU2006 suite are used: *400.perlbench* application which has high IPC and the largest instruction cache working set; *445.gobmk*, an application with hard-to-predict branches; *429.mcf* and *471.omnetpp* memory-bound applications which stress memory latency; and *433.milc*, a memory-bound application which stresses memory bandwidth.

⁷ In a modern processor, the Completed Instruction Buffer holds instructions that have been speculatively executed. Associated with each executed instruction in the buffer are its results in rename registers and any exception flags. The retire unit removes these executed instructions from the buffer in program order at a rate of up to four instructions per cycle. The retire unit updates the architected registers with the computed results from the rename registers. The retirement rate measures the rate of these updates.

4.4 Hypervisors

A hypervisor securely partitions resources of a computer system into one or more VMs. A *guest OS* is an operating system that runs under the control of a hypervisor rather than directly on the hardware. A hypervisor runs in kernel mode, while a guest OS runs in user mode; sometimes, the hardware supports a third mode of execution for the guest OS. Hypervisors allow several operating systems to run concurrently on a single hardware platform and enforce isolation among these systems, thus better security. A hypervisor controls how the guest OS uses the hardware resources; the events occurring in one VM do not affect any other VM running under the same hypervisor. At the same time, the hypervisor enables:

- Multiple services to share the same platform.
- The movement of a service from one platform to another called *live migration*.
- System modification while maintaining backward compatibility with the original system.

When a guest OS attempts to execute a privileged instruction, the hypervisor traps the operation and enforces the correctness and safety of the operation. The hypervisor guarantees the isolation of the individual VMs and, thus, ensures security and encapsulation, a major concern in cloud computing. At the same time, the hypervisor monitors the system performance and takes corrective actions to avoid performance degradation. For example, the hypervisor may swap out a Virtual Machine to avoid thrashing; to do so, the hypervisor copies all pages of the evicted VM from real memory to disk and makes the real memory frames available for paging by other VMs.

A hypervisor virtualizes the CPU and the memory. For example, the hypervisor traps interrupts and dispatches them to the individual guest operating systems; if a guest OS disables interrupts, the hypervisor buffers such interrupts until the guest OS enables them. The hypervisor maintains a *shadow page table* for each guest OS and replicates any modification made by the guest OS in its own shadow page table. This shadow page table points to the actual page frame, and it is used by the hardware component called the Memory Management Unit (MMU) for dynamic address translation.

Memory virtualization has important implications on the performance. Hypervisors use a range of optimization techniques; for example, the VMware systems avoid page duplication among different VMs, they maintain only one copy of a shared page, and use copy-on-write policies, while Xen imposes total isolation of the VM and does not allow page sharing. Hypervisors control the virtual memory management and decide what pages to swap out; for example, when the ESX VMware Server wants to swap out pages, it uses a *balloon process* inside a guest OS and requests it to allocate more pages to itself and, thus, swaps-out pages of some of the processes running under that VM. Then it forces the balloon process to relinquish control of the free page frames.

4.5 Execution of coarse-grained data-parallel applications

The distinction between fine-grained and coarse-grained parallelism introduced in Chapter 3 is important for understanding cloud software organization. Application developers have used the SPMD (Same-Program-Multiple-Data) paradigm for several decades for exploiting coarse-grained parallelism. The name SPMD illustrates perfectly the idea behind the concept—a large dataset is split into several segments processed independently, and often concurrently, using the same program.

For example, converting a large number of images, e.g., 10^9 , from one format to another, can be done by splitting the set into 1 000 thousand segments with 10^6 images each, and then running concurrently the conversion program on 1 000 processors. In this example, 1 000 processors cut the computing time by almost three orders of magnitude.

It is easy to see that such applications are ideal for cloud computing since they need a large computing infrastructure and can keep the systems busy for a fair amount of time. To CSP's delight, such jobs increase system utilization because no scheduling decisions have to be made until the time-consuming job has finished. This was noticed early on and, in 2004, the MapReduce idea was born [129].

MapReduce and Apache Hadoop, an open-source software framework consisting of a storage part, the Hadoop Distributed File System (HDFS), and the processing part called MapReduce discussed in Chapter 11 expand on the SPMD concept. MapReduce is a two-phase process. Input data is first segmented, then the computations on data segments are carried out during the Map phase; partial results are then merged during the Reduce phase.

This extends the scope of SPMD for computations that are not totally independent, the so-called *embarrassingly parallel* applications. In the previous example, instead of converting the 10^9 images, we now search during the Map phase for an object that could appear in any of them; after the search in each segment is completed, during the Reduce phase, we combine the partial results and further refine the information about the object from the set of images selected during the Map phase.

Dryad is a general-purpose distributed execution engine developed in 2007 by Microsoft for coarse-grained data-parallel applications. Microsoft wanted to use Dryad for running Big Data applications on its clustered server environment as a proprietary alternative to Hadoop, a widely used platform for coarse-grained data-parallel applications. A Dryad application combines computational *vertices* with communication links to form a dataflow graph [255]. Then, it runs the application by executing the vertices of this graph on a set of available computers, communicating through files, TCP pipes, and shared-memory.

The system is centrally controlled by a *Job Manager* (JM) running either on one of the nodes of the cluster or outside the cluster, on the user's computer. The JM uses a *Name Server* (NS) to locate the nodes of the cluster where the work is actually done. The JM uses an application-specific description to construct the dataflow graph of the application. A *daemon* running on each cluster node communicates with the JM and controls the execution of the code for the vertex of the graph assigned to that node. Daemons communicate directly among themselves without the intervention of the JM and use the information provided by the dataflow graph to carry out the computations.

A detailed description of the Dryad dataflow graph given in [255] presents a set of simpler graphs used to construct more complex one. The graph nodes are annotated to show the input and output datasets. Two connection operations are the point-wise composition and the complete bipartite composition.

The DryadLINQ (DLNQ) system is the product of a related Microsoft project [530]. It exploits the Language INtegrated Query (LINQ), a set of .NET constructs for performing arbitrary side effect-free transformations on datasets to automatically translate data-parallel codes into an workflow for distributed execution. The distributed execution is then executed on the Dryad platform. The development of DryadLINQ was motivated by the fact that parallel databases implement only declarative sides of SQL queries and do not support imperative programming. Dryad is not scalable and to no one's surprise, soon after announcing plans to release Windows Azure- and Windows Server-based implementations of open source Apache Hadoop, Microsoft discontinued the project.

4.6 Fine-grained cluster resource sharing in Mesos

Mesos is a light-weight framework for fine-grained cluster resource sharing developed in late 2010s at UC Berkeley. Mesos consists of only some 10 000 lines of C++ code [240]. Mesos runs on Linux, Solaris and OS X and supports frameworks written in C++, Java, and Python. Mesos can use Linux containers to isolate task CPU cores and memory.

A novelty of the system is a two-level scheduling strategy for large clusters with workloads consisting of a mix of frameworks. The term “framework” in this context means a large consumer of CPU cycles and widely-used software systems such as Hadoop, discussed in Chapter 11, and MPI (Message Passing Interface), a standardized and portable message-passing system used by the parallel-computing community since the 1990s. Another novelty is the concept of *resource offer*, an abstraction for a bundle of resources that a framework can allocate on a cluster node to run its tasks.

The motivation for Mesos is that *centralized scheduling is not scalable due to its complexity* and cannot perform well for fine-grained resource sharing. Framework jobs consisting of short tasks are mapped to resource *slots*, and fine-grained matching has a high overhead and prevents sharing across frameworks. Mesos allows multiple frameworks to share resources in a fine-grained manner and achieve data locality. It can isolate a production framework from experiments with a new version undergoing testing and experimentation. Each framework has a *scheduler* that receives resource offers from the master and an *executor* on each machine to launch the tasks of the framework. Scheduler functions are:

- *callbacks*,⁸ such as *resourceOffer*, *offerRescinded*, *statusUpdate*, and *slaveLost*, and
- *actions*, such as *replyToOffer*, *setNeedsOffers*, *setFilters*, *getGuaranteedShare*, and *killTask*.

The executer functions are *callbacks*, such as *launchTask* and *killTask*, and *actions*, such as *sendStatus*.

Framework requests differentiate mandatory from preferred resources. A resource is *mandatory* if the framework cannot run without it, e.g., a GPU is a mandatory resource for applications using CUDA.⁹ A resource is *preferred* if a framework performs better using a certain resource but could also run using another one.

This two-level scheduling strategy keeps Mesos simple and scalable and, at the same time, gives the frameworks the power to optimally manage a cluster. The system is flexible and supports pluggable *isolation modules* to limit CPU, memory, and network bandwidth, and the I/O usage of a process tree. Allocation modules can select framework-specific policies for resource management. For example, the killing of a task of a greedy or buggy framework is aware whether the tasks of a framework are interdependent, as in case of MPI, or independent, as in case of MapReduce.

Mesos is organized as follows: A master process manages daemons running on all cluster nodes, while frameworks run the tasks on cluster nodes. The master implements the fair-sharing of resource offers among frameworks and lets each framework manage resource sharing among its tasks. The system is robust, i.e., there are replicated masters in hot-standby state. When the active master fails, a ZooKeeper service¹⁰ is used to elect a new master. Then, the daemons and the schedulers reconnect to the newly elected master.

⁸ A callback is executable code passed as an argument to other code; the callee is expected to execute the argument either immediately or at a later time for synchronous and, respectively, asynchronous callbacks.

⁹ CUDA is a parallel-computing platform supporting graphics processing units (GPUs) for general-purpose processing.

¹⁰ ZooKeeper is a distributed coordination service implementing a version of the Paxos consensus algorithm, see Chapter 11.

The limitations of the distributed scheduling implemented by Mesos are also discussed in [240]. Sometimes, the collection of frameworks is not able to optimize bin packing and a centralized scheduler. Another type of fragmentation occurs when the tasks of a framework request relatively small quantities of resources and, upon completion, resources released by the tasks are insufficient to meet the demands for large quantities of resources by tasks from a different framework. Resource offers could increase the complexity of framework scheduling, but centralized scheduling is not immune to this problem, either.

Porting Hadoop to run on Mesos required relatively few modifications. Indeed, the *JobTracker* and the *TaskTrackers* components of Hadoop map naturally as Mesos framework scheduler and executor, respectively. Apache Mesos is an open-source system adopted by some 50 organizations including Twitter, Airbnb, and Apple (see <http://mesos.apache.org/documentation/latest/powerd-by-mesos/>).

Several frameworks based on Mesos have been developed over the years. Apache Aurora was developed in 2010 by Twitter and now is open-source. Chronos is a cron-like¹¹ system, elastic and able to express dependencies between jobs. Apple uses a Mesos framework called Jarvis¹² to support Siri. Jarvis is an internal PaaS cloud service to answer IOS user's voice queries. A fair scheduler at Facebook allocates the resources of a 2 000 node cluster dedicated to Hadoop jobs. MPI and MapReduceOnline (a streaming of MapReduce discussed in Chapter 11) jobs for ad targeting need the data stored on the Hadoop cluster, but the frameworks cannot be mixed.

The utilization analysis of a production cluster with several thousand servers used to run production jobs at Twitter shows that the average CPU utilization is below 20% and the average memory utilization is below 50% [240]. Reservations improve the CPU utilization up to 80%. Some 70% of the reservations overestimate the resources they need by one order of magnitude, while 20% of the reservations underestimate resource needs by a factor of five.

In summary, *one can view Mesos as the opposite of virtualization. A VM is based on an abstraction layer encapsulating an OS together with an application inside a physical machine. Mesos abstracts physical machines as pools of indistinguishable servers and allows a controlled and redundant distribution of tasks to these servers.*

4.7 Cluster management with Borg

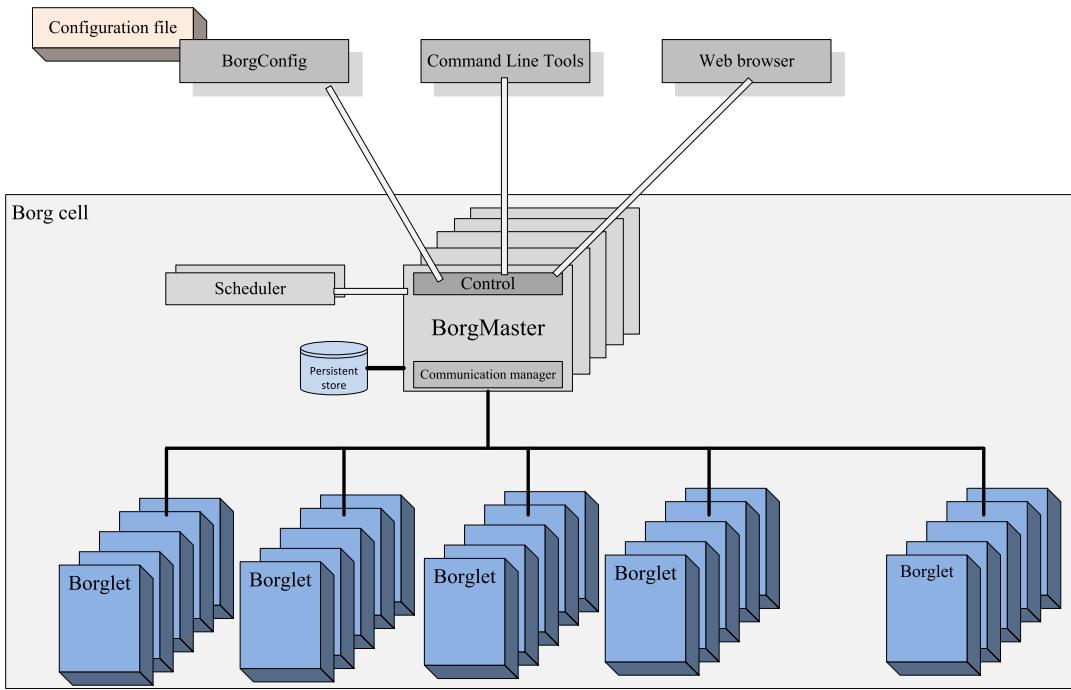
A computer cluster may consist of tens of thousands of processors. For example, a cell of a WSC is in fact a cluster consisting of multiple racks, each with tens of processors, as shown in Fig. 4.1. There are two sides of cluster management; One reflects the views of application developers who need simple means to locate resources for an application and then to control the use of resources; the other is the view of service providers concerned with system availability, reliability, and resource utilization.

These views drove the design of Borg, a cluster management software developed at Google [495]. Borg's design goals were:

- Manage effectively workloads distributed to a large number of systems; be highly reliable and available.

¹¹ Cron is a job scheduler for Unix-like systems used to periodically schedule jobs; often, it is used to automate system maintenance and administration.

¹² Jarvis is short for *Just A Rather Very Intelligent Scheduler*.

**FIGURE 4.3**

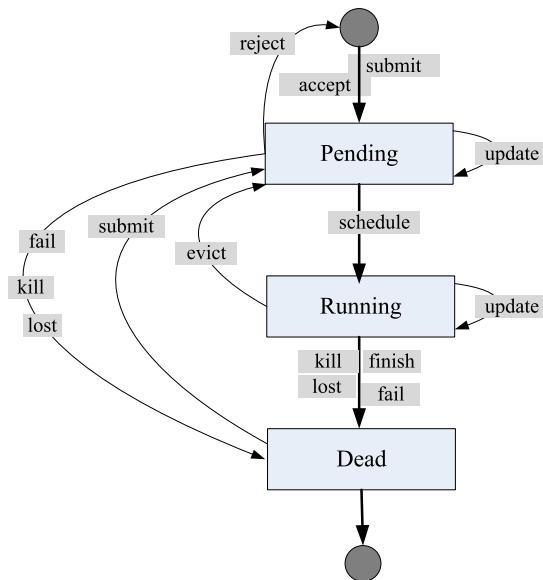
Borg architecture. A replicated *BorgMaster* interacts with *Borglets* running on each machine in the cell.

- Hide the details of resource management and failure handling and allow users to focus on application development. This requirement is important because the servers of a cluster differ in terms of processor type and performance, number of cores per processor, RAM size, secondary storage, network interface, and other capabilities.
- Support long-running, highly-dependable applications, including interactive production jobs and batch non-production jobs.

A Borg cluster consists of tens of thousands of machines co-located and interconnected by a data center-scale network fabric. A cluster managed by Borg is called a *cell*. The architecture of the system shown in Fig. 4.3 consists of a logically centralized controller, the *BorgMaster*, and a set of processes running on each machine in the cell, the *Borglets*. All Borg components are written in C++.

The main *BorgMaster* has five replicas, and each replica maintains an in-memory copy of the state of the cell. The state of a cell is also recorded in a Paxos-based store on local disks of each replica. An elected master serves as the Paxos leader and handles operations that change the state of a cell, e.g., submit a job or terminate a task.

The master process performs several actions: (i) handles client RPCs that change state or lookups data; (ii) manages state machines for machines, tasks, allocs, and other system objects; (iii) communicates with *Borglets*; and (iv) responds to requests submitted to a web-based user interface. *Borglets*

**FIGURE 4.4**

The states of a Borg task. Task state changes as a result of either user requests or system actions.

start, stop, and restart failing tasks, manipulate the OS kernel setting to manage local resources, and report the local state to the *BorgMaster*.

Users interact with running processes by means of RPCs to the *BorgMaster* and trigger task state transitions. Users request actions such as *submit*, *kill*, and *update*. The task state is also changed by system actions such as: *reject*, *evict*, and *lost*, see Fig. 4.4.

The *scheduler* is the other main component of the *BorgMaster*. The scheduler scans periodically in a round-robin order a priority queue of pending tasks. The *feasibility* component of the scheduling algorithm attempts to locate systems where to run tasks. The *scoring* component identifies the machine(s) to actually run the task.

Alloc and *alloc sets* reserve resources on a machine and, respectively, on multiple machines. Jobs have priorities: Distinct *priority bands* are defined for activities such as monitoring, production, batch, and testing. A *quota system* for job scheduling uses a vector including the quantity of resources such as CPU, RAM, and disk for specified periods of time. Higher-priority quota cost more than lower-priority ones. To simplify resource management and balance the load, a system similar with the one described in [22] is used to generate a single cost value per vector and minimize the cost, while balancing the workload and leaving room for spikes in demand.

Production jobs are allocated about 70% of CPU resources and 55% of the total memory [495]. A Borg job could have multiple tasks and runs in a single cell. The majority of jobs do not run inside a VM. Tasks map to Linux processes running in containers.

To manage large cells, the scheduler spawns several concurrent processes to interact with the *BorgMaster* and *Borglets*. These processes operate on cached copies of the cell state. Several other

design decisions are important for system scalability. For example, to avoid frequent time-consuming machine and task scoring, Borg caches the score until the properties of the machine or task change significantly.

To avoid determining the feasibility of each pending task on every machine, Borg computes feasibility and scoring per equivalence classes of tasks with similar requirements. Moreover, this evaluation is not done for every machine in the cell but on random machines until enough suitable machines have been found.

The state of the *BorgMaster* can be saved as a *checkpoint* file and later used for studies of system performance and effectiveness, or to restore the state to an earlier point in time. A simulator, the *Faux-Master*, designed to identify system errors and performance problems by replaying checkpoint files, facilitated the effort to improve the Borg system.

Results collected for a 12 000-server cluster at Google show an aggregate CPU utilization of 25–35% and an aggregate memory utilization of 40%. A reservation system raises these figures to 75% and 60%, respectively [412].

4.8 Evolution of a cluster management system

It is interesting to examine the evolution of cluster management systems over the years and see if cloud service providers, such as Amazon, Google, and Microsoft, have optimized their scheduling strategies to increase the resource utilization and efficiency of large data centers. Such studies are possible only if trace data covering extended periods of time are made available to researchers from industry and academia. Fortunately, in 2011, Google published a one-month trace from a cluster managed with Borg [411]. Eight years later, in May 2019, Google published again detailed job scheduling information from eight Borg cells [510], enabling a group of researchers to draw some conclusions regarding changes in workloads, scheduling strategies, and resource utilization.

The analysis in [476] starts with a comparison of hardware used by the two traces, 2019 versus 2011: eight versus one cell, 96 400 servers versus 12 600, 7 versus 3 hardware platforms, 21 versus 10 server configurations, and about the same number of servers per cell, 12 000 versus 12 000. As we can see, trace data collected in 2019 involves a considerably larger number of more diverse servers. The software is also different; it provides schedulers with a wealth of information to optimize resource utilization and better support high priority jobs. The 2019 Borg supports alloc sets, i.e., resources reserved for a job and distributed across multiple machines, job dependencies, batch queuing, and vertical scaling; the system also allows a much larger range of priority values 0–450 instead of 0–11. Vertical scaling is now supported by the Autopilot discussed in [426].

2018 workloads differ from those for 2011. The mean job-arrival rate increased from 964 to 3 360 jobs per hour, and the job submission rate was 3.7 times higher. The median task-scheduling rate increased 3.6 times; many task scheduling events are for rescheduling indicating that there is more “churn” in the system, i.e., more task completion events.

The 2019 scheduler is more agile and has to work harder, the number of jobs to be scheduled is seven times higher, and most of the workload has moved into the high-priority best-effort batch tier. The heavy-tailed distribution of job sizes and the power law Pareto distribution with heavy tail of computing cycles and memory add to the scheduling challenges of 2019 Borg. Significantly improved scheduling decisions led to an increase of both average CPU and memory utilization: 20–40% for CPU

and 3–30% for memory. The utilization variance is lower, and there were fewer servers with utilization larger than 80%. Median scheduling delays decreased, but the tail is longer for the last 28% of jobs.

The large variability of resource consumption is reflected by the spread of the coefficient of variation C^2 of two measures, Normalized Compute Unit-hours (NCU-hours) and Normalized Memory Unit-hours (NMU-hours) used per job. $C^2 = 23\,000$ for compute consumption because the variance is about 33 300 with a mean of 1.2 NCU-hours and $C^2 = 43\,000$ for memory consumption. Both values are extremely large, $C^2 = \text{variance}/\text{mean}^2$ is invariant to data normalization and $C^2 = 1$ for an exponential distribution, and a study of workloads in several supercomputing centers found that C^2 ranged from 28 to 256, depending on the supercomputing center. Trace data analysis shows that compute consumption and memory consumption follow almost the same distribution, and the two metrics are correlated [476]. One of the conclusions of this research is that resource-prediction algorithms appear more efficient than users manually defining their own resource requirements.

4.9 Shared state cluster management

Could multiple independent schedulers do a better cluster management job than monolithic or two-level schedulers? The designers of the Omega system understood that efficient scheduling of large clusters is a very hard problem due to the scale of the system, combined with the workload diversity, and that only a novel approach should be considered [442].

The workload of Google systems targeted by Omega was the mix of production/service and batch jobs discussed in Section 4.7. More than 80% of the workload are short batch jobs, spawning a large number of tasks. A larger share of resources, 55–80%, is allocated to production jobs running for extended periods of time and with fewer tasks than the batch jobs. The scheduling requirements are: short turnaround time for batch jobs and strict availability and performance targets for production jobs.

The scheduler workload increases with cluster size and with the granularity of the tasks to be scheduled. The finer the task granularity, the more scheduler decisions have to be made; thus the likelihood of spatial and temporal resource fragmentation and lower resource utilization is higher. The solution adopted in Omega's design is to allow multiple independent schedulers to access a *shared cluster state* protected with a lock-free optimistic concurrency control algorithm.

In Omega, there is no central resource allocator, and each scheduler has access to all cluster resources. A scheduler has its own private and frequently updated copy of the *shared cluster state*, a resilient master copy of the state of all cluster resources. Whenever it makes a resource allocation decision, a scheduler updates the shared cluster state in an *atomic transaction*.

In case of conflict among tasks, at most one commit succeeds, and the resource is allocated to the winner. Then, the shared cluster state re-syncs with local copies of all schedulers. The losers may then retry at a later time to gain access to the resource and can be successful after the resource has been released by the task holding it.

Multiple schedulers may attempt to allocate the same resource at the same time, thus there is the possibility of conflict. An optimal solution to this problem depends upon the frequency of conflicts. An optimistic approach increases parallelism and assumes that conflicts seldom occur and, when detected, they can be resolved efficiently. A pessimistic approach used by Mesos is to ensure that a resource is available to only one framework scheduler at a time.

Table 4.2 A side-by-side comparison of four types of schedulers. The schedulers differ in terms of: (a) scope, the set of resources controlled; (b) possibility of conflict and conflict resolution method; (c) allocation granularity, gang scheduling versus task-by-task; and (d) scheduling policy.

Scheduler	Resources	Conflict	Granularity	Policy
Monolithic e.g., Borg	All available resources	None	Global	Priority & preemption
Static partition e.g., Dryad	Fixed subset of resources	None	Per-partition policy	Scheduler-dependent
Two-level e.g., Mesos	Dynamic subset of resources	Pessimistic	Gang scheduling	Strict fairness
Shared-state e.g., Omega	All available	Optimistic	Per-scheduler policy	Priority & preemption

Jobs typically spawn multiple tasks, and the next question is whether all tasks of a job should be allocated all the resources they need at the same time when the job starts execution, a strategy called co-scheduling or gang scheduling. The alternative is to allocate resources only at the time when a task needs them. In the former case, resources end up being idle until the tasks actually need them, and the average resource utilization decreases. In the latter case, there is a chance of deadlock as some tasks need resources allocated to other tasks, while those holding these resources need the resources held by the first group of tasks.

Several metrics are useful to compare the effectiveness of large cluster schedulers. If we view scheduling as a service and a scheduling request as a transaction, then the time elapsed from the instance a job is submitted until the scheduler attempts to schedule the job is the waiting time; the time required to schedule a job is the service time. The waiting time and the service time are two important metrics for scheduler effectiveness. Conflict resolution is a component of the scheduler service time for a shared-state scheduler like Omega. The *conflict fraction* is the average number of conflicts per successful transaction.

The service time has two components, t_{sch_job} , the overhead for scheduling a job, and t_{sch_task} , the time to schedule a task of the job; thus, the total time to make a scheduling decision for a job with n tasks is:

$$t_{sch} = t_{wait} + t_{sch_job} + n \times t_{sch_tks}. \quad (4.1)$$

A monolithic scheduler using a centralized scheduling algorithm does not scale up. Another solution is to statically partition a cluster into sub-clusters allocated to different types of workloads. This policy is far from optimal due to fragmentation because the balance between different types of workload changes over time. A two-level scheduler has its own limitations, as we have seen in Section 4.6. Table 4.2 provides a side-by-side comparison of monolithic schedulers, schedulers for static-partitioned clusters, two-level schedulers, and multiple, independent, shared-state schedulers.

A light weight simulator was used for a comparative study of Omega and the other schedulers. The two-level scheduling model in this simulator emulates Mesos and achieves fairness by alternately offering all available cluster resources to different schedulers. It assumes low-intensity tasks, thus resources become available frequently and scheduler decisions are quick. As a result, a long scheduler decision time means that subsets of cluster resources are unavailable to other schedulers for extended periods of

time. The simulation results show that Omega is scalable, and at realistic workloads, there is little interference among independent schedulers. The simulator was also used to investigate gang scheduling [34] useful for MapReduce applications.

A more accurate, trace-driven scheduler can be used to gain further insights into scheduler conflicts. Trace-driven simulation is quite challenging; moreover, a large number of simplifying assumptions limit the accuracy of the simulator results. Neither the machine failures, nor the disparity between resource requests and the actual usage of those resources in the traces, are simulated by the trace-driven simulator designed for Omega. The results produced by the trace-driven simulator are consistent with the ones provided by the light-weight simulator.

4.10 QoS-aware cluster management

Quality of Service (QoS) guarantees are important for the designers of cloud applications and for the users of computer clouds who wish to enforce a well-defined range of response time, execution time, or other significant performance metrics for their cloud workloads. Enforcing workload constraints is far from trivial, ergo few cluster management systems could legitimately claim adequate QoS support.

Two aspects of cluster management, resource allocation and resource assignment, play a key role for supporting QoS guarantees. *Resource allocation* is the process of determining the amount of resources needed by a workload, while *resource assignment* means identifying the location of resources that satisfy an allocation. Both aspects of resource management require the ability to classify a given workload as a member of one of several distinct classes. Once classified, the steps are to allocate to the workload precisely the amount of resources typical for that class, assign the resources, monitor the workload execution, and adjust this amount if needed.

QoS and workload classification. Workload classification is a challenging problem due to the wide spectrum of system workloads. Thus, an effective filtering mechanism is needed to support a classification algorithm capable to make real-time decisions. Classification is widely used by recommender systems such as the one used by Netflix. A recommender system seeks to predict the preference of a user for an item by filtering information from multiple users regarding the item. Such systems are used to recommend research articles, books, movies, music, news, and any other imaginable item.

A short diversion to the Netflix Challenge [56,135] could help understanding the basis of the classification technique for a QoS-aware cluster management. The Netflix Challenge uses Singular Value Decomposition (SVD) and PQ-tree¹³ reconstruction [53,501]. SVD input is a sparse matrix A of rank r describing a system of n viewers and m movies:

$$A_{m,n} = \begin{pmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,n} \\ a_{2,1} & a_{2,2} & \dots & a_{2,n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m,1} & a_{m,2} & \dots & a_{m,n} \end{pmatrix} = U \Sigma V^t \quad (4.2)$$

¹³ A PQ tree is a tree-based data structure that represents a family of permutations on a set of elements, that is used to solve problems where the goal is to find an ordering that satisfies various constraints.

with

$$U_{m,r} = \begin{pmatrix} u_{1,1} & u_{1,2} & \dots & u_{1,r} \\ u_{2,1} & u_{2,2} & \dots & u_{2,r} \\ \vdots & \vdots & \ddots & \vdots \\ u_{m,1} & u_{m,2} & \dots & u_{m,r} \end{pmatrix} \quad \text{and} \quad V_{r,n} = \begin{pmatrix} v_{1,1} & v_{1,2} & \dots & v_{1,n} \\ v_{2,1} & v_{2,2} & \dots & v_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ v_{r,1} & v_{r,2} & \dots & v_{r,n} \end{pmatrix}. \quad (4.3)$$

Σ , the diagonal matrix of singular values of matrix A , is:

$$\Sigma_{r,r} = \begin{pmatrix} \sigma_{1,1} & 0 & \dots & 0 \\ 0 & \sigma_{2,2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_{r,r} \end{pmatrix}. \quad (4.4)$$

The ratings are the matrix elements $a_{ij} \in A$, $1 \leq i \leq m$, $1 \leq j \leq n$. SVD decomposes matrix A as $A = U \cdot \Sigma \cdot V^t$ with U being the matrix of left singular vectors representing correlation between rows of A and the similarity values, Σ is the matrix of similar values, and V is the matrix of right singular vectors representing the correlation between columns of A and the similarity values.

PQ-reconstruction with Stochastic Gradient Descent (SGD) [53] uses matrices U and V to reconstruct the missing entries in matrix A . The initial reconstruction of A uses $P^t = \Sigma \cdot V^t$ and $V = U$. SGD iterates over the elements $R = [r_{ui}]$ of $R = Q \cdot P^t$ until convergence. The iteration process uses two empirically determined values, η and λ , the *learning rate* and the SGD *regularization factor*, respectively. Two parameters μ and b_u the average rating and a user bias that account for the divergence of some viewers from the norm, respectively, are also used. The error ϵ_{ui} and the values q_i and p_u at each iteration are computed as:

$$\begin{aligned} \epsilon_{ui} &\leftarrow r_{ui} - \mu - b_u - q_i \cdot p_u^t \\ q_i &\leftarrow q_i + \eta(\epsilon_{ui} p_u - \lambda q_i) \\ p_u &\leftarrow p_u + \eta(\epsilon_{ui} q_i - \lambda p_u) \end{aligned} \quad (4.5)$$

The iteration continues until the L2 norm of the error becomes arbitrarily small

$$\|\epsilon\|_{L2} = \sqrt{\sum_{u,i} |\epsilon_{ui}|^2} \leq \epsilon. \quad (4.6)$$

The convergence speed of stochastic gradient descent is limited by the noisy approximation of the true gradient. When the gains decrease too slowly, the variance of the parameter estimate decreases equally slowly. When the gains decrease too quickly, the expectation of the parameter estimate takes a very long time to approach the optimum. SVD complexity is $\mathcal{O}(\min(n^2m, m^2n))$, and the complexity of PQ-reconstruction with SGD is $\mathcal{O}(n \times m)$.

Quasar. A performance-centric approach for cluster management is implemented by two systems developed at Stanford University, Quasar [136] and its predecessor Paragon [134]. While Paragon handles only resource assignment, Quasar implements resource allocation, as well as resource assignment.

Quasar is implemented in C, C++ and Python, runs on Linux and OS X, and supports applications written in C/C++, Java, and Python. Applications run unaltered, there is no need to modify them for running under Quasar.

Quasar is based on several innovative ideas. One of them regards reservation systems and the realization that users are seldom able to accurately predict the resource needs of their applications. Moreover, performance isolation, though highly desirable, is hard to implement. As a result, the execution time and the resources used by an application can be affected by other applications sharing the same physical platform(s). This is the reason why reservation systems often lead to underutilization of resources and why QoS guarantees are seldom offered.

Quasar presents a high-level interface to allow users, as well as schedules integrated in widely used frameworks, to express constraints for the workloads they manage. These constraints are then translated into actionable resource-allocation decisions. Classification techniques are then used to evaluate the impact of these decisions on all system workloads.

Performance constraints differ for different types of workloads. For transaction processing systems, the system bandwidth, expressed as number of queries per second, represents a meaningful constraint because it reflects the response time experienced by users. For large batch workloads, e.g., the ones involving frameworks such as *Hadoop*, the execution time captures the expectations of the end-users.

To operate efficiently, the classification algorithms is based on four parameters: resources per node, number of nodes for allocation, server type, and degree of interference for assignment. The results of the four independent classifications are combined by a greedy algorithm used to determine as accurately as feasible the set of resources needed to satisfy the performance constraints. The system constantly monitors the workload performance and adjusts the allocations when feasible.

Rather than relying exclusively on the user's characterization of the workload, the classification system combines information gathered about the workload from prescreening with information from a database about past workloads. Once accepted in the system, a workload is profiled during a short execution on a few servers with two randomly selected scale-up allocations.

For example, *Hadoop* workloads are profiled for a small number of map tasks and two configurations of parameters, such as the number of mappers per node, the size of the Java AM heap, the block size, the amount of memory per task, the replication factor, and the compression factor. A configuration includes all relevant data for the workload. The configuration data is uploaded into a table with workloads as rows and scale-up configurations as columns. To constrain the number of columns, the entries are quantized to integer multiples of cores and blocks of memory and storage.

The classification engine and the scheduler. The classification engine distinguishes between the allocation of more servers to a workload, called *resource scale out*, and additional resources from servers already allocated to the workload, called *resource scale up*. The Quasar classification engine carries out, for each workload, four classifications for scale up, scale out, heterogeneity, and interference. Some workloads may require both types of scaling, others one type or another. For example, Quasar may monitor the number of queries per second and the latency for a web server and apply both types of scaling. Initially, Quasar was focused on compute cores, memory, and storage capacity with the expectation to soon also cover the network bandwidth.

In addition to scale up and scale out, there are two other types of classifications, heterogeneity and interference. To operate with a matrix of small dimensions and, thus, reduce computational complexity, the four types of classifications are done independently and concurrently. The greedy scheduler combines data from the four classifications.

The scale-up classification evaluates how the number of cores, the cache, and the memory size affect performance. The scale-out classification is only applied to several types of workloads that can use multiple servers, and profiling is done with the same parameters as for the scale-up classification. The heterogeneity classification is the result of profiling the workload on several randomly selected servers. Lastly, interference classification reflects the sensitivity and tolerance of other workloads using shared cores, cache, memory, and communication bandwidth.

The objective of the greedy scheduler is to allocate to each workload the least amount of resources enabling it to meet its SLO (Service Level Objectives).¹⁴ For each allocation request, the scheduler ranks available servers based on the resource quality, e.g., the sustainable throughput combined with minimal interference. First, it attempts vertical scaling, allocating more resources on each node, and then, if necessary, switches to horizontal scaling allocating additional nodes, while keeping the total number of nodes as low as feasible.

Resources are allocated to applications on a FCFS basis. This could lead to sub-optimal assignments, but such assignments can be easily detected by sampling a few workloads. The scheduler also implements admission control to prevent oversubscription.

In summary, Quasar provides QoS guarantees and, at the same time, increases resource utilization. The process starts with an initial profiling of an application with short runs. The information from the initial profiling is expanded with information regarding four factors that can affect performance, scale up, scale out, heterogeneity, and interference. Then, the greedy scheduler uses the classification output to allocate resources that enable SLO compliance and maximize resource utilization.

4.11 Resource isolation

A recurring theme of this chapter is that a cluster management systems must perform well for a mix of applications and deliver the performance promised by the strict Service Level Objectives (SLOs) for each workload. The dominant components of this application mix are *latency-critical* workloads, e.g., web search, and *best-effort* batch workloads, e.g., Hadoop. The two types of workloads share the servers and compete with each other for their resources.

The resource management systems discussed up to now act at the level of a cluster, but cannot be very effective at the level of individual servers or processors. First, they cannot have accurate information simply because the state of processors changes rapidly and communication delays prohibit a timely reaction to these changes. Secondly, a centralized or even a distributed system for fine-grained, server-level resource tuning would not be scalable.

Each server should react to changing demands and dynamically alter the balance of resources used by co-located workloads. A system with feedback is needed to implement an *iso-latency policy*, in other words, to supply sufficient resources so that SLOs are met. More bluntly, this means allowing Latency Critical (LC) workloads to expand their resource portfolio at the expense of co-located best-effort workloads.

This is the basic philosophy of the Heracles system developed at Stanford and Google [312]. In this section, we discuss the real-time mechanisms used by Heracles controller to isolate co-located

¹⁴ A service level objective is a key element of a SLA. SLOs are agreed upon as a means of measuring the performance of the CSP and are a way of avoiding disputes between the users and the CSP based on misunderstanding.

workloads. In this context the term “isolate” means to prevent the best-effort workload to interfere with the SLO of the latency-critical workload.

Latency-critical workloads. A closer look at three Google latency-critical workloads, *websearch*, *ml_cluster*, and *memkeyval*, helps us better understand why resource isolation is necessary for co-located workloads. The first, *websearch*, is the query component of the web search service. Every query has a large fan-out to thousands of leaf nodes, each one of them processing the query on a shard of the search index stored in DRAM. Each leaf node has strict SLO of tens of ms. This task is compute intensive because it has to rank search hits and has a small working set of instructions, a large memory footprint, and a moderate DRAM bandwidth.

The second, *ml_cluster*, is a stand-alone service using machine-learning for assigning a snippet of text to a cluster. Its SLO is also of tens of ms. It is slightly less CPU intensive, requiring a larger memory bandwidth and lower network bandwidth than *memkeyval*. Each request for this service has a small cache footprint, but a high rate of pending requests puts pressure on the cache and DRAM.

The third, *memkeyval*, is an in-memory, key-value store used by the back-end of the web service. Its SLO latency is of hundreds of ms. The high request rate makes this service compute intensive mostly due to the CPU cycles needed for network protocol processing.

Sharing resources of individual servers is complicated because the intensity of the LC workloads at any given time is unpredictable; therefore, their latency constraints are unlikely to be satisfied at times of peak demand unless special precautions are taken. Resource reservation at the level needed for peak demand of LC workloads may come to mind first. But this naive solution is wasteful: It leads to low or extremely low resource utilization and, thus, the need for better alternatives.

Processor resources. Processor resources subject to dynamic scaling and the mechanisms for resource isolation for each one of them are discussed next. Physical cores, cache, DRAM, the power supplied to the processor, and network bandwidth are all resources that affect the ability of an LC workload to satisfy the SLO constraints. Individual resource isolation is not sufficient; cross-resource interactions deserve close scrutiny. For example, contention for cache affects DRAM bandwidth; a large network bandwidth allocated to query processing affects CPU utilization because communication protocols consume a large number of CPU cycles.

Processor cores are the engine delivering CPU cycles and an obvious target of dynamic rather than static allocation for co-located workloads. This problem is complicated by Hyper-Threading (HT) in multicore Intel processors. HT is a proprietary form of SMT (simultaneous multi-threading) discussed in Section 4.3. HT takes advantage of superscalar architecture and increases the number of independent instructions in the pipeline. For each physical core the OS uses two virtual cores and shares the workload between them when possible. This interferes with the instruction execution, shared caches, and TLB (translation look-aside buffer)¹⁵ operations.

Dynamic frequency scaling is a technique for adjusting the clock rate for cores sharing a socket. The higher the frequency, the more instructions are executed per unit of time by each core, and the larger is the processor power consumption. Clock frequency is related to the operating voltage of the processor. The *dynamic voltage scaling* is a power conservation technique often used together with frequency scaling, thus the name *dynamic voltage and frequency scaling* (DVFS).

¹⁵ TLB can be viewed as a cache for dynamic address translation; it holds the physical address of recently used pages in virtual memory.

The *overclocking* techniques based on DVFS opportunistically increase the clock frequency of processor cores above the nominal rate when the workload increases. To allow the cores of an Intel processor to adjust their clock frequency independently, the *Enhanced Intel SpeedStep* technology option should be enabled in the BIOS.¹⁶ Heracles reacts to lower best-effort workloads by reducing the number of cores assigned to best-effort tasks.

The cycle stalls limit the effective IPC (instructions per clock cycle) of individual cores. This means that the shared *last-level cache*¹⁷ is another critical resource shared by the LC and best-effort co-located workloads and should be dynamically allocated. Lastly, the DRAM bandwidth can greatly affect the performance of applications with a large memory footprint.

One answer to the question of how to implement an isolation mechanism allowing LC workloads to scale up could be to delegate this task to the local scheduler. Why not use existing work-conserving¹⁸ real-time schedulers such as SCHED_FIFO or CFS, the Completely Fair Scheduler? A short detour into the world of real-time schedulers used by most operating systems should convince us that these schedulers are designed to support data streaming and cannot satisfy the SLO requirement of LC tasks.

The SCHED_FIFO scheduler allocates the CPU to a high-priority process for as long it wants it, subject only to the needs of higher-priority realtime processes. It uses the concept of “real-time bandwidth,” *rt_bandwidth*, to mitigate the conflicts between several classes of processes; once a process exceeds its allocated *rt_bandwidth*, it is suspended. The bandwidth of LC tasks changes so SCHED_FIFO scheduler cannot accomplish what we need.

CFS uses a red–black tree¹⁹ in which the nodes are structures derived from the general *task_struct* process descriptor with additional information. CFS is based on the “sleeper fairness” concept enforcing the rule that interactive tasks that spend most of their time waiting for user input or other events get a comparable share of CPU time as other processes, when they need it. As threads are spawned for every new query request, we see that this approach is not satisfactory either.

Communication bandwidth sharing inside the server is controlled by the OS. Linux can be configured to guarantee outgoing bandwidth for the latency-critical workload. For incoming traffic it is necessary to throttle the core allocation until flow-control mechanisms of communication protocols are triggered. Communication among servers is supported by the cluster interconnection fabric and can be ensured by communication protocols that give priority to short messages typical for the latency-critical workloads.

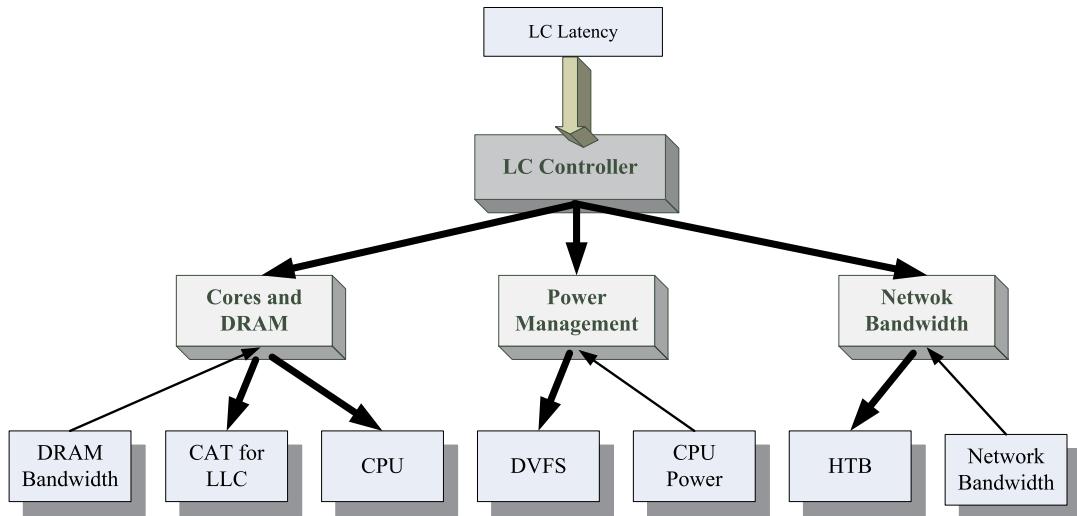
Architectural support is needed for workload isolation at the microarchitecture level. Newer generations of Intel processors, such as the Xeon E5-2600 v3 family, provide the hardware framework to manage a shared resource, like a last-level cache through Cache Monitoring Technology (CMT) and Cache Allocation Technology (CAT). CMT enables an OS or a hypervisor to determine the usage of

¹⁶ The basic input/output system (BIOS) is invoked after a computer system is powered up to load the OS and later to manage the data flow between the OS and devices such as keyboard, mouse, disk, video adapter, and printer.

¹⁷ Multicore processors have multiple level caches. The last level cache (LLC) is the cache called before accessing memory. Each core has its own L1 I-cache (instruction cache) and D-cache (data cache). Sometimes, two cores share the same unified (instruction+ data) L2 cache and all cores share an L3 cache. In this case, the highest shared LLC is L3.

¹⁸ A work-conserving scheduler tries to keep the resources busy, if there is work to be done, while a non-work conserving scheduler may leave resources idle while there is work to be done.

¹⁹ A red–black tree is a self-balancing binary search tree where each node has a “color” bit (red or black) to ensure the tree remains approximately balanced during insertions and deletions.

**FIGURE 4.5**

Heracles organization. The system runs on each server and controls isolation of a single-LC workload and multiple best-effort jobs. LC controller acts based on information regarding latency constraints of LC workload and manages three resource controllers for: (a) Core and DRAM—uses CAT for LLC management and acts based on DRAM bandwidth data; (b) Power management—acts on DVFS using information on CPU power consumption; and (c) Network bandwidth—enforces bandwidth limits for outgoing traffic from best-effort tasks using *qdisc* scheduler with HTB (token bucket queuing) discipline with information provided by network bandwidth-monitoring system.

cache by applications running on the platform. It assigns a Resource Monitoring ID (RMID) to each of the applications or VMs scheduled to run on a core and monitors cache occupancy on a per-RMID basis. CAT allows access to portions of the cache according to the *class of service* (COS).

Heracles organization and operation. Heracles runs as a separate instance on each server and manages the local interactions between the latency-critical and best-effort jobs. Fig. 4.5 shows the latency-critical controller and the three resource controllers for cores, caches, and DRAM for power management and for communication bandwidth. The controller uses the *slack*, the difference between the SLO target and the tail of the measured performance index. A negative value of the latency slack means that the time-sensitive workload has increased in intensity and is getting close to exceeding its SLO latency, and, thus, it requires more resources.

The latency-critical controller is activated every 15 units of time and uses as input the latency-critical application latency and its load. When the slack is negative or when the latency-critical load is larger than 80% of capacity, the best-effort application is disabled. If the slack is less than 10%, the best-effort task is not allowed to grow, and when the slack further decreases to 5%, then two cores allocated to best-effort tasks are removed. Best-effort is enabled when the latency-critical load is less than 80%.

The operation of the latency-critical controller is described by the following pseudocode:

```

while True
    latency = Poll_Latency-critical - AppLatency
    load   = Poll_latency-critical -AppLoad
    slack  = (target - latency)/target
    if slack < 0
        Disable Best-effort
        EnterCooldown()
    elseif load > 0.85
        Disable best-effort
    elseif load < 0.80
        Enable Best-effort
    else if slack < 0.10
        Dissallow Best-effort Growth
        if slack < 0.05
            Best-effort_core.RemoveTwoCores
    sleep(15)

```

There is a strong interaction among the number of cores allocated to a workload, its LLC cache, and DRAM requirements. This strong correlation explains why a single specialized controller is dedicated to the management of cores, LLC, and DRAM. The main objective of this controller is to avoid memory bandwidth saturation. The high-water mark that triggers action is 90% of the peak streaming DRAM bandwidth, measured by the value of hardware counters for per-core memory traffic. When this limit is reached, cores are removed from best-effort tasks.

When the DRAM bandwidth is not saturated, the gradient descent method is used to find the largest number of cores and cache partitions by alternating between increasing the number of cores and increasing the number of cache partitions allocated to best-effort tasks subject to the condition that the SLO of latency-critical tasks are satisfied. The power controller determines if there is sufficient power slack to guarantee latency-critical SLO with the lowest clock frequency. The system was evaluated with three latency-critical workloads, and the results show an average 90% resource utilization without any SLO violations.

4.12 In-memory cluster computing for Big Data

The distinction between system and application software is blurred in cloud computing. The software stack includes components based on abstractions that combine aspects of applications and system management. Two systems developed at UC Berkeley, Spark [533] and Tachyon [302], are perfect examples of such elements of a cloud software stack.

It is unrealistic to assume that very large clusters could accommodate in-memory storage of petabytes or more in the foreseeable future. Even if storage costs will decline dramatically, the intensive communication among the servers will limit the performance. There are iterative and other classes of Big Data applications where a stable subset of the input data is used repeatedly. In such cases, dramatic performance improvements can be expected if a *working set* of input data is identified, loaded in memory, and kept for future use.

Obvious examples of such applications are those involving multiple databases and multiple queries across them and interactive data mining involving multiple queries over the same subset of data. Another example of an iterative algorithm is the *PageRank* algorithm [72], where data sharing is more complex. At each iteration, a document with rank $r^{(i)}$ and n neighbors sends a contribution of $\frac{r^{(i)}}{n}$ to each one of them and updates its own rank as

$$r^{i+1} = \frac{\alpha}{N} + (1 - \alpha) \sum_{j=1}^n c_j \quad (4.7)$$

with α as the dumping factor, and N is the number of the documents in the database, and the sum over all contributions it received.

A distributed shared-memory (DSM) is a solution to in-memory data reuse. DSM allows fine-grained operations, yet access to individual data elements is not particularly useful for the class of applications discussed in this section. DSM does not support effective fault-recovery and data distribution and does not lead to significant performance improvements. Ad hoc solutions to in-memory data reuse for various frameworks have been implemented, e.g., HaLoop [77] for MapReduce.

The question is whether a data sharing abstraction suitable for a broad class of applications and use cases can be developed for supporting a restricted form of shared memory based on *coarse-grained* transformations. This abstraction should provide a simple, yet expressive, user-interface allowing the end-user to describe data transformations, as well as powerful behind-the-scene mechanisms to carry out the data manipulations in a manner consistent with the system configuration and the current state of the system.

A data sharing abstraction. The concept of *Resilient Distributed Dataset* (RDD), for fault-tolerant, parallel data structures was introduced in [532]. RDD allows a user to keep intermediate results and optimizes their placement in the memory of a large cluster. The user interface of RDD exposes: (1) partitions, atomic pieces of the dataset; (2) dependencies on the parent RDD; (3) a function for constructing the dataset; and (4) metadata about data location.

Spark provides a set of operators to efficiently manipulate such persistent datasets using a set of coarse-grained operations such as *map*, *union*, *sample*, and *join*. *Map* creates an object with the same partitions and preferred locations as its parent, but applies the function used as an argument to the call to the *iterator* method applied to the parent's records. *Union* applied to two RDDs returns an RDD whose partitions are the union of the partitions of the two parents. *Sampling* is similar to *map*, but the RDD stores for each partition a random number generator to deterministically sample parent records. *Join* creates an RDD with either two narrow, two wide, or mixed dependencies.

The driver program created by the user launches at run-time multiple workers that read data from a distributed file system such as HDFS and distributes the data across multiple RDD partitions. Tasks locality is ensured by the delay scheduling algorithm [532] used by the *Spark*'s scheduler. If a task fails, the system restarts it on a different node if the parent is available. Partitions too large to fit in memory are stored on the secondary storage, possibly on solid-state disks, if available.

Spark [533] and RDDs are restricted to *I/O* intensive applications performing bulk writing. Spark and Tachyon [302], discussed later in this section, share the concept of *lineage* to support error recovery without the need to replicate the data. Lineage means to trace back a descendent from a common ancestor and, in the context of this discussion, it means that lost output is recovered by again carrying out the tasks that created the lost data in the first place.

Spark driver programs. Imagine that an application wants to access a large log file stored in HDFS as a collection of lines of text. We wish to create a persistent dataset called *errors* of lines starting with the prefix “ERROR” distributed over the memory of the cluster, count the number of lines in this persistent dataset, and then carry out two more actions: (1) count errors containing the string “MySQL”; and (2) return the time of the errors, assuming that time is the third field in a tab-separated format of an array called HDFS. The following Spark code will do the job:

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(_.startsWith("ERROR"))
errors.persist()
errors.count
errors.filter(_.contains("MySQL")).count()
errors.filter(_.contains("HDFS")).map(_.split('\t')(3)).collect()
```

Notice that the *lines* dataset is not stored, only the much smaller *errors* dataset is stored in memory and used for the three actions. Behind the scenes, the Spark scheduler will dispatch a set of tasks carrying the last two transformations to the nodes where the cached partitions of *errors* reside.

The Spark code for the *PageRank* algorithm summarized in Eq. (4.7) is:

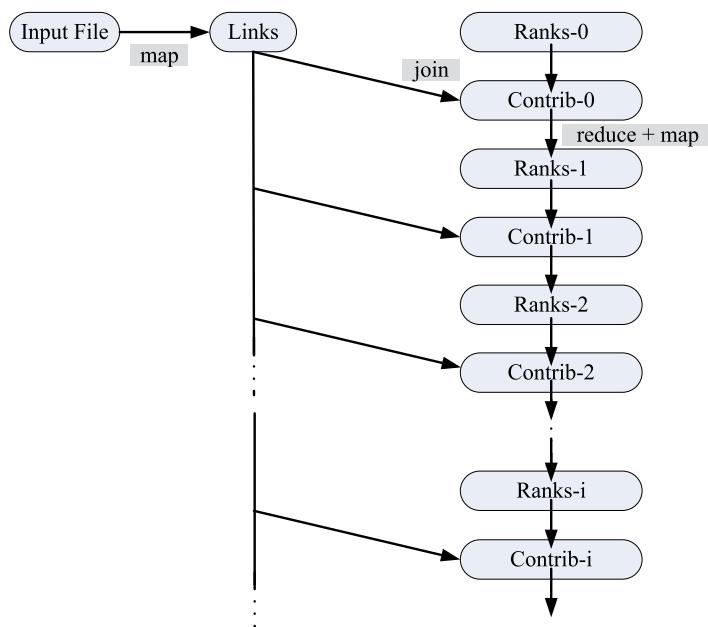
```
// Load graph as an RDD of (URL, outlinks) pairs
val links = spark.textFile(...).map(...).persist()
var ranks = // RDD of (URL, rank) pairs
for (i <- 1 to ITERATIONS) {
    // Build RDD of (targetURL, float) pairs with contributions sent by each page
    val contribs = links.join(ranks).flatMap {
        (url, (links, rank)) => links.map(dest => (dest, rank/links.size))
    }
    // Sum contributions by URL and get new ranks
    ranks = contribs.reduceByKey((x,y) => x+y).mapValues(sum => a/N + (1-a)*sum)
}
```

The *map*, *reduce*, and *join* operations and the lineage datasets of the graph of the several iterations of the *PageRank* algorithm are shown in Fig. 4.6.

A new *ranks* dataset is created at each iteration, and it is wise to call the *persist* action to save the dataset on secondary storage using the *RELIABLE* argument and reduce the recovery time in case of system failure. The *ranks* can be partitioned in the same way as the *links* to ensure that the *join* operation requires no additional communication.

Spark dependencies. An important design decision in Spark was to distinguish between *narrow* and *wide* dependencies. The former allow only one partition of RDD to use the parent RDD and enable pipelined execution on one cluster to compute all parent partitions, for example, enable the application of a *map* followed by a *filter* on an element-by-element basis. In addition, recovery after a node failure is more efficient for narrow dependencies because only the lost parent needs to be recomputed and this can be done in parallel.

On the other hand, wide dependencies enable multiple children to depend on a single parent. Data from all parent partitions must be available and shuffled across the nodes for MapReduce operations.

**FIGURE 4.6**

The lineage dataset for several iterations of the *PageRank* algorithm and the *map*, *reduce*, and *join* operations.

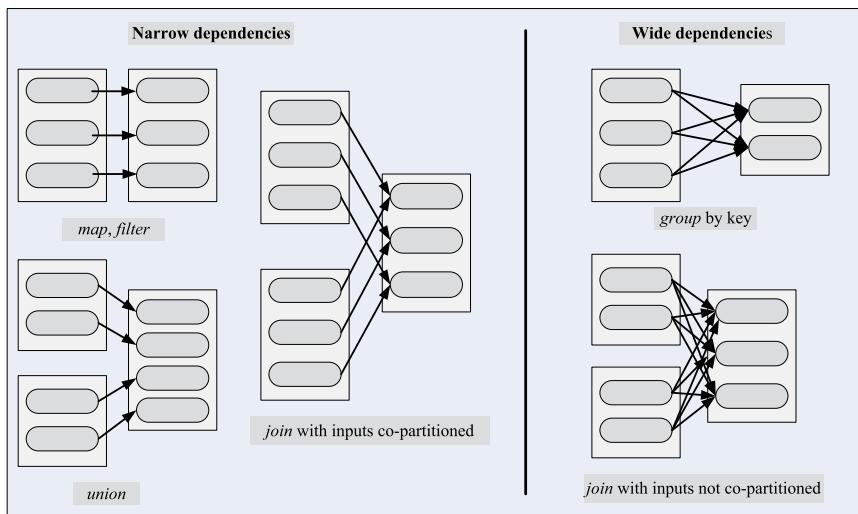
This also complicates recovery after a node failure. Fig. 4.7 shows the effect of dependencies on RDD partitions as a result of *map*, *filter*, *union*, *join*, and *group* operations.

Persistent RDD can be stored in memory either as deserialized Java objects, or as serialized data, and can also be stored on disk. Lineage is a very effective tool to recover RDDs after a node failure. Spark also supports checkpointing, and this is particularly useful for large lineage graphs.

According to [533] “Spark is up to 20 times faster than Hadoop for iterative applications, speeds-up a data analytics report 40 times and can be used interactively to scan a 1 TB dataset with 5–7 seconds latency.” It is also very powerful; only 200 lines of Spark code implement the HaLoop model for MapReduce applications. HaLoop [77] extends MapReduce with programming support for iterative applications and improves efficiency by adding various caching mechanisms and by making the task scheduler loop-aware.

These results show that caching improves dramatically Big Data applications performance running on storage systems that support only append operations, such as HGFS. Lost data can be recovered by lineage; there is no need for replication of immutable data. On the other hand, fault-tolerance of applications involving write operations is more challenging. Fault tolerance based on data replication incurs a significant performance penalty when a data item is written on multiple nodes of a cluster.

The write bandwidth throughput for both hard disks and solid-state disks is three orders of magnitude lower than the memory bandwidth. Only, the random access latency of solid-state disks is much lower than the latency of hard disks, and their sequential I/O bandwidth is not larger. The network

**FIGURE 4.7**

Narrow and wide dependencies in Spark. Arrows show the transformation of the partitions of one RDD due to *map*, *filter*, two RDDs for *union*, and two RDDs for *join* with inputs co-partitioned for narrow dependencies, when only one partition of RDD is allowed to use the parent RDD. For wide dependencies, when multiple children depend upon a single parent, two transformations are presented: a *group* by key of one RDD and *join* with input not co-partitioned for two RDDs.

bandwidth is also orders of magnitude lower than the memory bandwidth. The system discussed next addresses precisely the problem of supporting fault-tolerance for in-memory datasets where both read and write operations must be supported.

Tachyon. The system was designed for high throughput in-memory storage for applications performing both extensive reads and writes [302]. The name of the system targeting Big Data workloads discussed in Chapter 12, “Tachyon,”²⁰ was most likely chosen to reflect the performance of the system—in Greek “tachy” means “fast.” To recover lost data, the system exploits the lineage concept used also by Spark and avoids data replication, which could dramatically affect its performance.

Fault-tolerant, in-memory caching of datasets, while supporting both read and write operations, requires answers to several challenging questions:

1. How to recover the lost data due to server failures?
2. How to limit the resources and the time needed to recover lost data?
3. How to identify frequently used files and recover them with high priority when lost?
4. How to avoid recovering temporary files?

²⁰ Tachyon is also the name given to a hypothetical particle that moves faster than the speed of light. In modern physics, “tachyon” refers to imaginary mass fields rather than to faster-than-light particles.

5. How to share resources among the two activities, running the jobs and re-computations?
6. How to ensure the system's fault-tolerance?
7. How to manage the storage for binaries necessary for data recovery?
8. How to choose files to be evicted when their cumulative size exceeds the available storage space?
9. How to deal with file-name changes?
10. How to deal with changes in the cluster's run-time environment?
11. How to support different frameworks?

The answers to these questions given by the designers of the system are discussed next. Checkpointing alone is not a solution for re-computation of lost data because periodic checkpointing leads to an unbounded recovery time. Lineage alone is also not feasible because the depth of the lineage graphs keeps growing, and the time to recompute the entire path to the leaf of the graph is prohibitive. The solution adopted by Tachyon is based on combined checkpointing and lineage.

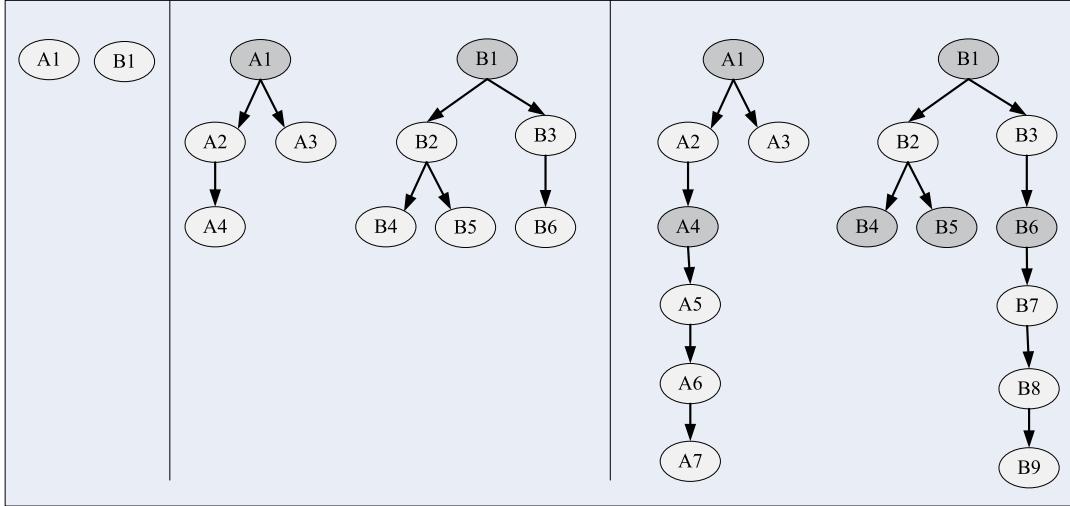
The Edge algorithm introduced in [302] checkpoints only the leaves of the lineage directed acyclic graph (DAG). This strategy reduces the number of checkpointed files and limits the resources needed for recovery. The implicit assumption of this approach is that data is immutable between consecutive checkpoints. Data should be versioned if modified between two consecutive checkpoints, and different files produced from the same parent should have different IDs.

A read counter associated with every file is used as the file priority. Frequently, read files have a high priority, while temporary files with a low read count are avoided. Re-computation is done asynchronously in the background using the lineage, thus, the interference with jobs running on the cluster is minimized and their SLOs can be guaranteed. The system is controlled by a *Tachyon Master* (TM) with several stand-by replicas ready to take over if the current master fails. If the master fails, a Paxos algorithm is used to elect the next master.

Computing the lineage of a file requires re-running the binaries of all applications executed from the instance the parent was created until the lost data was created. The workflow manager, a component of the TM, uses a DAG for each file and does a depth-first search (DFS) of nodes to reach the targeted files; it stops as soon as it reaches a node representing a file already in storage. So, it is fair to ask how much storage is necessary for the binaries of all jobs that can potentially be executed during the recovery process. Data gathered by Microsoft shows that a typical data center running some 1 000 jobs daily needs about 1 TB of storage for all the binaries executed over a period of one year [214].

Data eviction is based on a LRU (Least Recently Used) policy. This policy is justified by the access frequency and the by temporal locality. According to a cross-industry study [103], the file access in a large data center often follows a Zipf-like distribution, and 75% of re-accessing occurs within six hours. A file is uniquely identified by an immutable ID in the lineage information record to address file name changes. This ensures that the re-computation done according to the ordered list prescribed by lineage reads the same files and in the same order as in the original execution.

Among the most frequent changes in the cluster's runtime environment are changes of the version of a framework used to re-compute lost data and changes of the OS version. To address this problem, the system runs in a *Synchronous mode* before any such change when all un-replicated files are checkpointed and the new data is saved. Once this is done, this mode is disabled. Lastly, Tachyon requires a program written in a framework to provide information before writing a new file. This information is used to decide if the file should be in memory only and to recover a lost file using its lineage.

**FIGURE 4.8**

Edge algorithm. Dark-filled ellipses represent checkpointed files and light-filled ones represent un-checkpointed files. Lineage is represented by a DAG with files as vertices and edges representing the transformation of a parent to all of its descendants. Shown are the lineage DAGs of two files, A_1 and B_1 . At each stage, only the leaf nodes of the lineage DAG are checkpointed. A_1 and B_1 are checkpointed first, then A_4 , B_4 , B_5 , and B_6 . In the next stage, not shown in the figure, only A_7 and B_9 will be checkpointed.

Edge algorithm. The algorithm assumes that the lineage is represented by a DAG with files as vertices and edges representing the transformation of a parent to all of its descendants. The algorithm checkpoints the leafs of the graph. For example, assume a lineage chain for a file A_0 including files $\{A_0, A_1, A_2, \dots, A_i, \dots\}$. Then, if there is a checkpoint of A_i and A_j is lost, the re-computation starts with the latest checkpoint, in this case A_i , rather than A_0 . Fig. 4.8 shows the lineage DAGs of two files A_1 and B_1 and the leafs checkpointed at several instances of time; A_1 and B_1 are checkpointed first, then A_4 , B_4 , B_5 , and B_6 .

The Edge algorithm does not take into account priorities. A balanced algorithm alternates edge-driven checkpointing with priority-based checkpointing, allocating a fraction c of time to the former and a fraction $(1 - c)$ of time to the latter. To guarantee applications SLOs the recovery time for any file must be bounded.

Call W_i the time to checkpoint an edge i of the DAG and G_i to generate edge i from its ancestors. Then two bounds on the recovery time for any file are proven in [302]: (1) Edge checkpointing alone leads to the following bound of the recovery time

$$\mathcal{T}^{\text{edge}} = 3 \times M \text{ with } M = \max_i \{T_i\}, \text{ and } T_i = \max(W_i, G_i)^9. \quad (4.8)$$

This shows that the re-computation time is independent of the DAG's depth. (2) The bound for the recovery time for alternating edge and priority checkpointing is

$$\mathcal{T}^{\text{edge,priority}} = \frac{3 \times M}{c} \quad \text{with } M = \max_i\{T_i\}, \quad \text{and } T_i = \max(W_i, G_i). \quad (4.9)$$

Resource management. Resource management policies should address several concerns. When several files have to be recovered at the same time the system should consider data dependencies to avoid recursive task launching. Dynamic file checkpointing priorities are also necessary; low priority assigned to a file requested by a low priority job should be automatically increased when the same file is requested by a high priority job. The lineage record should be deleted after checkpointing to save space.

All resources should be dedicated to normal work when no recovery is necessary. Typical average server utilization rarely exceeds 30%, so most of the time, there are sufficient resources available for data recovery. But what to do when the system is near its capacity? Then the priority of the jobs and of the recovery come into play and is used by the cluster scheduler.

Tachyon could accommodate the two most frequently used scheduling policies for cluster resource management, priority based and fair-shared based scheduling. In case of *priority-based* scheduling, all re-computation jobs are given the lowest priority by default. Deadlock may occur unless precautions are taken.

For example, assume that a job \mathcal{J}_i has a higher priority than the file \mathcal{F} it needs to recover and that job J is scheduled to run. At the time when \mathcal{J}_i needs access to \mathcal{F} , the job $\mathcal{R}_{\mathcal{F}}$ needed to recover \mathcal{F} cannot run as it inherits the lower priority of \mathcal{F} . The solution is *priority inheritance*. In this case, \mathcal{J}_i and, implicitly the job $\mathcal{R}_{\mathcal{F}}$, should inherit the priority of \mathcal{J}_i that needs \mathcal{F} . If another job with an even higher priority needs the file \mathcal{F} its priority, hence the priority of $\mathcal{R}_{\mathcal{F}}$, should increase again.

Assume now a *fair-share scheduler*. In this case, $W_1, W_2, \dots, W_i, \dots$ are the weights of resources allocated to jobs $\mathcal{J}_1, \mathcal{J}_2, \dots, \mathcal{J}_i, \dots$, respectively. The minimal share unit is $W_g = 1$. When files $\mathcal{F}_{i,1}, \mathcal{F}_{i,2}, \mathcal{F}_{i,3}$ of job \mathcal{J}_i are lost the scheduler allocates the minimum weight W_g from the W_i to the three recovery jobs $\mathcal{R}_{\mathcal{J}_i, \mathcal{F}_1}, \mathcal{R}_{\mathcal{J}_i, \mathcal{F}_2}$ and $\mathcal{R}_{\mathcal{J}_i, \mathcal{F}_3}$. At the time job \mathcal{J}_i needs to access the file $\mathcal{F}_{i,2}$ the scheduler allocates the fractions $(1 - \alpha)$ and α of W_i to \mathcal{J}_i and $\mathcal{R}_{\mathcal{J}_i, \mathcal{F}_2}$, respectively.

Tachyon implementation. Tachyon has two layers: the *lineage layer* tracks the sequence of jobs that have created a file; the *persistence layer* manages data in storage, can be any replication-based storage such as HDFS, and is used for asynchronous checkpoints. Tachyon has a master-slave architecture with several passive replicas of the master and worker daemons running on every cluster node and managing local resources. The lineage information is tracked by a workflow manager running inside the master. The workflow manager computes the order of checkpoints and interacts with the cluster resource manager to get resources needed for re-computations.

Each worker uses a RAM disk for storing memory-mapped files. The concept of wide and narrow dependences inherited from Spark is used to carry out the operations discussed earlier in this section. The system was implemented in about 36 000 lines of Java code and uses the *Zookeeper* for election of a new master when one fails.

Tachyon lineage can capture the requirements of MapReduce, and SQL, as well as Hadoop and Spark can run on top of it. According to [302] Tachyon has a 110 times higher write throughput and for realistic workloads, improves end-to-end latency four times compared with in-memory HDFS. It

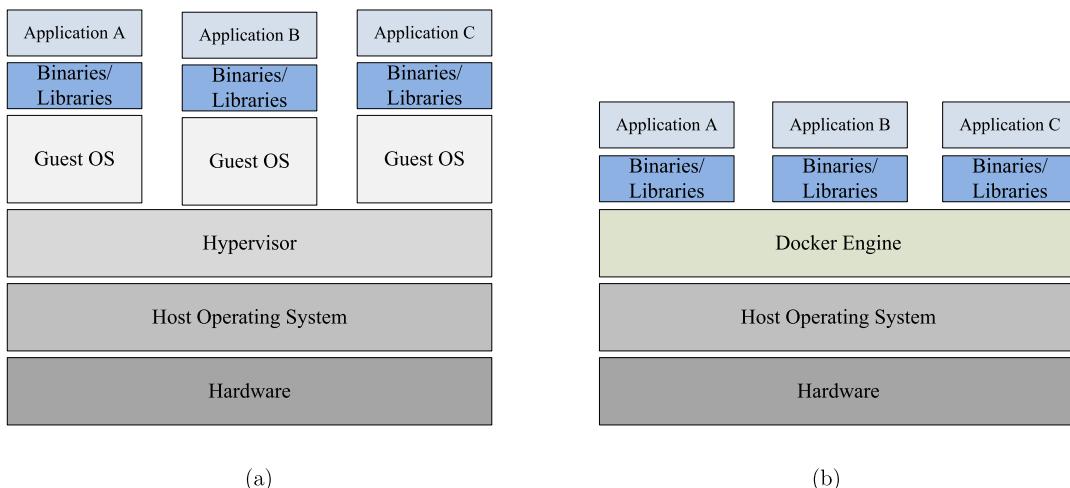


FIGURE 4.9

System organization for Virtual Machines and Docker containers. (a) VM; (b) Multiple Docker containers running on the same system share the OS kernel thus, have a smaller memory footprint and a shorter start-up time than VMs.

can also reduce network traffic by up to 50% because many temporary files are deleted before being checkpointed. Data from Facebook and Bing show that it consumes not more than 1.65% of cluster resources for re-computations.

4.13 Containers; Docker containers

This idea of containers was extended from file systems supported by *chroot* to other namespaces including the process IDs. Initially, the *cgroups* concept was implemented at Google in 2006 in the Linux kernel to isolate, control, limit, prioritize, and account for the resources (CPU, memory, disk I/O, networks) available to a set of processes, by freezing groups of processes and managing checkpointing and restarting. Resource limiting enforces the target set for resource utilization, while prioritization allows a group of processes to get a larger share of CPU cycles or a higher disk I/O throughput.

Docker created containers along with a set of tools with standard APIs; it also made the same container portable across all environments. “Docker containers wrap a piece of software in a complete filesystem that contains everything needed to run: code, runtime, system tools, system libraries—anything that can be installed on a server. This guarantees that the software will always run the same, regardless of its environment,” according to Docker. Fig. 4.9 depicts the organization of both VM- and container-based systems.

Containers are light weight, cost-effective in a public cloud environment, provide better performance, and require fewer hardware resources for a private cloud. Containers isolate an application from the underlying infrastructure and from other applications and support performance and security.

isolation. Multiple containers running on the same machine share the OS kernel, have a smaller memory footprint, and a shorter start-up time than VMs.

Another advantage of Docker containers is increased productivity. Containerization allows developers to choose the most suitable programming languages and software systems and eliminates the need to make copies of production code and install the same configuration in different environments. It also supports efficient up and down scaling of an application.

Docker ecosystem is built around a few concepts discussed next. An *image* is a blueprint of an application. A *container* consists of one or more images and runs the actual application. A *daemon* is a background service running on the host that manages building, running, and distributing Docker containers. The daemon is the process that runs under the operation system to which clients talk to. A *client* is a command line tool used by a user to interact with the daemon. A *hub* is a registry of Docker images, a directory of all available Docker images.

There are two types of images, base and child. *Base images* have no parent image, usually images with an OS like Ubuntu, or Busybox.²¹ *Child images* are build on base images with additional functionality. *Official images* are maintained and supported by Docker, and have one word name; for example, *python*, *ubuntu* are base official images. *User images* are created and shared by users. A *Dockerfile* is a facility to automate the image creation, a text-file including Linux-like commands invoked by clients to create an image.

Docker Swarm exposes standard Docker API. Docker tools including Docker CLI, Docker Compose, Dokku, and Krane, work in this native Docker clustering. The distribution of it is packed as a Docker container and to set it up one only needs to install one of the service discovery tools and run the *swarm* container on all nodes, regardless of the OS.

Cloud computing has embraced containerization. Containers-as-a-Service (CasS) is geared toward efficiently running a single application. Several CSPs including Heroku, OpenShift, dotCloud and CloudFoundry use containers to support PaaS delivery model. Amazon, Google, Microsoft, OpenStack, Cloudstack and other CSPs offering IaaS implementations support containers. Container support by Amazon, Google, and Microsoft is overviewed next.

ECS is *Amazon EC2 Container Service*. It is straightforward for AWS users to create and manage ECS clusters, as ECS is integrated with existing services such as IAM for permissions, CloudTrail to get data regarding resources used by a container, CloudFormation for cluster launching, and other services. AWS uses a custom scheduler/cluster manager for containers. Container hosts are regular EC2 instances. To deploy a containerized application on AWS one has to first publish the image on an AWS accessible registry e.g., the Docker Hub. Amazon ECS is free, while Google Container Engine, discussed next, is free up to five nodes. The charge granularity at AWS is one hour while at Google or Microsoft the charge is for the actual time used.

Google Container Engine (GKE). GKE is based on Kubernetes, an open source cluster manager available to Google developers and the community of outside users. Google's approach to containerization is slightly different, its emphasis is more on performance than ease of use as discussed in Section 4.14. Two billion containers are started at Google every week according to http://www.theregister.co.uk/2014/05/23/google_containerization_two_billion/. GKE is integrated with other

²¹ BusyBox is software providing several stripped-down Unix tools in a single executable file and running in environments such as Linux, Android, FreeBSD, or Debian.

services including Google Cloud Logging. Google users have access by default to private Docker registry and a JSON-based declarative syntax for configuration. The same syntax can be used to define what happens with the hosts. GKE can launch and terminate containers on different hosts as defined in the configuration file.

Microsoft Azur Container Service. The Azure Resources Manager API supports multiple orchestrations including Docker, and Apache Mesos.

In recent years the Open Container Initiative (OCI) was involved in an effort to create industry-standard container format and runtime systems under the Linux Foundation. The 40+ members of the OCI work together to make Docker more accessible to the large cloud users community. OCI's approach is to break Docker into small reusable components. In 2016 OCI has announced Docker Engine 1.11 which uses a daemon, *containers* to control *runc*²² or other OCI compliant run-time systems to run containers. Users access the Docker Engine via a set of command and a user interface.

4.14 Kubernetes

Kubernetes is a cluster manager for containers. The origin of the word Kubernetes is in ancient Greek language: “kubernan” means to steer and “kubernetes” is helmsman. Kubernetes is an open source software system developed and used at Google for managing containerized applications in a clustered environment. Kubernetes bridges the gap between a clustered infrastructure and assumptions made by applications about their environments. Kubernetes is written in *Go*²³ and it is designed to work with operating systems that offer lightweight virtual computing nodes. The system is lightweight, modular, portable and extensible. Mesos is augmenting Kubernetes and expected to support Kubernetes API.

Kubernetes is an open-ended system and its design allowed a number of other systems to be build atop Kubernetes. The same APIs used by its control plane are also available to developers and users who can write their own controllers, schedulers, etc., if they choose so. The system does not limit the type of applications, does not restrict the set of runtimes or languages supported and allows users to choose the logging, monitoring, and alerting systems of their choice.

Kubernetes does not provide built-in services including message buses, data-processing frameworks, databases, or cluster storage systems and does not require a comprehensive application configuration language. It provides deployment, scaling, load balancing, logging, and monitoring. Kubernetes is not monolithic, and default solutions are optional and pluggable.

Kubernetes organization. A master server manages a Kubernetes cluster and provides services to manage the workload and to support communication for a large number of relatively unsophisticated *minions* servers that do the actual work. *Etcd* is a lightweight, distributed key-value store to share configuration data with the cluster nodes. The master also provides an API service; an HTTP/JSON API is used for service discovery. The Kubernetes scheduler tracks resources available and those allocated to the workloads on each host.

²² *runc* is an implementation of the Open Containers Runtime specification and the default executor bundled with Docker Engine. Its open specification allows developers to specify different executors without any changes to Docker itself.

²³ *Go* or *Golang* is open source compiled, statically typed language like Algol and C, has garbage collection, limited structural typing, memory safety features, and CSP-style concurrent programming.

Minions use a *docker* service to run encapsulated application containers. A *kuberlet* service allows minions to communicate with the master server and with the *etcd* store to get configuration details and update the state. A proxy service of the minions interacts with the containers and provides a primitive load balance.

In Kubernetes, *pods* are groups of containers scheduled onto the same host and serve as units of scheduling, deployment, horizontal scaling, and replication. Pods share fate and share resources such as storage volumes. The *run* command is used to create a single container pod and the *Deployment* which monitors that pod. All applications in a pod share a network namespace, including the IP address and the port space thus can communicate with each other using *localhost*.

Pods manage co-located support software including: content management systems, controllers, managers, configurators, updaters, logging and monitoring adapters, and event publishers. Pods also manage file and data loaders, local cache managers, log and checkpoint backup, compression, rotation, snapshotting, data change watchers, log trailers, proxies, bridges, and adapters.

Arguably, pods are preferable to running multiple applications in a single Docker container for several reasons: (i) efficiency—containers can be lightweight as the infrastructure takes on more responsibility; (ii) transparency and user convenience—containers within the pod are visible to the infrastructure and allow it to provide process management, resource monitoring, and other services; (iii) users do not need to run their own process managers, or deal with signal and exit-code propagation; and (iv) decoupling software dependencies—individual containers may be versioned, rebuilt, and redeployed independently.

Kubernetes *replication controllers*—handle the lifecycle of containers. The pods maintained by a replication controller are automatically replaced if they fail, get deleted, or are terminated. A replication controller supervises multiple pods across multiple nodes. *Labels* provide the means to find and query containers and *services* identify a set of containers performing a common function.

Kubernetes has different CLI, API, and YAML²⁴ definitions than Docker and has a steep learning curve. Kubernetes setup is more complicated than Docker Swarm and its installation differs for different operating systems and service providers.

4.15 Further readings

There is little doubt that Amazon has a unique position in the cloud computing world. There is a wealth of information on how to use AWS services on the AWS web site, but little has been published about the algorithms and the mechanisms for resource allocation and the software used by AWS. The effort to maintain a shroud of secrecy probably reflects the desire to maintain Amazon's advantage over its competitors. There are only a few papers, including [255,530], describing research results related to Microsoft's Azur cloud platform.

In stark contrast to Amazon and Microsoft, Google's research teams publish often and have made a significant contribution to understanding the challenges posed by very large systems. The evolution of ideas and Google's perspective in cloud computing is presented in [131]. An early discussion of Google

²⁴ CLI stands for Command Line Interface and provides the means for a user to interact with a program; YAML is a human-readable data serialization language.

cluster architecture is presented in [50]. The current hardware infrastructure and the Warehouse Scale Computers are analyzed in a 2013 book [52] and in a chapter of a classical computer architecture book [232]. A very interesting analysis of WSC performance is due to a team including Google researchers [265]. The discussion of multicores best suited for typical Google workloads is presented in [242].

There is a wealth of information regarding cluster management and the systems developed at Google: Borg [495], Omega [442], Quasar [136], Heracles [312], and Kubernetes [80]. Controlling latency is discussed in [130]. Performance analysis of large-scale systems using Google trace data is reported in [412].

Important research results related to cloud computing have been reported at UC Berkeley, where Mesos was designed [240], Stanford [134,136,311,312], and Harvard [265]. Cloud energy consumption is analyzed in many publications, including [9,30,45,51,52,331,494,499].

A very positive development is the dominance of open software. Open software is available from Apache, Linux Foundation, and others. Docker software and tutorials can be downloaded from <https://www.docker.com/> and <https://www.digitalocean.com/community/tags/docker?type=tutorials> and <http://prakhar.me/docker-curriculum/>, respectively.

Detailed information about Kubernetes is at <http://kubernetes.io/> and <https://www.digitalocean.com/community/tutorials/an-introduction-to-kubernetes>. The Open Containers Initiative of the Linux Foundation has developed technologies for container-based applications, see <https://www.opencontainers.org/news/news/2016/04/docker-1.11-first-runtime-built-containerd-and-based-oci-technology>. A distributed main memory processing system is presented in [263].

4.16 Exercises and problems

Problem 1. The average CPU utilization is an important measure of performance for the cloud infrastructure. [52] reports a median CPU utilization in the 40–70% range, while [265] reports a median utilization around 10% consistent with utilization reported for the CloudSuite [170]. Read the three references. Discuss the results in [265] and explain the relation between CPU utilization and memory bandwidth and latency.

1. Discuss the effects of cycle stalls and of the ILP (Instruction Level Parallelism) on processor utilization. Analyze the data on cycle stalls and ILP reported in [265].
2. Identify the reasons for the high ratio of cache stalls and the low ILP for data-intensive WSC workloads reported in [265].
3. Why do WSC workloads exhibit the high ratio of cache stalls and the low ILP?
4. What conclusions regarding memory bandwidth and latency can be drawn from the results reported in [265]? Justify your answers.

Problem 2. Discuss the results regarding simultaneous multithreading (SMT) reported in [265].

1. For what type of workloads is SMT most effective? Explain your answer.
2. The efficacy of SMT can be estimated by comparing specific per-hyperthread performance counters with ones aggregated on a per-core basis. This is very different from measuring the speedup that a single application experiences from SMT. Why?
3. Why is it difficult to measure the SMT efficiency in a cloud environment?

Problem 3. Mesos is a cluster management system designed to be robust and tolerant to failure. Read [240] and answer the following questions:

1. What are the specific means to achieve these design goals?
2. It is critical to make the Mesos *master* fault-tolerant because all frameworks depend on it. What are the special precautions to make the *master* fault-tolerant?

Problem 4. Borg system is a cluster manager that runs hundreds of thousands of jobs, from many thousands of different applications, across a number of clusters, each with up to tens of thousands of machines. Read [495] and answer the following questions:

1. Does Borg have an admission control policy? If so, describes the admission control mechanism.
2. What are the elements that make the Borg scheduler scalable?
3. How does the job mix on Borg cells affect the CPI (Cycles per Instruction)?

Problem 5. Omega is a scalable cluster management system based on a parallel scheduler architecture built around shared state, using lock-free optimistic concurrency control. Read [442] and answer the following questions:

1. What scheduler performance metrics are used for Omega, why is each one of them relevant, and how it is actually measured?
2. Trace-driven simulation was used to gain insights into the system. What are the benefits of trace-driven simulation, and how was it used to investigate conflicts?
3. One of the simulation results refers to gang scheduling. What is gang scheduling, and why it is beneficial for MapReduce applications?

Problem 6. Quasar is a QoS-aware cluster management system using a fast classification techniques to determine the impact of different resource allocations and assignments on workload performance.

1. Read [135] to understand the relationship between the Netflix challenge and the cluster resource allocation problem.
2. Quasar [136] classifies resource allocation for scale up, scale out, heterogeneity, and interference. Why are these classification criteria important, and how are they applied?
3. What are stragglers, and how does Quasar deal with them?

Problem 7. Cluster management systems must perform well for a mix of applications and deliver the performance promised by the SLOs for each workload. Resource isolation is critical for achieving strict SLOs.

1. What are the mechanisms used by Heracles [312] for mitigating interference?
2. Discuss the results related to latency of Latency Sensitive workload (LS) from Heracles experiments.
3. Discuss the results related to Effective Machine Utilization (EMU) from Heracles experiments.

Problem 8. Effective cloud resource management requires understanding the interaction between the workloads and the cloud infrastructure. An analysis of both sides of this equation uses a trove of trace data provided by Google. This analysis is reported in [412].

1. What conclusions can you draw from the analysis of the trace data regarding the schedulers used for cluster management?
2. What are the reasons for scheduler behavior revealed by the trace analysis?
3. What characteristics of the Google workloads are most notable?

Problem 9. Tachyon is a distributed file system enabling reliable data sharing at memory speed across cluster computing frameworks. Read [302] and answer the following questions:

1. Does the evolution of memory, storage, and networking technologies support the argument that cloud storage systems should achieve fault-tolerance without replication?
2. What is file popularity, and what is the distribution of file popularity for big data workloads?
3. How is this distribution used?
4. For what type of events is this distribution particularly important?

Cloud resource virtualization

5

Virtualization abstracts the underlying physical resources of computer and communication systems and simplifies their use, isolates users from one another, and supports replication, which, in turn, increases system elasticity and reliability. Resource virtualization, the technique analyzed in this chapter, is a ubiquitous alternative to traditional data center operation. Virtualization, a basic tenet of cloud computing, simplifies some of the resource management tasks. For example, a *virtual machine* (VM) running under a hypervisor can be migrated to another server to balance the load and save energy. At the same time, virtualization allows users to operate in environments they are familiar with, rather than forcing them to work in idiosyncratic environments.

Interpreters, memory, and communication channels are fundamental abstractions necessary to describe the operation of a computing system [430]. Their respective physical realizations are: processors to transform information; primary and secondary memory to store information; and communication systems enabling all functional units to interact with one another. Processors, memory, and communication systems have different bandwidths, latencies, reliabilities, and other physical characteristics.

System software transforms physical implementations of the three abstractions into computer systems able to run diverse applications and manage all system resources. The traditional modus operandi of a data center is to install an operating system (OS) on individual computers and rely on conventional OS techniques to ensure resource sharing, application protection, and performance isolation.

Cloud service providers (CSPs), as well as cloud users, face multiple challenges in such a setup. CSPs are under pressure to optimize accounting, security, and system resource management. Users are under pressure to develop and optimize their application performance for one system and have to start over again when the application is moved to another data center with different hardware, system software, and libraries. Virtualization alleviates these problems for CSPs and cloud users, but at a price.

Resource sharing in a VM environment requires not only ample hardware support and, in particular, powerful processors and fast interconnects but also architectural support for multilevel control because resource sharing occurs at multiple levels. Resources, such as CPU cycles, memory, secondary storage, and I/O and communication bandwidth, are shared among several VMs. VM resources are shared among processes or threads of an application.

This chapter starts with a discussion of virtualization principles and the motivation for virtualization in Section 5.1. Section 5.2 is focused on performance and security isolation. Alternatives for the implementation of virtualization are analyzed in Section 5.3. Two distinct approaches to processor virtualization, *full virtualization* and *paravirtualization*, are presented in Section 5.4. Full virtualization is feasible when the hardware abstraction provided by the hypervisor is an exact replica of the physical hardware. In this case any operating system running on the hardware will run without modifications under the hypervisor. Paravirtualization requires modifications of the guest operating systems because the hardware abstraction provided by the hypervisor does not support all hardware operations.

Traditional processor architectures were conceived for one level of control as they support two execution modes, the kernel and the user mode. In a virtualized environment, all resources are under the control of a hypervisor, and a second level of control is exercised by the guest OS. While a two-level scheduling for sharing CPU cycles can be easily implemented, sharing of resources such as cache, memory, and I/O bandwidth is more intricate. In 2005 and 2006, the x86 processor architecture was extended to provide hardware support for virtualization, as discussed in Section 5.5. Nested virtualization allows hypervisors to run inside a VM, complicating even further the virtualization landscape.

Section 5.6 covers an open-source emulator and virtualizer. Several hypervisors are used nowadays. Xen and an optimization of its network performance are presented in Sections 5.8 and 5.9. KVM, a virtualization infrastructure of the Linux kernel, is discussed in Section 5.7. Nested virtualization is analyzed in Section 5.10, followed by the presentation of a trusted kernel virtualization in Section 5.11.

High-performance processors have multiple functional units, but do not provide explicit support for virtualization, as discussed in Section 5.12 which covers Itanium paravirtualization. System functions critical for the performance of a VM environment are cache and memory management, handling of privileged instructions, and I/O handling.

Cache misses are an important source of performance degradation in a VM environment as we shall see in Section 5.13. An overview of open-source software platforms for virtualization is presented in Section 5.14. The potential risks of virtualization are the subject of Section 5.15, and virtualization software is discussed in Section 5.16.

5.1 Resource virtualization

Virtualization has been used successfully since the late 1950s; a virtual memory based on paging was first implemented on the Atlas computer at the University of Manchester in the United Kingdom in 1959.

Virtualization simulates the interface with a physical object by any one of four means [430]:

1. *Multiplexing*: create multiple virtual objects from one instance of a physical object. For example, a processor is multiplexed among a number of processes or threads.
2. *Aggregation*: Create one virtual object from multiple physical objects. For example, a number of physical disks are aggregated into a RAID disk.
3. *Emulation*: Construct a virtual object from a different type of a physical object. For example, a physical disk emulates a Random Access Memory.
4. *Multiplexing and emulation*. Examples are: virtual memory with paging multiplexes main memory and secondary storage; a virtual address emulates a real address; the TCP protocol emulates a reliable bit pipe and multiplexes a physical communication channel and a processor.

Virtualization is a critical aspect of cloud computing, equally important for the providers and the consumers of cloud services, and plays an important role for: (i) system security because it enables isolation of services running on the same hardware; (ii) performance and reliability because it enables applications to migrate from one platform to another; (iii) the development and management of services offered by a provider; and (iv) performance isolation.

In a cloud computing environment, a hypervisor or virtual-machine monitor runs on the physical hardware and exports hardware-level abstractions to one or more guest operating systems. A guest OS

interacts with the virtual hardware in the same manner it would interact with the physical hardware, but under the watchful eye of the hypervisor that traps all privileged operations and mediates the interactions of the guest OS with the hardware. For example, a hypervisor would control I/O operations for two virtual disks implemented as two different sets of tracks on a physical disk. New services can be added without the need to modify an operating system.

User convenience is a necessary condition for the success of the utility computing paradigm; one of the multiple facets of user convenience is the ability to run remotely using the system software and libraries required by the application. User convenience is a major advantage of a VM architecture versus a traditional operating system. For example, an AWS user could submit an Amazon Machine Image (AMI) containing the applications, libraries, data, and the associated configuration settings; a user could choose the operating system for the application, then start, terminate, and monitor as many instances of AMI as needed, using web service APIs and performance monitoring and management tools provided by AWS.

There are side effects of virtualization, notably the *performance penalty* and the *hardware costs*. All privileged operations of a VM must be trapped and validated by the hypervisor, which ultimately controls the system behavior. The increased overhead introduced by the hypervisor has a negative impact on the performance.

The cost of a system running multiple VMs is higher than the cost of a system running a traditional OS. In the former case, the physical hardware is shared among a set of guest operating systems, and it is typically configured with faster and/or multicore processors, more memory, larger disks, and additional network interfaces as compared to a system running a traditional operating system.

5.2 Performance and security isolation in computer clouds

To exhibit predictable performance an application must be isolated from applications it shares resources with. *Performance isolation* is a critical condition for QoS guarantees in cloud computing where all system resources are shared. If the run-time behavior of an application is affected by other applications running concurrently and competing for CPU cycles, cache, main memory, disk, and network access, it is rather difficult to predict application completion time and to optimize its performance. Performance unpredictability is a “deadly sin” for real-time operation and for embedded systems.

Several operating systems including Linux/RK [369], QLinux [466], and SILK [54] support some performance isolation. In spite of the efforts to support performance isolation, interactions among applications sharing the same physical system, often described as *QoS crosstalk*, still exist [473]. Accounting for consumed resources and the overhead of system activities, including context switching and paging to individual users, is challenging.

Processor virtualization presents multiple copies of the same physical processor or core to applications. Processor virtualization is different from *processor emulation* when the machine instructions of guest system are “emulated” in software running on the host system. Emulation is much slower than virtualization. For example, Microsoft’s VirtualPC was designed to run on x86 processor architecture. VirtualPC was running on emulated x86 hardware until Apple adopted Intel chips.

Traditional operating systems multiplex multiple processes or threads, while virtualization supported by a hypervisor multiplexes full operating systems. A hypervisor executes directly on the hardware a subset of frequently used machine instructions generated by the application and emulates

privileged instructions including device I/O requests. The subset of instructions executed directly by the hardware includes arithmetic instructions, memory access, and branching instructions. The overhead of context switching carried out by the hypervisor is larger and leads to a significant performance.

Operating systems use the process abstraction not only for resource sharing but also to support isolation. Unfortunately, this is not sufficient from a security perspective because once a process is compromised, it is rather easy for an attacker to penetrate the entire system.

Applications running under a VM can only access virtual devices emulated by the software. This layer of software has the potential to provide a level of isolation nearly equivalent to the isolation presented by two different physical systems. Therefore, virtualization can be used to improve security in a cloud computing environment.

A hypervisor is a much simpler and better specified system than a traditional OS. For example, the Xen hypervisor discussed in Section 5.8 has approximately 60 000 lines of code, while the *Denali* hypervisor [514] has only about half, i.e., 30 000 lines of code.

The security vulnerability of hypervisors is considerably reduced because the systems expose a much smaller number of privileged functions. For example, Xen can be accessed through 28 hypercalls, while a standard Linux allows hundreds, e.g., Linux 2.6.11 allows 289 system calls. A traditional OS supports special devices, e.g., `/dev/kmem`, and many privileged third-party programs, e.g., `sendmail` and `sshd`, in addition to a plethora of system calls.

5.3 Virtual machines

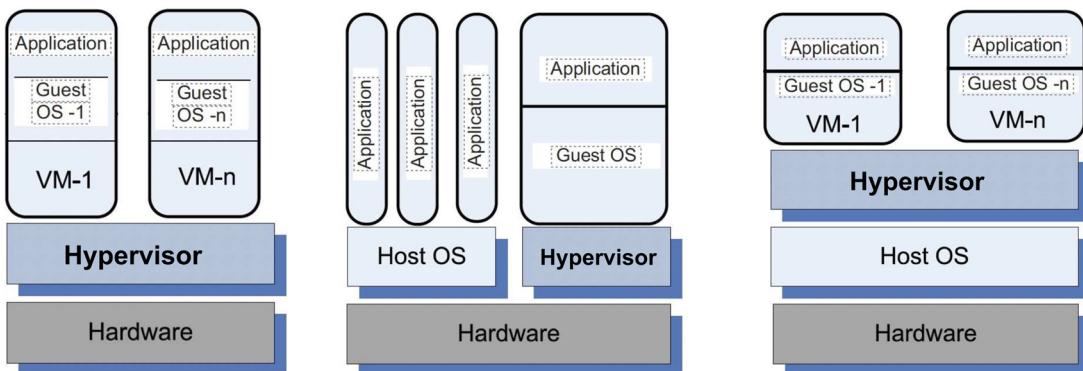
A VM is an isolated environment with access to a subset of the physical resources of a computer system. Each VM appears to be running on the bare hardware, giving the appearance of multiple instances of the same computer, though all are supported by a single physical system. VM history can be traced back to the 1960s.¹ In the early 1970s, IBM released its widely used VM 370 system, followed in 1974 by the MVS (Multiple Virtual Storage) system.

There are two types of VMs, process and system. A *process VM* is a virtual platform created for an individual process and destroyed once the process terminates. Virtually all operating systems provide a process VM for each one of the applications running, but the more interesting process VMs are those that support binaries compiled on a different instruction set. A *system VM* supports an OS together with many user processes. When the VM runs under the control of a normal OS and provides a platform-independent host for a single application, we have an *application VM*, e.g., Java Virtual Machine (JVM).

A *system VM* provides a complete system: Each VM can run its own OS, which in turn can run multiple applications. Systems such as Linux Vserver <http://linux-vserver.org>, OpenVZ (Open Virtual-iZation) [376], FreeBSD Jails [415], and Solaris Zones [403] based on Linux, FreeBSD, and Solaris, respectively, implement *OS-level virtualization technologies*.

The OS-level virtualization enables a physical server to run multiple isolated OS instances subject to several constraints; the instances are known as containers, Virtual Private Servers, or Virtual Envi-

¹ In 1963, MIT announced Project MAC (Multiple Access Computer) and chose the GE-645 as the mainframe for its Multics project. IBM was GE's competitor, realized that there is a demand for such systems, and designed the CP-40 mainframe, followed by the CP-67, also called CP/CMS mainframes. CP was a program running on the mainframe used to create VMs running a single-user OS called CMS.

**FIGURE 5.1**

Hypervisor's role in traditional, hybrid, and hosted VM architectures. (Left) Traditional VMs; hypervisor supports multiple VMs and runs directly on the hardware. (Center) Hybrid VM; hypervisor shares the hardware with a host OS and supports multiple VMs. (Right) Hosted VM; hypervisor runs under a host OS.

ronments. For example, OpenVZ requires both the host and the guest OS to be Linux distributions. These systems claim performance advantages over the systems based on a hypervisor, such as Xen or VMware. According to [376], there is only a 1% to 3% performance penalty for OpenVZ as compared to a standalone Linux server. OpenVZ is licensed under the GPL version 2.

A hypervisor enables several VMs to share a physical system. Several organizations of the software stack supporting virtualization are possible:

- *Traditional*—VM also called a “bare metal” hypervisor—a thin software layer that runs directly on the host machine hardware; its main advantage is performance, Fig. 5.1(b). Examples: VMware ESX, ESXi Servers, Xen, OS370, and Denali.
- *Hybrid*—the hypervisor shares the hardware with an existing OS, Fig. 5.1(c). Example: VMWare Workstation.
- *Hosted*—the VM runs on top of an existing OS, Fig. 5.1(d); the main advantage of this approach is that the VM is easier to build and install. Another advantage is that the hypervisor could use several components of the host OS, such as the scheduler, the pager, and the I/O drivers, rather than providing its own. A price to pay for this simplicity is the increased overhead and the associated performance penalty; indeed, the I/O operations, page faults, and scheduling requests from a guest OS are not handled directly by the hypervisor, instead they are passed to the host OS. Performance, as well as the challenges to support complete isolation of VMs, make this solution less attractive for servers in a cloud computing environment. User-mode Linux is an example of a hosted VM.

Processor virtualization is not without its own challenges. As pointed out in [100] services provided by a VM “operate below the abstractions provided by the guest OS... it is difficult to provide a service that checks file system integrity without the knowledge of on-disk structure.” The hypervisors discussed in Section 4.4 manage resource sharing among the VMs running on a physical system.

5.4 Full virtualization and paravirtualization

In 1974, Popek and Goldberg gave a set of sufficient conditions for a computer architecture to support virtualization and allow a hypervisor to operate efficiently. Their crisp description of these conditions in [402] is a major contribution to the field:

1. A program running under the hypervisor should exhibit a behavior essentially identical to that demonstrated when running on an equivalent machine directly.
2. The hypervisor should be in complete control of the virtualized resources.
3. A statistically significant fraction of machine instructions must be executed without the intervention of the hypervisor.

One way to identify an architecture suitable for a VM is to distinguish *sensitive machine instructions* that require special precautions at execution time from the ones that can be executed without special precautions. In turn, sensitive instructions are: (i) *control sensitive*, i.e., instructions that attempt to change either the memory allocation, or operate in kernel mode; and (ii) *mode sensitive*, i.e., instructions whose behavior is different in kernel mode.

An equivalent formulation of the conditions for efficient virtualization can be based on this classification of machine instructions: *A hypervisor for a third or later generation computer can be constructed if the set of sensitive instructions is a subset of the privileged instructions of that machine.* To handle nonvirtualizable instructions, one could resort to two strategies:

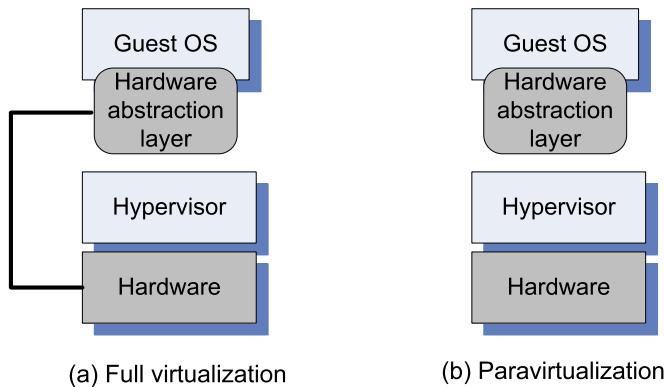
Binary translation. The hypervisor monitors the execution of guest operating systems; nonvirtualizable instructions executed by a guest OS are replaced with other instructions.

Paravirtualization. The guest OS is modified to use only instructions that can be virtualized.

There are two basic approaches to processor virtualization, see Fig. 5.2: *full virtualization* when each VM runs on an exact copy of the actual hardware; and *paravirtualization* when each VM runs on a slightly modified copy of the actual hardware. Paravirtualization is often adopted for several reasons: (1) some aspects of the hardware cannot be virtualized; (2) has better performance; and (3) presents a simpler interface.

Full virtualization requires a virtualizable architecture. The hardware is fully exposed to the guest OS which runs unchanged, and it is necessary to ensure that this execution mode is efficient. Systems such as the VMware EX Server support full virtualization on x86 architecture and have to address several problems, including the virtualization of the MMU. Privileged x86 instructions executed by a guest OS fail silently, thus, traps must be inserted whenever privileged instructions are issued by a guest OS. The system must also maintain shadow copies of system control structures, such as page tables, and trap every event affecting the state of these control structures. Therefore, the overhead of many operations is substantial.

Computer architectures such as x86 are not easily virtualizable as we shall see in Section 5.5. Paravirtualization is the alternative, though it has its own problems. Paravirtualization requires modifi-

**FIGURE 5.2**

(a) Full virtualization requires the hardware abstraction layer of the guest OS to have some knowledge about the processor architecture. The guest OS runs unchanged, thus, this virtualization mode is more efficient than paravirtualization. (b) Paravirtualization is used when the processor architectures is not easily virtualizable. In this case the hardware abstraction layer of the guest OS does not have knowledge about the hardware. The guest OS is modified to run under the hypervisor and must be ported to individual hardware platforms.

cations to a guest OS. Moreover, the code of the guest OS has to be ported to each individual hardware platform. Xen [48] and Denali [514] are based on paravirtualization.

Generally, the virtualization overhead negatively affects the performance of applications running under a VM. Sometimes, an application running under a VM can perform better than one running under a classical OS. This is the case of *cache isolation* when the cache is divided among VMs. In this case it is beneficial to run workloads competing for cache in two different VMs [449]. Often, the cache is not equally partitioned among processes running under a classical OS, and one process may use the cache space better than the other. For example, in the case of two processes, one write-intensive and the other read-intensive, the cache may be aggressively filled by the first.

The I/O performance of applications running under a VM depends on factors, such as: the disk partition used by the VM, the CPU utilization, the I/O performance of the competing VMs, and the I/O block size. The discrepancies between the optimal choice and the default ones on a Xen platform are between 8% and 35% [449].

5.5 Hardware support for virtualization

In early 2000, it became obvious that hardware support for virtualization was necessary, and Intel and AMD started working on the first generation virtualization extensions of the x86² architecture. In 2005,

² x86-32, i386, x86, and IA-32 refer to the Intel CISC-based instruction architecture, now supplanted by x86-64, which supports vastly larger physical and virtual address spaces. The x86-64 specification is distinct from Itanium, initially known as IA-64 architecture.

Intel released two Pentium 4 models supporting *VT-x*, and in 2006, AMD announced Pacifica and then several Athlon 64 models.

The Virtual Machine Extension (VMX) was introduced by Intel in 2006, and AMD responded with the Secure Virtual Machine (SVM) instruction-set extension. The *Virtual Machine Control Structure* (VMCS) of VMX tracks the host state, and the guest VMs as control is transferred between them. Three types of data are stored in VMCS:

- *Guest state*. Holds virtualized CPU registers (e.g., control registers or segment registers) automatically loaded by the CPU when switching from kernel mode to guest mode on *VMEntry*.
- *Host state*. Data used by the CPU to restore register values when switching back from guest mode to kernel mode on *VMExit*.
- *Control data*. Data used by the hypervisor to inject events, such as exceptions or interrupts into VMs and to specify which events should cause a *VMExit*; it is also used by the CPU to specify the *VMExit* reason.

VMCS is shadowed in hardware to overcome the performance penalties of nested hypervisors discussed in Section 5.10. This enables the guest hypervisor to access VMCS directly, without disrupting the root hypervisor in the case of nested virtualization. VMCS shadow access is almost as fast as a nonnested hypervisor environment. VMX includes several instructions [253]:

VMXON—enter vmx operation;
VMXOFF—leave vmx operation;
VMREAD—read from the VMCS;
VMWRITE—write to the VMCS;
VMCLEAR—clear VMCS;
VMPTRLD—load VMCS pointer;
VMPTRST—store VMCS pointer;
VMLAUNCH/VMRESUME—launch or resume a VM; and
VMCALL—call to the hypervisor.

A 2006 paper [361] analyzes the challenges to virtualizing Intel architectures and presents VT-x and VT-i, virtualization architectures for x86 and Itanium architectures, respectively. Software solutions at that time addressed some of the challenges, but hardware solutions could not only improve performance but also security and, at the same time, simplify the software systems. The problems faced by virtualization of x86 architecture are:

1. *Ring deprivileging*. x86 architecture provides four protection rings, 0–3. A hypervisor forces a guest VM, including an OS, to run at a privilege level greater than 0. Two solutions are then possible:
 - 1.1. The (0/1/3) mode when the hypervisor, the guest OS, and the application run at privilege levels 0, 1 and 3, respectively; this mode is not feasible in 64-bit mode, as we shall see shortly.
 - 1.2. The (0/3/3) mode when the hypervisor, a guest OS, and applications run at privilege levels 0, 3, and 3, respectively.

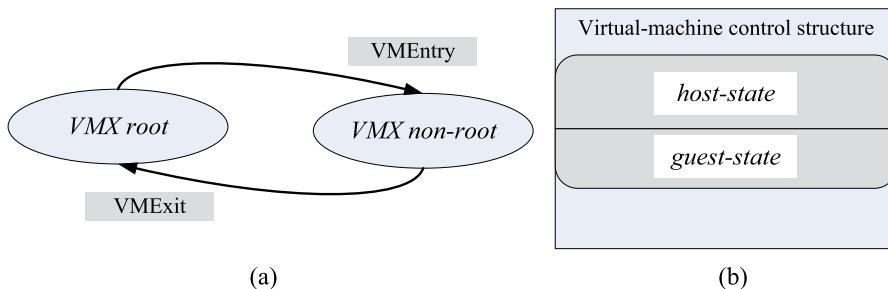
2. *Ring aliasing.* Such problems are created where a guest OS is forced to run at a privilege level other than that it was originally designed for. For example, when the CS register³ is *PUSHed*, the current privilege level in the CR is stored on the stack [361].
3. *Address space compression.* A hypervisor uses parts of the guest address space to store several system data structures such as the interrupt-descriptor table and the global-descriptor table. Such data structures must be protected, but the guest software must have access to them.
4. *Nonfaulting access to privileged state.* LGDT, SIDT, SLDT, and LTR instructions load registers GDTR, IDTR, LDTR, and TR, respectively, and can only be executed by software running at privileged level 0 because they point to data structures that control CPU operation. When executed at a privilege level other than 0, the four instructions *fail silently*, and a guest OS does not realize that any of these instructions has failed.
5. *Guest system calls.* Two instructions, SYSENTER and SYSEXIT, support low-latency system calls. The first causes a transition to privilege level 0, while the second causes a transition from privilege level 0 and fails if executed at a level higher than 0. The hypervisor must then emulate every guest execution of either of these instructions and that has a negative impact on performance.
6. *Interrupt virtualization.* In response to a physical interrupt, the hypervisor generates a “virtual interrupt” and delivers it later to the target guest OS. But every OS has the ability to mask interrupts,⁴ thus the virtual interrupt could only be delivered to the guest OS when the interrupt is not masked. Keeping track of all guest OS attempts to mask interrupts greatly complicates the hypervisor and increases the overhead.
7. *Access to hidden states.* Elements of the system state, e.g., descriptor caches for segment registers, are hidden; there is no mechanism for saving and restoring the hidden components when there is a context switch from one VM to another.
8. *Ring compression.* Paging and segmentation are the two mechanisms to protect hypervisor code from being overwritten by guest OS and applications. Systems running in 64-bit mode can only use paging, but paging does not distinguish between privilege levels 0, 1, and 2, thus the guest OS must run at privilege level 3, the so-called (0/3/3) mode. Privilege levels 1 and 2 cannot be used, thus the name ring compression.
9. *Frequent access to privileged resources increases hypervisor overhead.* The task-priority register (TPR) is frequently used by a guest OS; the hypervisor must protect the access to this register and trap all attempts to access it. That can cause a significant performance degradation.

A major architectural enhancement provided by the VT-x is the support for two modes of operation and a new data structure, Virtual Machine Control Structure (VMCS), including *host-state* and *guest-state* areas, see Fig. 5.3: *VMX root*: intended for hypervisor operations, and very close to the x86 without VT-x and *VMX nonroot*: intended to support a VM.

When executing a *VMEntry* operation, the processor state is loaded from the *guest-state* of the VM scheduled to run; then, the control is transferred from the hypervisor to the VM. A *VMExit* saves the processor state in the *guest-state* area of the running VM; it loads the processor state from the *host-state* area and finally transfers control to the hypervisor.

³ The x86 architecture supports memory segmentation with a segment size of 64 K. The CR (code-segment register) points to the code segment. *MOV*, *POP*, and *PUSH* instructions serve to load and store segment registers, including CR.

⁴ The interrupt flag, IF in the EFLAGS register is used to control interrupt masking.

**FIGURE 5.3**

(a) Two VT-x operation modes and the two transitions from one to the other; (b) VMCS includes *host-state* and *guest-state* areas controlling *VM entry* and *VM exit* transitions.

All *VMExit* operations use a common entry point to the hypervisor. Each *VMExit* operation saves the reason for the exit and eventually some qualifications in VMCS. Some of this information is stored as bitmaps. For example, the *exception bitmap* specifies which one of 32 possible exceptions caused the exit. The *I/O bitmap* contains one entry for each port in a 16-bit I/O space.

The VMCS area is referenced with a physical address, and its layout is not fixed by the architecture, but can be optimized by a particular implementation. VMCS includes control bits that facilitate virtual interrupts implementation. For example, *external-interrupt exiting* when set causes the execution of a *VM exit* operation; moreover, the guest is not allowed to mask these interrupts. When the *interrupt window exiting* is set, a *VM exit* operation is triggered if the guest is ready to receive interrupts.

Processors based on two new virtualization architectures, VT-d⁵ and VT-c have been developed. The first supports the I/O Memory Management Unit (I/O MMU) virtualization and the second supports the network virtualization. Also known as *PCI pass-through*, the I/O MMU virtualization gives VMs direct access to peripheral devices. VT-d supports:

- DMA address remapping: address translation for device DMA transfers.
- Interrupt remapping: isolation of device interrupts and VM routing.
- I/O device assignment: devices can be assigned by an administrator to a VM in any configuration.
- Reliability features: it reports and records DMA and interrupt errors that may otherwise corrupt memory and impact VM isolation.

5.6 QEMU

QEMU (Quick EMULATOR) is a open-source machine emulator and virtualizer that emulates the host processor architecture through dynamic binary translation for several architectures, including x86-64, PowerPC, RISC-V, ARMv7, and ARMv8. It supports a set of hardware and device models for the host, enabling it to run a variety of guest operating systems. QEMU has four operating modes:

⁵ The corresponding AMD architecture is called AMD-Vi.

- a. *User-mode emulation.* Used primarily for fast cross-compilation and cross-debugging. Runs Linux or macOS code compiled with different instruction sets. Calls a subroutine to deal with endianness⁶ and 32/64 bit addressing mismatch. This process called *thunk* is used to inject an additional calculation into another subroutine; a thunk is primarily used to delay a calculation until its result is needed, or to insert operations at the beginning or the end of the other subroutine.
- b. *System emulation.* Supports virtual hosting of several VM on a physical system and emulates a full system including I/O devices. QEMU can boot several guest operating systems, including Linux, Solaris, Microsoft Windows, DOS, and BSD; it supports emulating the ISAs mentioned above.
- c. *KVM hosting.* Supports setting up and migration of KVM images. It emulates the hardware, but the execution of the guest is done by KVM at the request of QEMU.
- d. *Xen hosting.* Emulates the hardware, and the execution of the guest is done within Xen, totally hidden from QEMU.

QEMU can simulate multiple CPUs running as a symmetric multiprocessor. The VM can interface with different host hardware devices. Virtual disk images can be stored in a special format that only takes up as much disk space as the guest OS allows. The disk QCOW2 format allows creation of overlay images and allows reverting the emulated disk contents to an earlier state. For example, a stable image could hold a fresh install of an operating system known to work and be reverted to it when one of the overlay images become unusable due to a virus attack or some other cause.

QEMU can emulate network cards and share host connectivity using network address translation and can also connect to network cards of other instances of QEMU. QEMU does not depend on the presence of graphical output methods on the host system. Instead, it can allow one to access the screen of the guest OS via an integrated VNC server. A VNC (Virtual Network Server) transmits keyboard and mouse input from one computer to another over a network, relaying graphical-screen updates using the Remote Frame Buffer protocol (RFB).

QEMU does not require administrative rights to run, unless additional modules require it. For example, KQEMU, a Linux kernel module which speeds up emulation of x86 or x86-64 guests on platforms with the same CPU architecture, requires administrative privileges. There are several OEMS that support hot pluggable NUMA hardware.

QEMU emulates ARMv7 instruction set. Xilinx Cortex A9-based Zynq SoC includes models for: the ARM Cortex-A9 CPU, ARM Cortex-A9 MPCore, DDR Memory Controller, Static Memory Controller (NAND/NOR Flash), DMA Controller, Gigabit Ethernet Controller, USB Controller, and other systems. QEMU version 6.0.0, released in 2021, includes support for ARMv8.1-M “Helium” architecture and Cortex-M55 CPU and for ARMv8.4 TTST, SEL2, and DIT extensions.

5.7 Kernel-based Virtual Machine

Kernel-based Virtual Machine (KVM) [288] is a virtualization infrastructure of the Linux kernel.⁷ KVM was released as part of the 2.6.20 Linux kernel in 2007. KVM technology is implemented as two

⁶ Endianness determines if the least significant byte of a word to be stored in memory will go to lowest or the highest address of the assigned memory space.

⁷ Up-to-date information about KVM can be found at <http://www.linux-kvm.org>.

components: (i) the KVM-loadable module installed in the Linux kernel, which provides management of the virtualization hardware, exposing its capabilities through the `/proc` file system; and (ii) platform emulation, provided by a modified version of QEMU, which executes as a user-space process, coordinating with the kernel for guest operating system requests.

KVM was originally designed for x86-32 and has since been ported to x86-64, ARM, and several other architectures. Running multiple guest operating systems on x86 architecture was quite difficult before the introduction of VMX and SVM extensions of Intel architecture. These extensions enable the hypervisor to run within the privileged ring –1 and enable KVM to provide VMs with an execution environment nearly identical to the physical hardware. KVM executes the guest VM's instructions directly on the host. Each guest OS is isolated and runs in a different instance of the execution environment.

KVM provides hardware-assisted virtualization for several guest operating systems, including Linux, BSD, Solaris, Windows, and macOS. KVM provided paravirtualization support for Linux, OpenBSD, FreeBSD, NetBSD, and Windows guests using the *VirtIO* API. *Virtio* is the main platform for I/O virtualization in KVM. The idea behind *VirtIO* is to have a common framework for I/O virtualization for different hypervisors.

KVM supports hot plug vCPUs, dynamic memory management, and live migration. KVM inherits Linux kernel support for *cpu-hot plugging*; this means that one can enable or disable a CPU or a core without the need to reboot. This mechanism can be used to replace defective components and can be exploited for dynamic partitioning when running multiple Linux partitions. As the workloads change, it is useful to move CPUs from one partition to the next without rebooting or interrupting the workloads.

Live migration is the process of moving a running VM or an application between different physical machines without disconnecting the client or application. With live migration, the memory, storage, and network connectivity of the VM are transferred from the original guest machine to the destination.

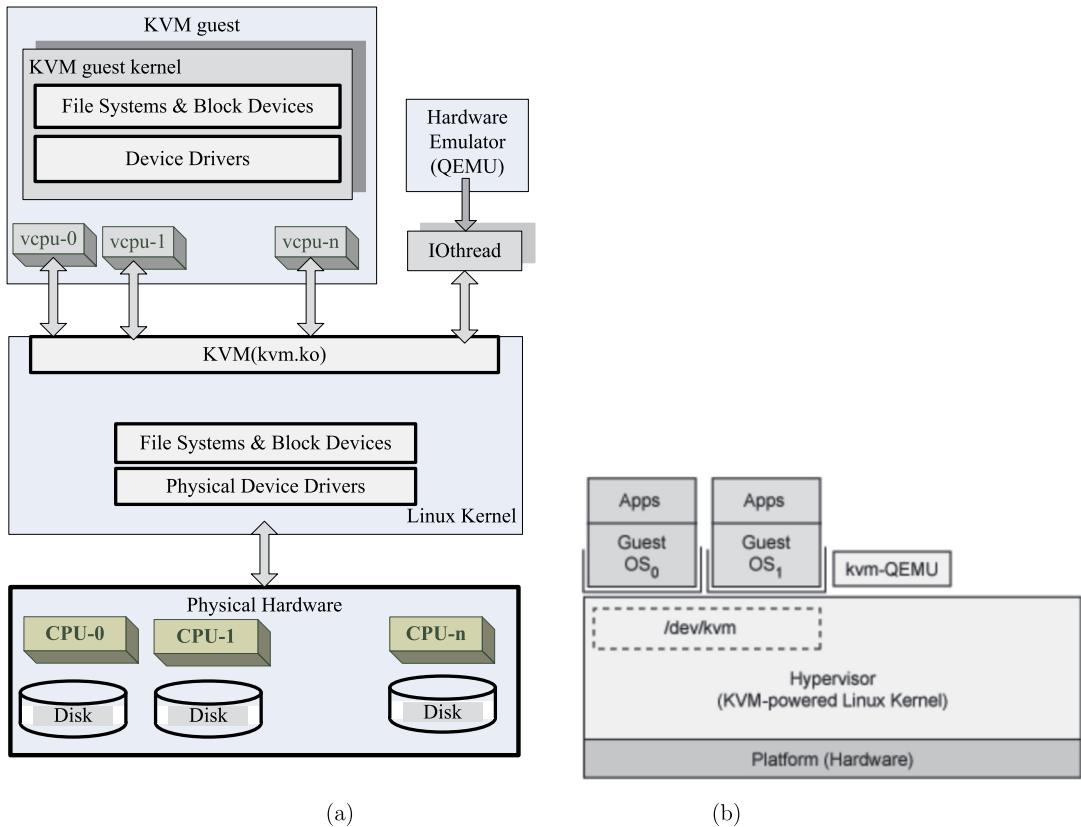
KVM converts Linux into a bare-metal hypervisor and inherits from the Linux kernel operating system-level components, such as a memory manager, process scheduler, input/output (I/O) stack, device drivers, security manager, and the network stack. Every VM is implemented as a regular Linux process, scheduled by the standard Linux scheduler, with dedicated virtual hardware like a network card, graphics adapter, CPU(s), memory, and disks, Fig. 5.4(b).

KVM does not run as a normal program inside Linux, but it relies on the Linux kernel infrastructure to run. Its organization is shown in Fig. 5.4.(a) As opposed to Xen which runs on the physical hardware, KVM runs inside Linux as a driver handling the new virtualization instructions exposed by hardware. KVM has several components:

1. A generic host kernel module exposing the architecture-independent functionality.
2. An architecture-specific kernel module for the host system.
3. A user-space emulation of the VM hardware that the guest OS runs on.
4. A guest OS performance-optimization additions.

Instead of emulating several hardware systems, KVM uses higher-level client applications such as QEMU, *crosvm*, or *Firecracker* for device emulation. *kvm-userspace* is a fork of the QEMU project; it short-circuits the emulation code to allow only x86-on-x86 and use the KVM API for running the guest OS on the host CPU. When the guest OS performs a privileged operation, the CPU exits and KVM takes over. If KVM itself can service the request, it runs it and then gives control back to the guest.

Firecracker is an open-source virtualization technology built for creating and managing secure, multitenant container and function-based services, such as AWS Lambda and AWS Fargate, a serverless

**FIGURE 5.4**

(a) KVM organization; KVM runs inside Linux as a driver handling the new virtualization instructions exposed by hardware; the *IOthread* generates requests on the guest's behalf to the host; it also handles events. (b) Multiple VMs running under KVM.

compute engine for containers working with Amazon Elastic Container Service (ECS) and Amazon Elastic Kubernetes Service (EKS). Firecracker uses KVM to manage microVMs.

KVM exposes the */dev/kvm* interface enabling a user-space host to: (a) setup the guest VM address space; (b) feed the guest simulated I/O; and (c) map the video display of the host. The host supplies a firmware image used by the guest to bootstrap into the host OS.

Several characteristics of KVM and other virtualization environments are displayed in Table 5.1. KVM has a number of important advantages over other virtualization environments:

1. Lower total cost of ownership and no vendor lock-in.
2. Excellent performance: apps run faster on KVM compared with other hypervisors.
3. The open-source advantage, the access source code, and the flexibility to integrate with anything.

Table 5.1 Virtualization software. Only a subset of architectures and operating systems are shown. Operating systems abbreviated as F, L, M, U, and W for FreeBSD, Linux, MacOS, Unix-like, and Windows, respectively. Software licenses are described at https://en.wikipedia.org/wiki/GNU_General_Public_License.

Name	Creator/Licences	Host CPU	Guest CPU	Host OS	Guest OS
KVM	Red Hat GPLV2	x86-64, ARM, Alpha	Same as host	L, F L, F, M, W	L, F W
QEMU	— GPL/LGPL /	x86-64, ARM Alpha	x86-64, ARM Alpha	L, F, M, W	Changes regularly
VMware ESX server	VMware / Proprietary	x86-64	x86-64	no host	L, F, W
VMware Fusion	VMware / Proprietary	x86-64	x86-64	M	L, F, W
VMWare Server	VMWare / Proprietary	x86-64	x86-64	L, W	L, F, W
Xen	Citrix / GNU, GPLv2	x86-64, ARM	Same as host	L, U	L, F, W
Xen Server	Citrix / GNU, GPLv2	x86-64, ARM	Same as host	No host	L, F, W

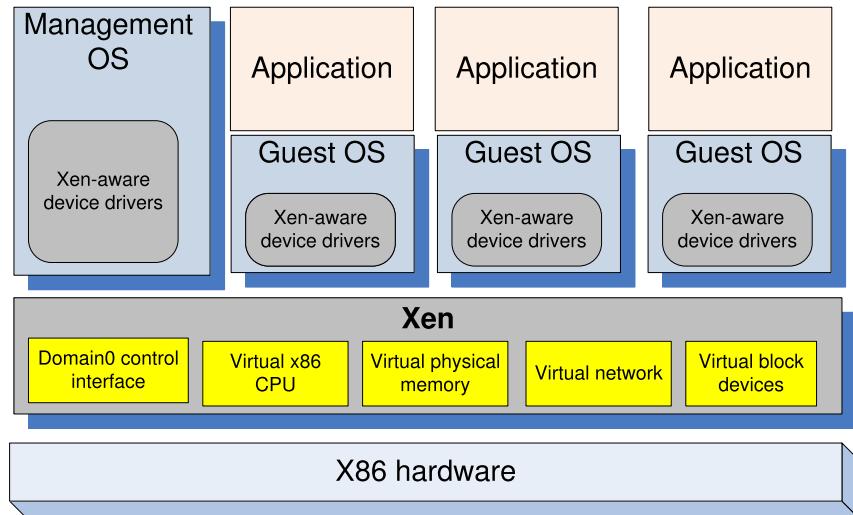
4. Cross-platform interoperability: users benefit from existing infrastructure investments.
5. Simplicity to create, start, stop, pause, migrate, and template a large number of VMs on a range of hardware and software.

KVM virtualization includes sVirt and SELinux (Security-Enhanced Linux) technologies developed to detect and prevent complex security threats. The stability, consistency, and compatibility of Linux is shared with the virtual layer since KVM has been part of the Linux kernel for a long time.

In conclusion, a hypervisor using Linux as the core has tangible benefits. An obvious one is the steady progression of Linux, a platform that continues to advance, from typical optimizations and bug fixes, scheduling, and memory-management innovations to support for multiple processor architectures. Another benefit is that one can take advantage of the Linux platform as an operating system in addition to a hypervisor. In addition to running multiple guest operating systems on the Linux hypervisor, one can run other traditional applications at that level.

5.8 Xen—a hypervisor based on paravirtualization

Xen is a hypervisor developed by the Computing Laboratory at the University of Cambridge, United Kingdom, in 2003. The goal of the Cambridge group was to design a hypervisor capable of scaling to about 100 VMs running standard applications and services without any modifications to the Application Binary Interface. Fully aware that x86 architecture does not support efficiently full virtualization, the designers of Xen opted for paravirtualization.

**FIGURE 5.5**

Xen for x86 architecture. The management OS dedicated to the execution of Xen control functions and privileged instructions resides in *Dom0*; guest operating systems and applications reside in *DomU*.

Since 2010, Xen has been a free software, developed by the community of users and licensed under the GNU General Public License (GPLv2). Several operating systems including Linux, Minix, NetBSD, FreeBSD, NetWare, and OZONE can operate as paravirtualized Xen guest operating systems running on x86, x86-64, Itanium, and ARM architectures.

We analyze next the original implementation of Xen for the x86 architecture discussed in [48]. The creators of Xen used the concept of *domain* (*Dom*) to refer to the ensemble of address spaces hosting a guest OS and address spaces for applications running under this guest OS. Each domain runs on a virtual x86 CPU. *Dom0* is dedicated to the execution of Xen control functions and privileged instructions, and *DomU* is a user domain, Fig. 5.5.

The most important aspects of Xen paravirtualization for virtual memory management, CPU multiplexing, and I/O device management are summarized in Table 5.2 [48]. Efficient management of TLB (Translation Look-aside Buffer), a cache for page table entries, requires either the ability to identify the OS and the address space of every entry or to allow software management of the TLB. Unfortunately, x86 architecture does not support either the tagging of TLB entries or the software management of the TLB. As a result, the address space switching when the hypervisor activates a different OS, requires a complete TLB flush. Flushing the TLB has a negative impact on performance.

The solution adopted was to load Xen in a 64-MB segment at the top of each address space and to delegate the management of hardware page tables to the guest OS with minimal intervention from Xen. The 64-MB region occupied by Xen at the top of every address space is not accessible or not remappable by the guest OS.

When a new address space is created, the guest OS allocates and initializes a page from its own memory, registers it with Xen, relinquishing control of the write operations to the hypervisor. Thus, a

Table 5.2 Paravirtualization strategies for virtual memory management, CPU multiplexing, and I/O devices for the original x86 Xen implementation.

Function	Strategy
Paging	A domain may be allocated discontinuous pages. A guest OS has direct access to page tables and handles pages faults directly for efficiency; page table updates are batched for performance and validated by Xen for safety.
Memory	Memory is statically partitioned between domains to provide strong isolation. <i>XenLinux</i> implements a <i>balloon driver</i> to adjust domain memory.
Protection	A guest OS runs at a lower priority level, in ring 1, while Xen runs in ring 0.
Exceptions	A guest OS must register with Xen a description table with the addresses of exception handlers previously validated.
System	To increase efficiency, a guest OS must install a “fast” handler
Interrupts	A lightweight event system replaces hardware interrupts; synchronous system calls from a domain to Xen use <i>hypercalls</i> and notifications are delivered using the asynchronous event system.
Multiplexing	A guest OS may run multiple applications.
I/O devices	Data is transferred using asynchronous I/O rings.
Disk access	Only <i>Dom0</i> has direct access to IDE and SCSI disks; all other domains access persistent storage through the Virtual Block Device (VBD) abstraction.

guest OS could only map pages it owns. On the other hand, a guest OS has the ability to batch multiple page-update requests to improve performance. A similar strategy is used for segmentation.

x86 Intel architecture supports four protection rings or privilege levels; virtually all OS kernels run at level 0, the most privileged one, and applications at level 3. In Xen the hypervisor runs at level 0, the guest OS at level 1, and applications at level 3.

Applications make system calls using the so-called *hypercalls* processed by Xen; privileged instructions issued by a guest OS are *paravirtualized* and must be validated by Xen. When a guest OS attempts to execute a privileged instruction directly, the instruction fails silently.

Memory is statically partitioned between domains to provide strong isolation. To adjust domain memory, *XenLinux* implements a *balloon driver* that passes pages between Xen and its own page allocator. Page faults are handled directly by the guest OS for efficiency.

Xen schedules individual domains using Borrowed Virtual Time (BVT), a work-conserving⁸ and low-latency wake-up scheduling algorithm discussed in Section 9.7. BVT uses a virtual-time warping mechanism to support low-latency dispatch to ensure timely execution whenever needed, for example, for timely delivery of TCP acknowledgments.

A guest OS⁹ must register with Xen a *description table* with the addresses of exception handlers for validation. Exception handlers are identical to the native x86 handlers; the only one that does not follow this rule is the page fault handler, which uses an extended stack frame to retrieve the faulty address because the privileged register *CR2*, where this address is found, is not available to a guest OS.

Each guest OS can validate and then register a “fast” exception handler executed directly by the processor without the interference of Xen. A lightweight event system replaces hardware interrupts;

⁸ A work-conserving scheduling algorithm does not allow the processor to be idle when there is work to be done.

⁹ In the original Xen implementation [48], a guest OS could be either *XenoLinux*, *XenoBSD*, or *XenoXP*.

notifications are delivered using this asynchronous event system. Each guest OS has a timer interface and is aware of “real” and “virtual” time.

XenStore is a *Dom0* process supporting a system-wide registry and naming service. It is implemented as a hierarchical key-value stor. A *watch* function of the process informs listeners of changes of the key in the storage they have subscribed to. XenStore communicates with guest VMs via shared memory using *Dom0* privileges, rather than grant tables.

Toolstack is another *Dom0* component responsible for creating, destroying, and managing the resources and privileges of VMs. To create a new VM, a user provides a configuration file describing memory and CPU allocations, as well as device configuration. Then, *Toolstack* parses this file and writes this information in *XenStore*. *Toolstack* takes advantage of *Dom0* privileges to map guest memory, to load a kernel and virtual BIOS, and to set up initial communication channels with *XenStore* and with virtual console when a new VM is created.

Xen defines abstractions for networking and I/O devices. *Split drivers* have a frontend in the DomU and the backend in *Dom0*; the two communicate via a ring in shared memory. Xen enforces access control for the shared memory and passes synchronization signals. Access Control Lists (ACLs) are stored in the form of *grant tables*, with permissions set by the owner of the memory.

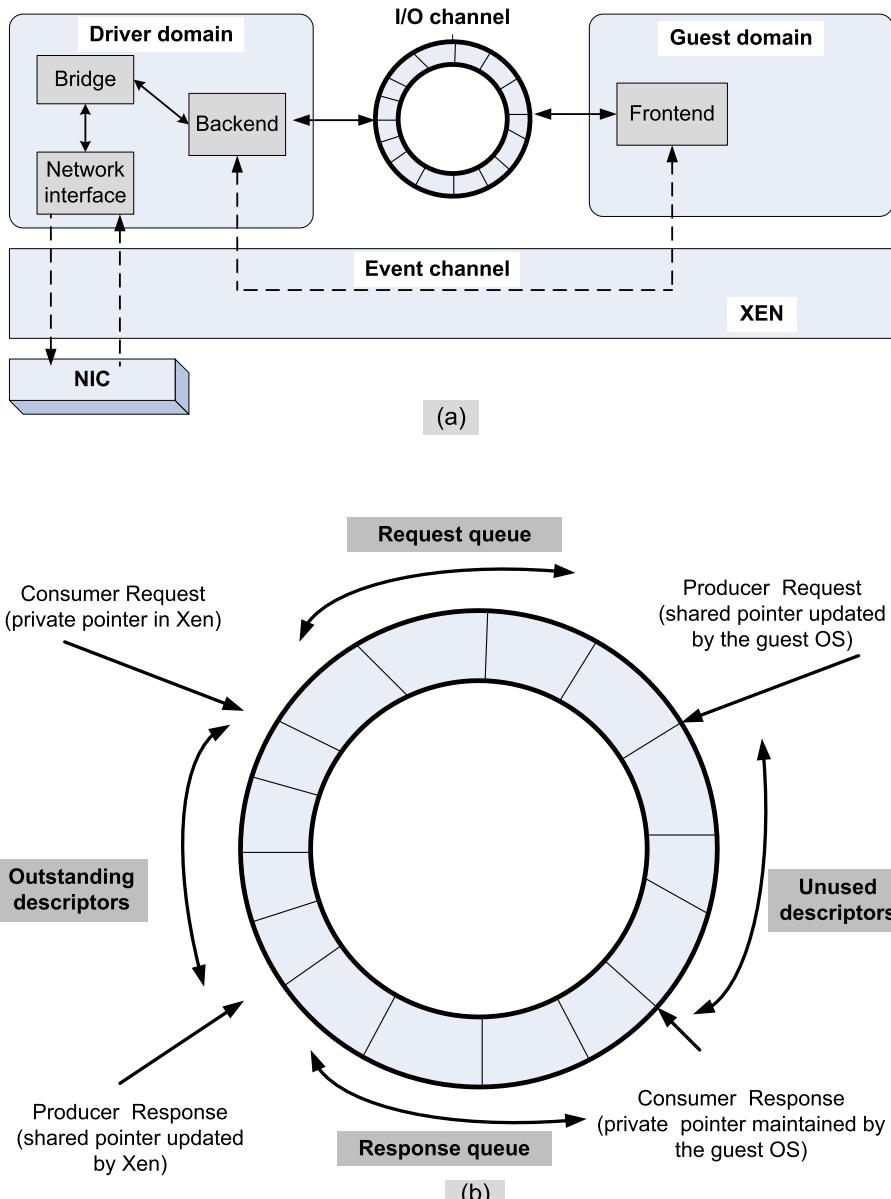
Data for I/O and network operations moves vertically through the system, very efficiently, using a set of I/O rings, see Fig. 5.6. A *ring* is a circular queue of descriptors allocated by a domain and accessible within Xen. Descriptors do not contain data, the data buffers are allocated off-band by the guest OS. Memory committed for I/O and network operations is supplied in a manner designed to avoid “crosstalk,” and the I/O buffers holding the data are protected by preventing page faults of the corresponding page frames.

Each domain has one or more Virtual Network Interfaces (VNIs), which support the functionality of a network interface card. A VNI is attached to a Virtual Firewall-Router (VFR). Two rings of buffer descriptors, one for packet sending and one for packet receiving, are supported. To transmit a packet, a guest OS enqueues a buffer descriptor to the send ring, then Xen copies the descriptor and checks safety, and finally copies only the packet header, not the payload, and executes the matching rules.

Rules of the form (*< pattern >*, *< action >*) require the *action* to be executed if the *pattern* is matched by the information in the packet header. The rules can be added or removed by *Dom0*; they ensure the demultiplexing of packets based on the destination IP address and port and, at the same time, prevent spoofing of the source IP address. *Dom0* is the only one allowed to access directly the physical IDE or SCSI disks. All domains other than *Dom0* access persistent storage through a Virtual Block Device (VBD) abstraction created and managed under the control of *Dom0*.

Xen includes a device emulator, QEMU, to support unmodified commodity operating systems. QEMU is a machine emulator; it runs unmodified OS images and emulates those architecture’s instructions for the host architecture it runs on. The project had several devices already emulated for the x86 architecture, including the chipset, network cards, and display adapters. QEMU emulates a DMA¹⁰

¹⁰ Direct Memory Access (DMA) is specialized hardware that allows the I/O subsystems to access the main memory without CPU intervention. It can also be used for memory-to-memory copying and can offload expensive memory operations, e.g., scatter-gather operations, from a CPU to dedicated a DMA engine.

**FIGURE 5.6**

Xen zero-copy semantics for data transfer using I/O rings. (a) The communication between a guest domain and the driver domain over an I/O and an event channel; NIC is the Network Interface Controller. (b) The circular ring of buffers.

and can map any page of the memory in a DomU. Each VM has its own instance of QEMU, which can run either as a *Dom0* process, or as a process of the VM.

Xen, initially released in 2003, has undergone significant changes in 2005, when Intel released the VT-x processors. In 2006, Xen was adopted by Amazon for its EC2 service, and in 2008, Xen, running on Intel's VT-d, passed the ACPI S3¹¹ test. Xen support for *Dom0* and *DomU* was added to the Linux kernel in 2011.

In 2008, the PCI pass-through was incorporated for Xen running on VT-d architectures. The PCI¹² pass-through allows a PCI device, be it a disk controller, Network Interface Card (NIC),¹³ graphic card, or Universal Serial Bus (USB), to be assigned to a VM. This avoids the overhead of copying and allows setting up of a *Driver Domain* to increase security and system reliability. A guest OS can exploit this facility to access the 3D acceleration capability of a graphics card. The BDF¹⁴ of a device must be known for pass-through.

An analysis of VM performance for I/O-bound applications under Xen is reported in [405]. Two web servers, each running under a different VM, share the same server running Xen; the workload generator sends requests for files of fixed size ranging from 1 KB to 100 KB. When the file size increases from 1 KB to 10 KB and to 100 KB, the CPU utilization, the throughput, data rate, and the response times are (97.5%, 70.44%, 44.4%), (1 900, 1 104, 1 112) requests/sec, (2 018, 11 048, 11 208) KBps, and (1.52, 2.36, 2.08) msec, respectively. From the first group of results, we see that for files 10 KB or larger, the system is I/O bound. The second set of results shows that the throughput measured in requests/second decreases by less than 50% when the system becomes I/O bound, but the data rate increases by a factor of five over the same range. The response time increases only about 10% when the file size increases by two orders of magnitude.

The paravirtualization strategy in Xen is different from the one adopted by the group at the University of Washington, the creators of the Denali system [514]. Denali was designed to support a number of VMs running network services one or more orders of magnitude larger than Xen. The design of the Denali system did not target existing application binary interfaces; it does not support some features of potential guest operating systems; for example, it does not support segmentation. Denali does not support application multiplexing, running multiple applications under a guest OS, while Xen does.

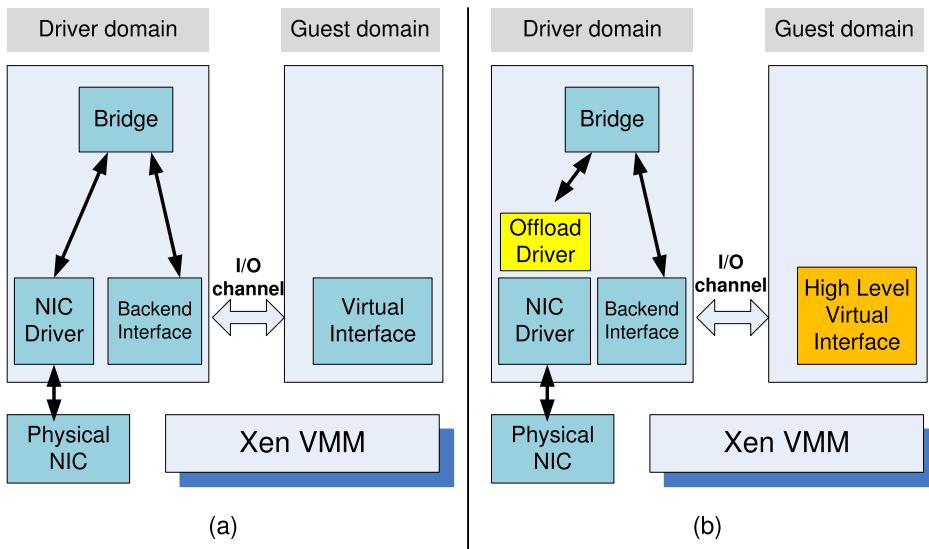
Finally, a few words regarding the complexity of porting commodity operating systems to Xen. It is reported that a total of about 3 000 lines of Linux code, or 1.36%, had to be modified; for Windows XP, this figure is 4 620, or about 0.04%, [48].

¹¹ Advanced Configuration and Power Interface (ACPI) specification is an open standard for device configuration and power management by the OS. It defines four Global “Gx” states and six Sleep “Sx” states. “S3” is referred to as Standby, Sleep, or Suspend to RAM.

¹² PCI stands for Peripheral Component Interconnect and describes a computer bus for attaching hardware devices to a computer. The PCI bus supports the functions found on a processor bus, but in a standardized format independent of any particular processor. The OS queries all PCI buses at startup time to identify the devices connected to the system and the memory space, I/O space, interrupt lines, and so on needed by each device present.

¹³ A Network Interface Controller is the hardware component connecting a computer to a LAN.

¹⁴ BDF stands for Bus.Device.Function, and it is used to describe PCI devices.

**FIGURE 5.7**

Xen network architecture: (a) Original architecture; (b) Optimized architecture.

5.9 Optimization of network virtualization in Xen 2.0

A hypervisor introduces a significant network communication overhead. For example, it is reported that the CPU utilization of a VMware Workstation 2.0 system running Linux 2.2.17 was 5 to 6 times higher than that of the native system (Linux 2.2.17) while saturating a 100 Mbps network [462]. In other words, the hypervisor executes a much larger number of instructions, 5 to 6 times larger, to saturate the network, while handling the same amount of traffic as the native system.

Similar overheads are reported for other hypervisors and, in particular, for Xen 2.0 [345,346]. To understand the sources of network overhead, we examine the basic network architecture of Xen, see Fig. 5.7(a). Privileged operations, including I/O, are executed by *Dom0* on behalf of a guest OS. In this context, we shall refer to *Dom0* as the *driver domain*. The *driver domain* is called in to execute networking operations on behalf of the *guest domain*. It uses the native Linux driver for the NIC (Network Interface Controller), which in turn, communicates with the physical NIC, also called the network adapter. The *guest domain* communicates with the *driver domain* through an I/O channel, see Section 5.8. More precisely, the guest OS in the guest domain uses a virtual interface to send and receive data to/from the backend interface in the driver domain.

A *bridge* uses broadcast communication to identify the MAC address¹⁵ of a destination system. Once identified, this address is added to a table. A bridge uses the link layer protocol to send a packet

¹⁵ A Media Access Control (MAC) address is a unique identifier permanently assigned to a network interface by the manufacturer.

Table 5.3 A comparison of send and receive data rates for a native Linux system, the Xen driver domain, an original Xen guest domain, and an optimized Xen guest domain.

System	Receive data rate (Mbps)	Send data rate (Mbps)
Linux	2 508	3 760
Xen driver	1 728	3 760
Xen guest	820	750
optimized Xen guest	970	3 310

to the proper MAC address, rather than broadcast it when the next packet for the same destination arrives.

The bridge in the driver domain performs a multiplexing/demultiplexing function. Packets received from the NIC are demultiplexed and sent to the VMs running under the hypervisor. Similarly, packets arriving from multiple VMs have to be multiplexed into a single stream before being transmitted to the network adaptor. In addition to bridging, Xen supports IP routing based on network address translation.

Table 5.3 shows the ultimate effect of this longer processing chain for the Xen hypervisor, as well as the effect of optimizations [346]; the receiving and sending rates from a guest domain are roughly 30% and 20%, respectively, of the corresponding rates of a native Linux application. Packet multiplexing/demultiplexing accounts for about 40% and 30% of the communication overhead for the incoming traffic and for the outgoing traffic, respectively.

The Xen network optimization discussed in [346] covers optimization of: (i) the virtual interface; (ii) the I/O channel; and (iii) the virtual memory. The effects of these optimizations are significant for the send-data rate from the optimized Xen guest domain, an increase from 750 to 3 310 Mbps and rather modest for the receive-data rate, 970 versus 820 Mbps.

We next examine each optimization area and start with the virtual interface. There is a tradeoff between generality and the flexibility, on one hand, and the performance on the other hand. The original virtual network interface provides the guest domain with a simple low-level network interface abstraction supporting sending and receiving primitives.

This design supports a wide range of physical devices attached to the driver domain but does not take advantage of the capabilities of some physical NICs, such as checksum offload, e.g., TSO,¹⁶ and scatter/gather DMA support. These features are supported by the High Level Virtual Interface of the optimized system, Fig. 5.7(b).

The next target of the optimization effort is the communication between the guest domain and the driver domain. Rather than copying a data buffer holding a packet, each packet is allocated space in a new page, and then the physical page containing the packet is remapped onto the target domain; for example, when a packet is received, the physical page is remapped to the guest domain. The optimization is based on the observation that there is no need to remap the entire packet.

For example, when sending a packet, the network bridge needs only to know the MAC header of the packet. As a result of this, the optimized implementation is based on an “out-of-band” channel used

¹⁶ TSO stands for TCP segmentation offload. This option enables the network adapter to compute the TCP checksum on *transmit* and *receive*, and to save the host CPU the overhead for computing the checksum; large packets have larger savings.

by the guest domain to provide the bridge with the packet MAC header. This strategy contributed to a better than four times increase in the send data rate compared with the nonoptimized version.

The third optimization covers virtual memory. The virtual memory in Xen 2.0 takes advantage of the *superpage* and *global page mapping* hardware features available on Pentium and Pentium Pro processors. A superpage increases the granularity of the dynamic address translation; a superpage entry covers 1 024 pages of physical memory, and the address translation mechanism maps a set of contiguous pages to a set of contiguous physical pages. This helps reduce the number of TLB misses.

All pages of a superpage belong to the same guest OS. When new processes are created, the guest OS must allocate read-only pages for the page tables of the address spaces running under the guest OS. This forces the system to use traditional page mapping, rather than the superpage mapping. The optimized version on network virtualization uses a special memory allocator to avoid this problem.

5.10 Nested virtualization

Nested virtualization describes a system organization when a *guest hypervisor* runs inside a VM, which is itself running under a *host hypervisor*. Fig. 5.8(a) illustrates an instance of nested virtualization where a KVM is the host hypervisor supporting resource sharing among three VMs. Two of the three VMs run guest hypervisors, Xen, and VMware's ESXi, and the third VM runs Windows. There are two VM running under guest hypervisors, one runs Linux under Xen and another runs Windows under ESXi.

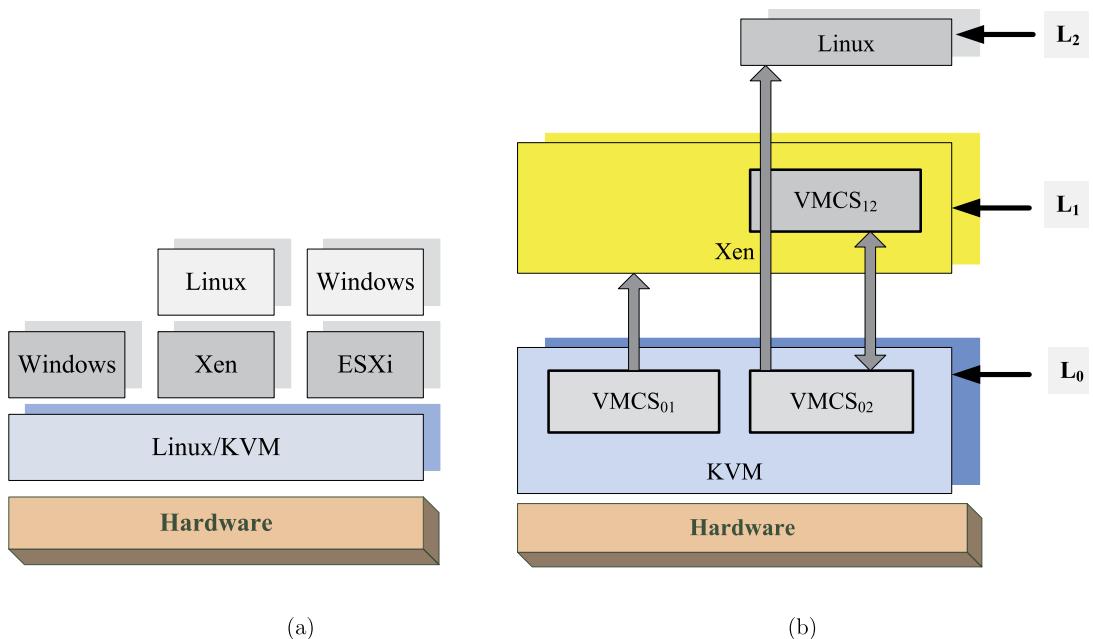
Nested virtualization is useful for experimenting with server setup or testing configurations. Nested virtualization enables IaaS users to run their own hypervisor as a VM. Nested virtualization can be also used for live migration of hypervisors together with their guest VMs for load balancing, for hypervisor-level protection, and for supporting other security mechanisms. Another use of nested virtualization is to experiment with cloud interoperability alternatives.

The x86 virtualization is based on the *trap and emulate* model. This model requires every sensitive instruction executed by either a guest hypervisor or OS to be handled by the most privileged hypervisor. Nested virtualization incurs a substantial performance price, unless the switching between the levels of the virtualization stack is optimized. It is thus not surprising that nested virtualization is not supported by many hypervisors, and not all operating systems can nest successfully with all hypervisors.

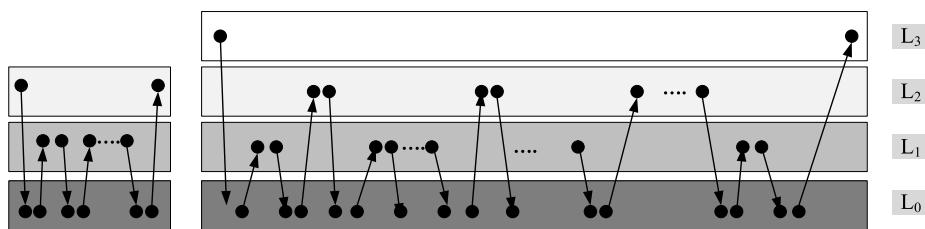
Nested virtualization is supported differently by Intel and AMD processors. Consider for example the Intel version discussed in [127] and illustrated in Fig. 5.8(b). In this example, KVM runs at level L₀ and controls the allocation of all resources. Xen runs at level L₁, and KVM uses VMCS₀₁ for the VM running Xen.

When Xen executes a *vmlaunch* operation to start a new VM, (see Section 5.5), a new VM control structure, VMCS₁₂ is created. Then, *vmlaunch* traps to L₀ and L₀ merges VMCS₀₁ with VMCS₁₂ and creates VMCS₀₂ to run Linux at level L₂. When an application running under Linux at level L₂ makes a system call, or when Linux itself executes a privileged instruction, L₂ traps and KVM decides whether to handle the trap itself or to forward it to Xen at level L₁ and, eventually, Xen resumes execution.

Nested virtualization is limited by the hardware support. When the hardware supports *multilevel nested virtualization*, each hypervisor handles all traps caused by sensitive instructions of guest hypervisors running directly above of it. Multilevel nested virtualization is supported by the IBM System Z architecture [381]. Intel and AMD processors support only *single-level nested virtualization*. This

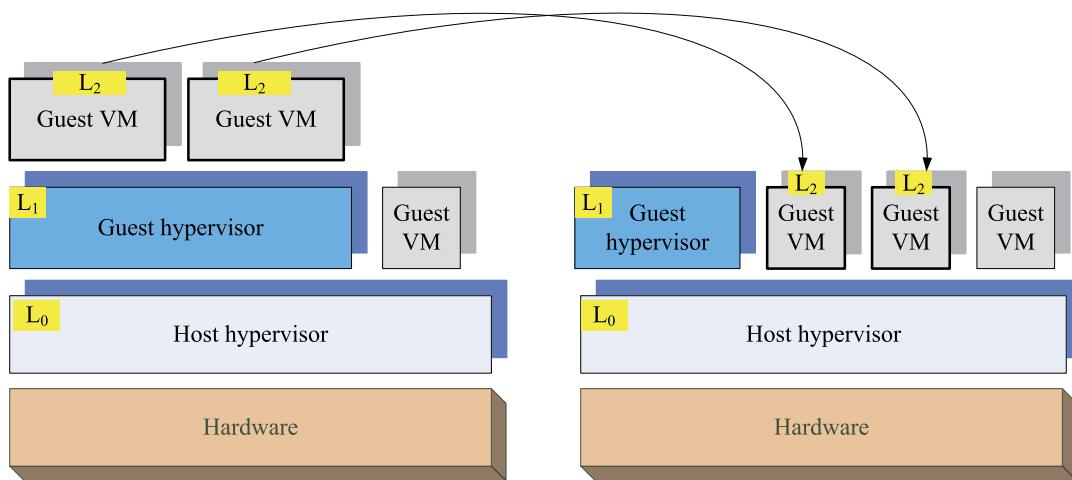
**FIGURE 5.8**

Nested virtualization [127]. (a) KVM allows three VM to run concurrently. Two VMs run hypervisors Xen and ESXi, and the third runs Windows. A VM runs Linux under Xen, and another VM runs Windows under ESXi. (b) Intel-supported nested virtualization. KVM runs at level L₀, Xen runs at level L₁, and KVM uses VMCS₀₁ for the VM running Xen.

**FIGURE 5.9**

Nested virtualization with single-level hardware virtualization support. A trap is handled by the L₀ trap handler regardless of the hypervisor where a trap occurs. Nested traps for: (Left) Two-level, L₀, L₁, and L₂ nested hypervisor(s); and (Right) Three-level, L₀, L₁, L₂, and L₃ nested hypervisor(s).

implies that the host hypervisor, the one running directly above the hardware and managing all system resources, handles all trapped instructions as shown in Fig. 5.9.

**FIGURE 5.10**

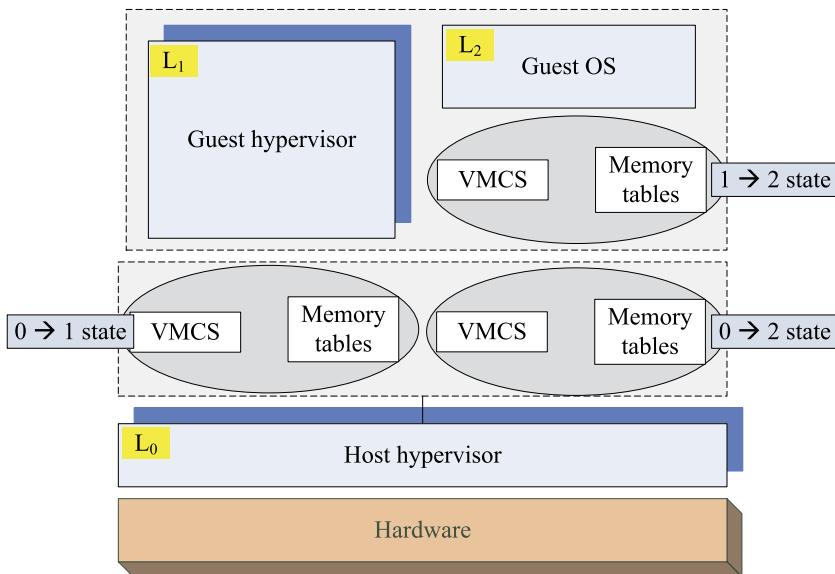
Multiple virtualization levels on the left are multiplexed into the single hardware virtualization level on the right, as described in [57]. A VMX instruction used by a guest hypervisor running in guest mode at level L_i is trapped and translated by the host hypervisor at level L_0 running in kernel mode into one that can be used to a VM at level L_{i+1} .

An in-depth discussion of the intricacies of nested virtualization on x86 architecture can be found in a paper describing the Turtle project from IBM Israel [57]. A guest hypervisor cannot use the hardware virtualization support because the x86 provides only a single-level hardware virtualization support. The aim of the project was to show that a 6–8% overhead is feasible for “unmodified binary-only hypervisors executing nontrivial workloads.”

VMX instructions can only be successfully executed in kernel mode. A guest hypervisor at level L_i operates in guest mode, and, whenever it executes a VMX instruction to launch a level L_{i+1} guest, the instruction is trapped and handled at level L_0 . Trapping execution exceptions enables the host hypervisor at level L_0 running in kernel mode to emulate VMX instruction executed by guest hypervisors at level L_i . This mechanism supports a critical technique for increasing efficiency of nested virtualization, the multiplexing multiple hypervisors, as shown in Fig. 5.10.

As long as the host hypervisor at level L_0 emulates faithfully the VMX instruction set, a guest hypervisor at level L_1 cannot distinguish if it is running directly on the hardware or not. It follows that the guest hypervisor at level L_1 can launch VMs using the standard mechanisms. The guest hypervisor at level L_1 does not run at the highest privileged level, and the action of starting a VM is trapped and handled by the trap handler at level L_0 . The specification of the new VM is then translated by the host hypervisor at level L_0 running in kernel mode into one that can be used to run L_2 directly on the hardware. This translation includes converting L_1 physical addresses to the physical address space of L_0 .

The guest hypervisor can use the same technique to give another guest hypervisor at level L_2 the same illusion, namely, that it is running directly on the hardware. The process can be extended, a hypervisor at level L_i giving the illusion that the one at level L_{i+1} is running directly on the hardware.

**FIGURE 5.11**

VMX extension for nested virtualization as described in [57].

L_0 , the host supervisor, manages resources allocated to L_1 , a guest hypervisor and L_2 , a guest OS, using $VMCS_{0 \rightarrow 1}$ and $VMCS_{0 \rightarrow 2}$ environment specification, respectively. L_1 creates $VMCS_{1 \rightarrow 2}$ within its own virtualized environment, and the processor uses it to emulate VMX for L_1 , as illustrated by Fig. 5.11. Switching from one level to another is emulated. For example, when an *VMExit* occurs while L_2 is running there are two possible paths:

- When an external interrupt, a nonmaskable interrupt, or any trappable event specified in $VMCS_{0 \rightarrow 2}$ that was not specified in $VMCS_{1 \rightarrow 2}$ occurs, L_0 handles the event and then L_2 execution is resumed.
- Trappable events specified in $VMCS_{1 \rightarrow 2}$ are handled by L_1 . The host hypervisor L_0 forwards the event to L_1 by copying $VMCS_{0 \rightarrow 2}$ fields updated by the processor to $VMCS_{1 \rightarrow 2}$ and then resuming L_1 . This makes the L_1 hypervisor believe there was a *VMExit* directly from L_2 to L_1 , handle the event and then resume L_2 by executing *VMLAUNCH* or *VMRESUME*, emulated by L_0 .

Another complication of nested virtualization is that the MMU must be virtualized to allow a guest hypervisor to translate guest virtual addresses to guest physical addresses. A multidimensional paging for multiplexing the three needed translation tables onto the two available in hardware is described in [57].

5.11 A trusted kernel-based virtual machine for ARMv8

Advanced RISC Machine (ARM) processors are widely used in mobile devices such as smartphones, tablets, and laptops. ARM processors are also used in embedded systems connected to the IoT. Such

systems require an increased level of security, therefore, it is not surprising that the latest generation of the ubiquitous ARM processors support the Trusted Execution Environment (TEE).

TEE functions are summarized at <http://www.globalplatform.org/> as: “TEE’s ability to offer isolated safe execution of authorized security software, known as trusted applications, enables it to provide end-to-end security by enforcing protected execution of authenticated code, confidentiality, authenticity, privacy, system integrity and data access rights.” Trusted applications, running in TEE, their assets, and data are isolated from the Rich Execution Environment where standard operating systems such as Linux run. TEE consists of several components:

1. A common abstraction layer, the Trusted Core Framework, providing OS functions, such as memory management, entry points for trusted applications, panic and cancelation handling, and trusted application properties access.
2. Inter-process communication used by rich execution environment applications to request services from TEE.
3. API for accessing services such as Trusted Storage for Data and Keys, TEE Cryptographic Operations, Time, and TEE Arithmetica.

AArch64, the 64-bit ARM architecture is compatible with AArch32. The members of the AArch64 family including ARMv8 Cortex-Axx ($xx = \{35, 53, 57, 72, 73\}$) processors share a number of features:

- Support a new instruction set, A64, with the same instruction semantics as AArch32, but with fewer conditional instructions. A64 includes major functional enhancements:
 1. 32 128-bit wide registers.
 2. Advanced SIMD supporting double-precision floating-point execution.
 3. Advanced SIMD supporting full IEEE 754¹⁷ execution.
- Include instruction-level support for cryptography, two encode and two decode instructions for AES, SHA-1, and SHA-256 support.
- Have 31 general-purpose registers accessible at all times.
- Provide revised exception handling in the AArch64 state.
- Support virtualization.
- Support the Trust Zone and the Global Trust TEE.

The *ARM Trust Zone* (ATZ) splits an ARM-based system into the Secure World, a trusted subsystem, responsible for the boot and the configuration of the entire system and the Non Secure World (NSW) intended for hosting operating systems such as Linux and Android, as well as user applications. A CPU has banked registers for each World.

Security-specific configurations can only be performed in the Secure World mode, while access to AMBA peripherals such as fingerprint readers, cryptographic engines, and others can be restricted to the Secure World. A secure context switch procedure routes interrupts either to the Secure or to the Non Secure World, depending upon the configuration and enables the two Worlds to communicate with

¹⁷ IEEE Standard for Floating-Point Arithmetic (IEEE 754) defines arithmetic formats, interchange formats, rounding rules, operations, and exception handling for floating-point numbers, see “IEEE Standard for Floating-Point Arithmetic” IEEE Computer Society (August 29, 2008). doi:[10.1109/IEEESTD.2008.4610935](https://doi.org/10.1109/IEEESTD.2008.4610935).

one another. ATZ is enabled by a set of hardware security extensions including: (i) a CPU with ARM Security Extensions (SE); (ii) a compliant Memory Management Unit (MMU); (iii) an AMBA system bus¹⁸, and (iv) interrupt and cache controllers.

T-KVM, a KVM-based trusted hypervisor for ARMv8, combines Trust Zone with GlobalPlatform TEE and SELinux [388]. T-KVM implements: (a) a trusted boot; (b) support for Trusted Computing inside a Virtual Machine; (c) a zero copy shared memory mechanism for data sharing between the two Trust Zone Worlds and between the VM and the host; and (d) a secure, ideally real-time, reliable, and error-free OS running in the Secure World.

The challenge of a secure boot is to eliminate vulnerabilities when the security mechanisms are not yet in place. The solution implemented in T-KVM is a four-stage boot process. A small program stored in the on-chip ROM, along with the public key needed for the attestation of the second stage loader, is activated in the first stage.

The second stage loads the microkernel in the Secure World zone and activates it. The third stage checks the integrity of the Linux kernel, a Non Secure World binary, and of its loader and, finally, the fourth stage runs it. The failure of any check in this chain of events brings the system to a secure state stop. T-KVM boot sequence is shown in Fig. 5.12(a).

The main challenge for supporting Trusted Computing inside a VM is the virtualization of the TEE APIs. To allow the TEE Client API to execute directly in the Guest OS, a specific QEMU device implements the TEE control plane and sets up its data plane, see Fig. 5.12(b). Requests for service, such as initialization/close session, command invocation, and notification of response, are sent to the TEE Device, which delivers them either to the trusted applications or to the client applications running on the guest OS. The data plane uses the shared memory. The TEE device notifies its driver upon receiving a response notification from the Trust Zone Secure World (TZSW), and the driver forwards the information to the Guest Client application.

Zero-copy shared memory is based on the fact that TAs can read/write VMs shared memory because TZSW can access the entire NSW workspace. The TEE Device control plane extends T-KVM shared memory mechanism, enabling it to send the shared memory address to the Secure World applications.

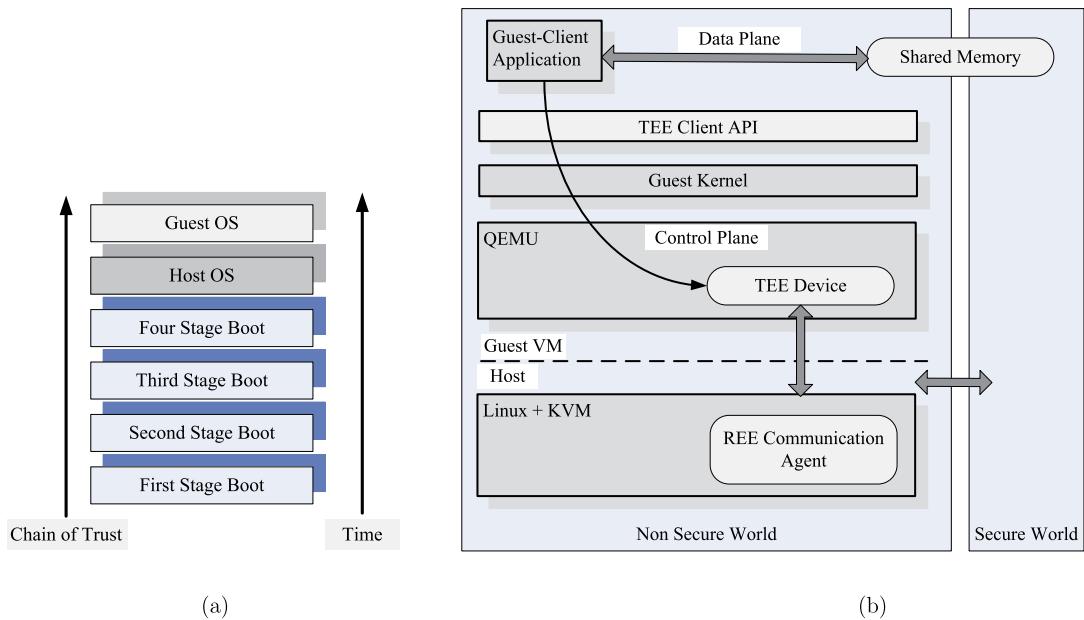
5.12 Paravirtualization of Itanium architecture

We now analyze some of the findings of a Xen project at HP-Laboratories [322]. This analysis will help us better understand the impact of the computer architecture on the ability to virtualize efficiently a given computer architecture. The goal of the project was to create a hypervisor for the Itanium family of IA64 Intel processors.

Itanium¹⁹ is a processor developed jointly by HP and Intel based on a new architecture, the Explicitly Parallel Instruction Computing. This architecture enables the processor to execute multiple instructions in each clock cycle and implements a form of Very Long Instruction Word (VLIW) archi-

¹⁸ The Advanced Microcontroller Bus Architecture (AMBA) is an open-standard, on-chip interconnect specification for connection and management of a large numbers of controllers and peripherals.

¹⁹ In the late 2000s, Itanium was the fourth-most deployed microprocessor architecture for enterprise-class systems. The first three were Intel's x86-64, IBM's Power Architecture, and Sun's SPARC.

**FIGURE 5.12**

T-KVM. (a) The boot sequence; a trusted application runs the Non Secure loader in stage three and in stage four a Non Secure OS is booted. TEE attestation and SELinux permissions are enforced while the host OS is running. Guest Client applications (CSs) use the Secure TEE Services while the Guest OS is running. (b) Communication between the Non Secure and the Secure Worlds, as described in [388].

tecture. In VLIW, a single instruction word contains multiple instructions, see http://www.dig64.org/about/Itanium2_white_paper_public.pdf.

The design mandated that the hypervisor should be capable of supporting execution of multiple operating systems in isolated protection domains with security and privacy enforced by the hardware. A hypervisor was also expected to support optimal server utilization and allow comprehensive measurement and monitoring for detailed performance analysis.

Virtualization of the IA64 architecture. The discussion in Section 5.3 shows that to be fully virtualizable, the ISA of a processor must conform to a set of requirements. Unfortunately, the IA64 architecture does not meet these requirements and that made the Xen project more challenging.

We first review the features of the Itanium processor important for virtualization and start with the observation that the hardware supports four *privilege rings*, PL0, PL1, PL2, and PL3. Privileged instructions can only be executed by the kernel running at level PL0, while applications run at level PL3 and can only execute nonprivileged instructions; PL1 and PL2 rings are generally not used.

The hypervisor uses *ring compression* and runs itself at PL0 and PL1, while forcing a guest OS to run at PL2. A first problem called *privilege leaking* is that several nonprivileged instructions allow an

application to determine the Current Privilege Level (CPL). As a result, a guest OS may not accept to boot or run, or may itself attempt to make use of all four privilege rings.

Itanium was selected because of its multiple functional units and multithreading support. The Itanium processor has 30 functional units: six general-purpose ALUs, two integer units, one shift unit, four data cache units, six multimedia units, two parallel shift units, one parallel multiply, one population count, three branch units, two 82-bit floating-point multiply–accumulate units, and two SIMD floating-point multiply–accumulate units. A 128-bit instruction word contains three instructions; the fetch mechanism can read up to two instruction words per clock from the L1 cache into the pipeline. Each unit can execute a particular subset of the instruction set.

The hardware supports 64-bit addressing. The processor has 32 64-bit general-purpose registers numbered from R0 to R31 and 96 automatically renumbered registers, R32 through R127, used by procedure calls. When a procedure is entered, the *alloc* instruction specifies the registers the procedure could access by setting the bits of a seven-bit field that controls the register usage; an illegal *read* operation from such a register out of range returns a zero value, while an illegal *write* operation to it is trapped as an illegal instruction.

The Itanium processor supports isolation of the address spaces of multiple processes with eight privileged *region* registers; the *processor abstraction layer* firmware allows the caller to set the values in the region register. The hypervisor intercepts the privileged instruction issued by the guest OS to its processor abstraction layer and partitions the set of address spaces among the guests OS to ensure isolation. Each guest is limited to 2^{18} address spaces.

The hardware has an *IVA register* to maintain the address of the *interruption vector table*; the entries in this table control both the interrupt delivery and the interrupt state collection. Different types of interrupts activate the interrupt handlers pointed at from this table, provided that the particular interrupt is not disabled. Each guest OS maintains its own version of this vector table and has its own IVA register; the hypervisor uses the guest OS IVA register to give control to the guest interrupt handler when an interrupt occurs.

CPU virtualization. When a guest OS attempts to execute a privileged instruction, the hypervisor traps and emulates the instruction. For example, when the guest OS uses the *rsm psr.i* instruction to turn off delivery of a certain type of interrupts, the hypervisor does not disable the interrupt but records the fact that interrupts of that type should not be delivered to the guest OS and, in this case, the interrupt should be masked.

There is a slight complication because the Itanium does not have an Instruction Register (IR) and the hypervisor has to use state information to determine if an instruction is privileged. Another complication is caused by the *register stack engine*, which operates concurrently with the processor and may attempt to access memory (load or store) and generate a page fault. Normally, the problem is solved by setting up a bit indicating that the fault is due to the register stack engine and, at the same time, the engine operations are disabled. The handling of this problem by the hypervisor is more intricate.

A number of *privileged-sensitive* instructions behave differently at different privilege levels. The hypervisor replaces each one of them with a privileged instruction during the dynamic transformation of the instruction stream. Among the instructions in this category are: (i) *cover*, saves stack information into a privileged register; the hypervisor replaces it with a *break.b* instruction; (ii) *tash* and *ttag*, access data from privileged virtual-memory control structures and have two registers as arguments. The hypervisor takes advantage of the fact that an illegal read returns a zero and an illegal write to a register in the range 32 to 127 is trapped and translates these instructions as *tash Rx = Ry* → *tpa Rx =*

$R(y + 64)$ and $ttag Rx = Ry \rightarrow tak Rx = R(y + 64)$, $0 \leq y \leq 64$; and (iii) access to performance data from performance data registers is controlled by a bit in the *Processor Status Register* with the *PSR.sp* instruction.

Memory virtualization. The virtualization is guided by the realization that a hypervisor should not be involved in most of memory read and write operations to prevent a significant degradation of the performance, but, at the same time, the hypervisor should exercise tight control and prevent a guest OS from acting maliciously. The Xen hypervisor does not allow a guest OS to access the memory directly, but instead, it inserts an additional layer of indirection called *metaphysical addressing* between virtual and real addressing.

A guest OS is placed in the metaphysical addressing mode. If the address is virtual, then the hypervisor first checks if the guest OS is allowed to access that address, and, if so, the hypervisor provides the regular address translation. The hypervisor is not involved when the address is physical. The hardware distinguishes between virtual and real addresses using bits in the Processor Status Register.

5.13 A performance comparison of virtual machines

There is well-documented evidence that hypervisors negatively affect the performance of applications [48,345,346]. The topic of this section is a quantitative analysis of the performance of VMs. The performance of two virtualization techniques is compared with the performance of a plain-vanilla Linux. The two VM systems are Xen and OpenVZ based [385] on paravirtualization and full virtualization, respectively.

OpenVZ, a system based on OS-level virtualization, uses a single-patched Linux kernel. The guest operating systems in different containers may be different software distributions, but must use the same Linux kernel version that the host uses. An OpenVZ container emulates a separate physical server; it has its own files, users, process tree, IP address, shared memory, semaphores, and messages. Each container can have its own disk quotas.

OpenVZ's lack of virtualization flexibility is compensated by a lower overhead. OpenVZ memory allocation is more flexible than in hypervisors based on paravirtualization. The memory not used in one virtual environment can be used by other virtual environments. The system uses a common file system; each virtual environment is a directory of files isolated using *chroot*. To start a new VM, one needs to copy the files from one directory to another, create a *config* file for the VM, and launch the VM.

OpenVZ has a two-level scheduler: at the first level, the fair-share scheduler allocates CPU time slices to containers based on *cpuunits* values. The second-level scheduler is a standard Linux scheduler deciding what process to run in that container. The I/O scheduler is also two-level; each container has an I/O priority, and the scheduler distributes the available I/O bandwidth according to priorities.

The discussion in [385] is focused on user's perspective, thus, the performance measures analyzed are the throughput and the response time. The general question is whether consolidation of the applications and the servers is a good strategy for cloud computing. The specific questions examined are: How does performance scale up with the load? What is the impact on performance of a mixture of applications? What are the implications of the load assignment on individual servers?

There is substantial experimental evidence that the load placed on system resources by a single application varies significantly in time. A time series displaying CPU consumption of a single application clearly illustrates this fact and justifies CPU multiplexing among threads and/or processes. The

concept of *application and server consolidation* is an extension of the idea of creating an aggregate load consisting of several applications and aggregating a set of servers to accommodate this load. The peak resource requirements of individual applications are unlikely to be synchronized, therefore, the aggregate average resource utilization is expected to increase.

The application used for comparison in [385] is a two-tier system consisting of an Apache web server and a MySQL database server. A client of the web application starts a session as the user browses through various items in the database, requests information about individual items, and buys or sells items. Each session requires the creation of a new thread; thus, an increased load means an increased number of threads. To understand the potential discrepancies in performance among the three systems, a performance-monitoring tool reports the counters that enable the estimation of: (i) the CPU time used by a binary; (ii) the number of L2-cache misses; and (iii) the number of instructions executed by a binary.

The experimental setup for three different experiments are shown in Fig. 5.13. The two tiers of the application, the web and the database, each run on a single server for the Linux, the OpenVZ, and the Xen systems in the first group of experiments. When the workload increases from 500 to 800 threads, the throughput increases linearly with the workload.

Response time increases only slightly for the base system and for the OpenVZ system, while it increases 600% for the Xen system. For 800 threads, the response time of the Xen system is four times larger than for OpenVZ. CPU consumption grows linearly with the load in all three systems. DB consumption represents only 1–4% of CPU consumption.

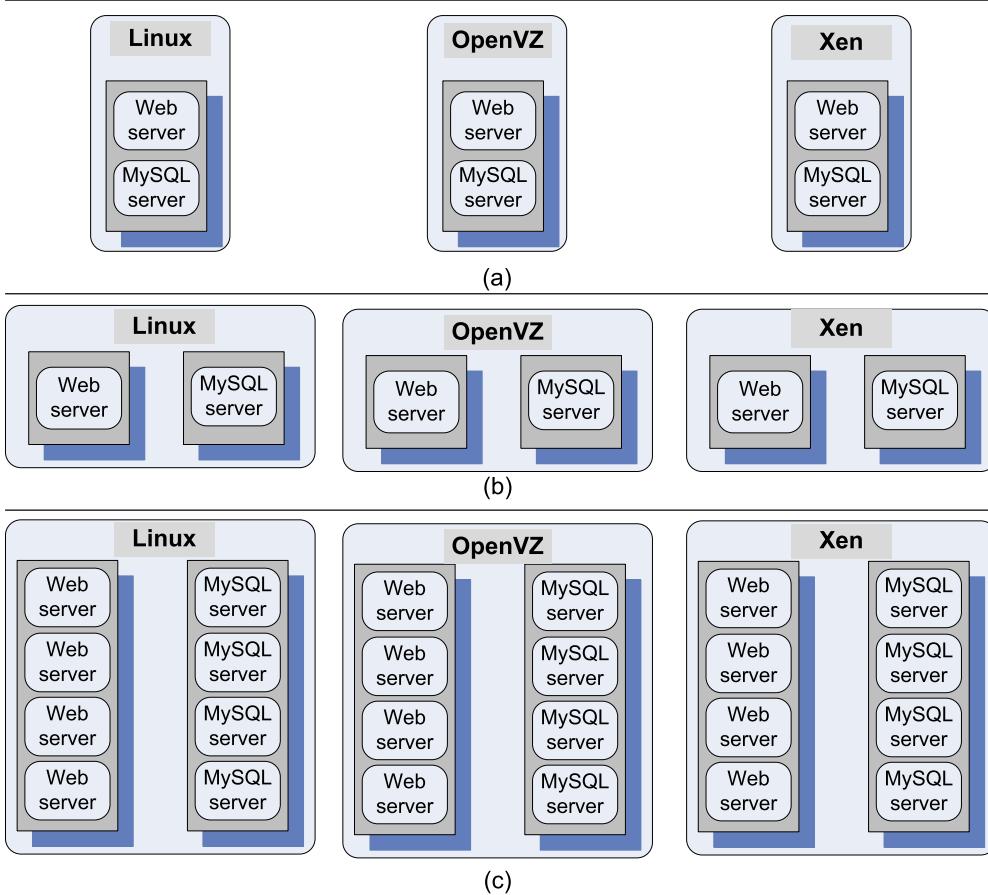
For a given workload the web-tier CPU consumption for the OpenVZ system is close to the base system, Linux, and it is about half of that for the Xen system. The performance analysis tool shows that the OpenVZ execution has two times more L2-cache misses than the base system, while the Xen *Dom0* has 2.5 times more, and the Xen application domain has nine times more cache misses.

The base system and the OpenVZ run a Linux OS, and the sources of cache misses can be compared directly, while Xen runs a modified Linux kernel. For the Xen-based system, the procedure *hypervisor_callback*, invoked when an event occurs, and the procedure *evtchn_do_upcall*, invoked to process an event, are responsible for 32% and 44%, respectively, of the L2-cache misses. The percentage of the instructions invoked by these two procedures are 40% and 8%, respectively.

Most L2-cache misses in OpenVZ and the base system occur in: (i) *do_anonymous_pages*, a procedure used to allocate pages for a particular application with the percentage of cache misses 32% and 25%, respectively; (ii) procedures *_copy_to_user_ll* and *_copy_from_user_ll* used to copy data from user to system buffers and back with the percentage of cache misses (12 + 7)% and (10 + 1)%, respectively. The first figure refers to copying from user to system buffers and the second to copying from system buffers to the user space.

The second group of experiments use two servers, one for the web and the other for the DB application, for each one of the three systems. When the load increases from 500 to 800 threads, the throughput increases linearly with the workload. The response time of the Xen system increases only 114%, compared with 600% reported for the first experiments.

The CPU time of the base system, the OpenVZ system, the Xen *Dom0*, and the *User Domain* are similar for the web application; for the DB application, the CPU time of the OpenVZ system is twice as large as that of the base system, while *Dom0* and the *User Domain* require CPU times of 1.1 and 2.5 times larger than the base system.

**FIGURE 5.13**

The setup for the performance comparison of a native Linux system with OpenVZ and the Xen systems. The applications are a web server and a MySQL database server. (a) In the first experiment, the web and the DB share a single system; (b) In the second experiment, the web and the DB run on two different systems; (c) In the third experiment, the web and the DB run on two different systems, and each has four instances.

The L2-cache misses for the web application relative to the base system are: the same for OpenVZ and 1.5 larger for *Dom0* of Xen and 3.5 times larger for the *User Domain*. The L2-cache misses for the DB application relative to the base system are: *two* times larger for the OpenVZ, 3.5 larger for *Dom0* of Xen, and seven times larger for the *User Domain*.

The third group of experiments uses two servers, one for the web and the other for the DB application, for each one of the three systems, but runs four instances of the web and the DB application on the two servers. The throughput increases linearly with the workload for the range used in the previous

two experiments, from 500 to 800 threads. The response time remains relatively constant for OpenVZ and increases five times for Xen.

The main conclusion drawn from these experiments is that the virtualization overhead of Xen is considerably higher than that of OpenVZ and that this is due primarily to L2-cache misses. Xen performance degradation is noticeable when the workload increases. Another important conclusion is that hosting multiple tiers of the same application on the same server is not optimal.

5.14 Open-source software platforms for private clouds

Private clouds provide a cost effective alternative for very large organizations. A private cloud has essentially the same structural components as a commercial one: the servers, the network, hypervisors running on individual systems, an archive containing disk images of VMs, a front-end for communication with the user, and a cloud control infrastructure. Open-source cloud computing platforms such as Eucalyptus [368], OpenNebula, and Nimbus can be used as a control infrastructure for a private cloud.

Schematically, a cloud infrastructure carries out the following steps to run an application: (i) retrieves the user input from the front-end; (ii) retrieves the disk image of a VM (Virtual Machine) from a repository; (iii) locates a system and requests the hypervisor running on that system to set up a VM; and (iv) invokes the DHCP (see Section 6.1) and the IP bridging software to set up a MAC and IP address for the VM.

Eucalyptus (<http://www.eucalyptus.com/>) can be regarded as an open-source counterpart of Amazon's EC2. The system supports several operating systems, including: CentOS 5 and 6, RHEL 5 and 6, Ubuntu 10.04 LTS, and 12.04 LTS. The components of the system are:

1. Virtual Machine. Runs under several hypervisors including Xen, KVM, and VMware.
2. Node Controller. Runs on every server/node designated to host a VM and controls the activities of the node. Reports to a cluster controller.
3. Cluster Controller. Controls a number of servers. Interacts with the node controller on each server to schedule requests on that node. Cluster controllers are managed by cloud controller.
4. Cloud Controller. Provides the cloud access to end-users, developers, and administrators. It is accessible through command line tools compatible with EC2 and through a web-based dashboard. Manages cloud resources, makes high-level scheduling decisions, and interacts with cluster controllers.
5. Storage Controller. Provides persistent virtual hard drives to applications. It is the correspondent of EBS. Users can create snapshots from EBS volumes. Snapshots are stored in Walrus and shared across availability zones.
6. Storage Service (Walrus). Provides persistent storage and, similarly to S3, allows users to store objects in buckets.

The system supports a strong separation between the user space and administrator space; users access the system via a web interface, while administrators need root access. The system supports a decentralized resource management of multiple clusters with multiple cluster controllers, but a single head node for handling user interfaces. It implements a distributed storage system called Walrus, the analog of Amazon's S3 system. The procedure to construct a VM is based on the generic one described in [445]:

Table 5.4 A side-by-side comparison of Eucalyptus and OpenNebula.

Eucalyptus		OpenNebula
Design	Emulate EC2	Customizable
Cloud type	Private	Private
User population	Large	Small
Applications	All	All
Customizability	Administrators limited users	Administrators and users
Internal security	Strict	Loose
User access	User credentials	User credentials
Network access	To cluster controller	—

- (i) The *euca2ools* front-end is used to request a VM.
- (ii) The VM disk image is transferred to a compute node.
- (iii) The disk image is modified for use by the hypervisor on the compute node.
- (iv) The compute node sets up network bridging to provide a virtual NIC with a virtual MAC address.
- (v) The head node the DHCP is set up with the MAC/IP pair.
- (vi) The hypervisor activates the VM.
- (vii) The user can now *ssh* directly into the VM.

The system can support a large number of users in a corporate enterprise environment. Users are shielded from the complexity of disk configurations and can choose for their VM from a set of five configurations of available processors, memory, and hard-drive space setup by the system administrators.

Open-Nebula (<http://www.opennebula.org/>) is a private cloud with users actually logging into the head node to access cloud functions. The system is centralized, and its default configuration uses the NFS filesystem. The procedure to construct a VM consists of several steps: (i) a user signs in to the head node using *ssh*; (ii) next, it uses the *onevm* command to request a VM; (iii) the VM template disk image is transformed to fit the correct size and configuration within the NFS directory on the head node; (iv) the *oned* daemon on the head node uses *ssh* to log into a compute node; (v) the compute node sets up network bridging to provide a virtual NIC with a virtual MAC; (vi) the files needed by the hypervisor are transferred to the compute node via the NFS; (vii) the hypervisor on the compute node starts the VM; and (viii) the user is able to *ssh* directly to the VM on the compute node.

According to the analysis in [445], the system is best suited for an operation involving a small-to-medium sized group of trusted and knowledgeable users who are able to configure this versatile system based on their needs.

Table 5.4 summarizes the features of the two systems [445]. Eucalyptus is best suited for a large corporation with its own private cloud because it ensures a degree of protection from user malice and mistakes; OpenNebula is best suited for a testing environment with a few servers.

OpenStack is an open source project started in 2009 at NASA in collaboration with Rackspace (<http://www.rackspace.com>) to develop a scalable cloud OS for farms of servers using standard hardware. Though recently NASA has moved its cloud infrastructure to AWS, in addition to Rackspace, several other companies including HP, Cisco, IBM, and Red Hat have an interest in *OpenStack*. The current

version of the system supports a wide range of features such as: APIs with rate limiting and authentication, live VM management to run, reboot, suspend, and terminate instances, role-based access control, and the ability to allocate, track, and limit resource utilization. The administrators and the users control their resources using an extensible web application called the *Dashboard*.

5.15 The darker side of virtualization

Can virtualization empower the creators of malware²⁰ to carry out their mischievous activities with impunity and minimal danger of being detected? How difficult is it to implement such a system? What are the means to prevent this type of malware to be put in place? The answers to these questions are discussed in this section.

It is well understood that, in a layered structure, a defense mechanism at some layer can be disabled by malware running at a layer below it. Thus, the winner in the continuous struggle between the attackers and the defenders of a computing system is the one in control of the lowest layer of the software stack, the one which controls the hardware.

A hypervisor allows a guest OS to run on virtual hardware; the hypervisor offers to the guest operating systems a hardware abstraction and mediates its access to the physical hardware. We argue that a hypervisor is simpler and more compact than a traditional OS, thus, it is more secure; but what if the hypervisor itself is forced to run above another software layer, and, thus, it is prevented from exercising direct control of the physical hardware?

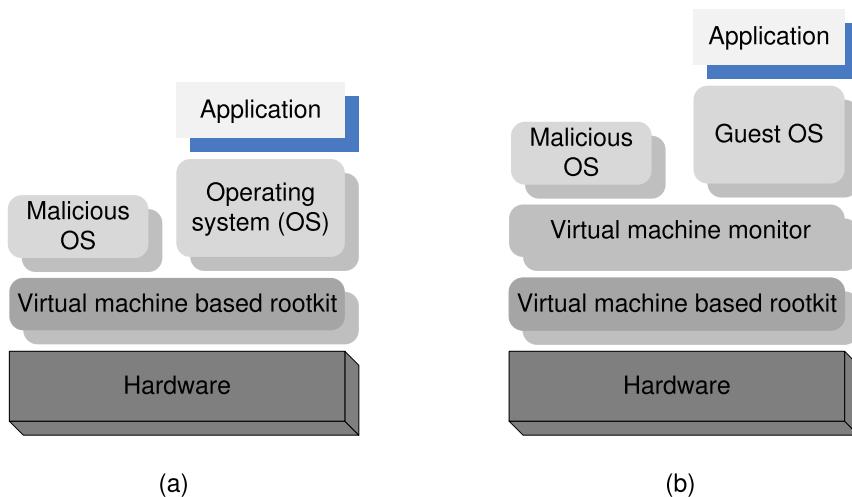
A 2006 paper [274] argues that it is feasible to insert a “rogue hypervisor” between the physical hardware and an OS, as shown in Fig. 5.14(a). Such a rogue hypervisor is called a *Virtual-Machine Based Rootkit* (VMBR). The term *rootkit* refers to malware with a privileged access to a system. The name comes from *root*, the most privileged account on a Unix system, and *kit*, a set of software components.

It is also feasible to insert the VMBR between the physical hardware and a legitimate hypervisor, as in Fig. 5.14(b). Since a VM running under a legitimate hypervisor sees a virtual hardware, the guest OS will not notice any change in the environment. The only trick is to present the legitimate hypervisor with a hardware abstraction, rather than allow it to run on the physical hardware.

Before we address the question how such an insertion is possible, we should point out that in this approach, the malware runs either inside a hypervisor or with the support of a hypervisor. A hypervisor is a very potent engine for the malware: It prevents the software of the guest OS or the application from detecting malicious activities. A VMBR can record key strokes, system states, data buffers sent to, or received from the network, and data to be written to, or read from the disk with impunity; moreover, it can change any data at will.

The only way for a VMBR to take control of a system is to modify the boot sequence and to first load the malware and only then load the legitimate hypervisor, or the OS; this is only possible if the attacker has root privileges. Once the VMBR is loaded, it must also store its image on the persistent storage.

²⁰ Malware, an abbreviation for *malicious software*, is software designed specifically to circumvent the authorization mechanisms and gain access to a computer system, gather private information, block access to a system, or disrupt the normal operation of a system; computer viruses, worms, spyware, and Trojan horses are examples of malware.

**FIGURE 5.14**

The insertion of a *Virtual-Machine Based Rootkit* (VMBR) as the lowest layer of the software stack running on the physical hardware; (a) below an OS; (b) below a legitimate hypervisor. The VMBR enables a malicious OS to run surreptitiously and makes it invisible to the genuine or the guest OS and to the application.

The VMBR can enable a separate malicious OS to run surreptitiously and make this malicious OS invisible to the guest OS and to the application running under it. Under the protection of the VMBR, the malicious OS could: (i) observe the data, the events, or the state of the target system; (ii) run services such as spam relays or distributed denial-of-service attacks; or (iii) interfere with the application.

A proof-of-concept VMBRs to subvert Windows XP and Linux and several services based on these two platforms are described in [274]. We should stress that modifying the boot sequence is by no means an easy task, and once an attacker has root privileges she is in total control of a system.

5.16 Virtualization software

Several virtualization software packages, including hypervisors, OS-level virtualization software, and desktop virtualization software, are available. There are two types of hypervisors, native and hosted. The set of *native hypervisors* includes:

Red Hat Virtualization (RHV)—enterprise virtualization based on KVM hypervisor.

Hyper—creates VMs on x86-64 systems running Windows.

z/VM—current version of IBM's VM operating systems.

VMware ESXi—enterprise-class, type-1 hypervisor from VMware.

Oracle VM Server for x86—server virtualization from Oracle Corporation. Incorporates the free, open-source Xen. Supports Windows, Linux, and Solaris guests.

Adeos—Adaptive Domain Environment for Operating Systems is a nanokernel hardware abstraction layer.

XtratuM—bare-metal hypervisor for embedded real-time systems. Available for the instruction sets x86, ARM Cortex-R4F processors, and others.

There are several *hosted independent hypervisors* including: (i) VMware Fusion—software hypervisor developed for Intel-based Macs to run Microsoft Windows, Linux, NetWare, or Solaris on VMs, along with the OS X OS, based on paravirtualization, hardware virtualization, and dynamic recompilation; (ii) PearPC—architecture-independent PowerPC platform emulator for PowerPC operating systems, including pre-Intel versions of OS X, Darwin, and Linux; (iii) Oracle VM VirtualBox—free and open-source hypervisor for x86 computers; and (iv) QEMU (Quick Emulator)—free and open-source hosted hypervisor. There are also *hosted specialized hypervisors* including: (i) coLinux—Cooperative Linux allows Microsoft Windows and the Linux kernel to run simultaneously; (ii) MoM—Mac-on-Mac is a port of Mac-on-Linux for Mac OS X; (iii) Mac-on-Linux—open-source VM for running the classic Mac OS or OS X on PowerPC computers running Linux; (iv) bhyve—a type-1 hypervisor included in FreeBSD running FreeBSD 9+, OpenBSD, NetBSD, Linux and Windows desktop, and Windows Server; and (v) L4Linux—a variant of Linux kernel running virtualized on L4 microkernel; L4Linux kernel runs a service on L4.

5.17 History notes and further readings

Virtual memory was the first application of virtualization concepts to commercial computers; it allowed multiprogramming and eliminated the need to tailor applications to the physical memory available on individual systems. Paging and segmentation are the two mechanisms supporting virtual memory. Paging was developed for the Atlas Computer built in 1959 at the University of Manchester. Independently, in 1961 Burroughs Corporation developed B5000, the first commercial computer with virtual memory; B5000 virtual memory used segmentation rather than paging.

In 1967, IBM introduced 360/67, the first IBM system with virtual memory expected to run on a new OS, called TSS. Before TSS was released, an operating system called CP-67 was created; CP-67 gave the illusion of several standard IBM-360 systems without virtual memory. The first hypervisor supporting full virtualization was the CP-40 system which ran on a S/360-40 that was modified at the IBM Cambridge Scientific Center to support Dynamic Address Translation, a key feature that allowed virtualization. In CP-40, the hardware's supervisor state was virtualized as well, allowing multiple operating systems to run concurrently in separate VM contexts.

Virtualization was driven by the need to share very expensive hardware among a large population of users and applications in the early age of computing. VM/370 system, released in early 1970s for large IBM mainframes, was very successful; it was based on a reimplementations of CP/CMS. In VM/370, a new VM was created for every user, and this VM interacted with the applications. The hypervisor managed hardware resources and enforced the multiplexing of resources. Modern-day IBM mainframes, such as the zSeries line, retain backwards-compatibility with the 1960s-era IBM S/360 line.

Microprocessor development coupled with advances in storage technology contributed to the rapid decrease of hardware costs and led to introduction of personal computers at one end of the spectrum and of large mainframes and massively parallel systems at the other end. The hardware and the operating

systems of 1980s and 1990s gradually limited virtualization and focused instead on efficient multitasking, user interfaces, and support for networking and security problems brought in by interconnectivity.

Advancements in computer and communication hardware and the explosion of the Internet partially due to the success of the World Wide Web in late 1990s renewed the interest in virtualization to support server security and isolation of services. In their review paper, Rosenbloom and Garfinkel write [423]: “hypervisors give OS developers another opportunity to develop functionality no longer practical in today’s complex and ossified operating systems, where innovation moves at a geologic pace.” Nested virtualization was first discussed in early 1970s by Popek and Goldberg [201,402].

Further readings. The text of Saltzer and Kaashoek [430] is a very good introduction to virtualization principles. Virtual machines are dissected in a paper by Smith and Nair [450], and architectural principles for virtual computer systems are analyzed in [200,201].

An insightful discussion of hypervisors is provided by the paper of Rosenblum and Garfinkel [423]. Several papers [48,345,346] discuss in depth the Xen hypervisor and analyze its performance, while [519] is a code repository for Xen. The Denali system is presented in [514].

Modern systems such as Linux Vserver (<http://linux-vserver.org/>), OpenVZ (Open VirtuAlization) [376], FreeBSD Jails [415], and Solaris Zones [403] implement *OS-level virtualization technologies*. Reference [385] compares the performance of two virtualization techniques with a standard OS.

A 2001 paper [100] argues that virtualization allows new services to be added without modifying the OS. Such services are added below the OS level, but this process creates a semantic gap between the VMs and these services. Reflections on the design of hypervisors are the subject of [101] and a discussion of Xen is reported in [108]. The state of the art and the future of nested virtualization are the subject of [127]. An implementation of nested virtualization for KVM is discussed in [57]. [541] surveys security issues in virtual systems and [283] covers reliability in virtual infrastructures. Virtualization technologies in HPC are analyzed in [410], and [452] provides a critical view on virtualization. [508] reports on IBM virtualization strategies.

5.18 Exercises and problems

- Problem 1.** Identify the milestones in the evolution of operating systems during the half century from 1960 to 2010 and comment on the statement from [423] “Hypervisors give OS developers another opportunity to develop functionality no longer practical in today’s complex and ossified operating systems, where innovation moves at a geologic pace.”
- Problem 2.** Virtualization simplifies the use of resources, isolates users from one another, and supports replication and mobility but exacts a price in terms of performance and cost. Analyze each one of these aspects for: (i) memory virtualization, (ii) processor virtualization, and (iii) virtualization of a communication channel.
- Problem 3.** Virtualization of the processor combined with virtual memory management pose multiple challenges; analyze the interaction of interrupt handling and paging.
- Problem 4.** In Section 5.2, we stated that a hypervisor is a much simpler and better specified system than a traditional OS. Hypervisor vulnerability is reduced because the systems expose a much smaller number of privileged functions. Compare the number of lines of code and of system calls for several operating systems including Linux, Solaris, FreeBSD, Unbuntu, AIX, and Windows with the corresponding figures for several system VMs.

- Problem 5.** In Section 5.4 we state that a hypervisor for a processor with a given ISA can be constructed if the set of *sensitive instructions* is a subset of the privileged instructions of that processor. Identify the set of sensitive instructions for the x86 architecture and discuss the problem each one of these instructions poses.
- Problem 6.** Table 5.3 summarizes the effects of Xen network performance optimization reported in [346]. The send-data rate of a guest domain is improved by a factor of more than four, while the improvement of the receive data rate is very modest. Identify several possible reasons for this discrepancy.
- Problem 7.** VMware EX Server supports full virtualization of x86 architecture. Analyze how VMware provides the functions discussed in Table 5.2 for Xen.
- Problem 8.** In 2012, Intel and HP announced that *Itanium* architecture will be discontinued. Review the architecture discussed in Section 5.12, and identify several possible reasons for this decision.
- Problem 9.** Read [385] and analyze the results of performance comparison discussed in Section 5.13.

Cloud access and cloud interconnection networks

6

Unquestionably, communication is at the heart of cloud computing. The decades-long evolution of microprocessor and storage technologies, computer architecture and software systems, parallel algorithms and distributed control strategies enabled cloud computing yet, only *interconnectivity supported by a continually evolving Internet made cloud computing feasible*.

This chapter presents basic concepts necessary for understanding the intricacies of the communication infrastructure used for cloud computing. Communication and computing are inseparable concepts as discussed in Chapter 3, where we have seen that, even at the scale of a single core, efficient communication between the CPU, the memory, and the I/O sub-system are critical for optimal performance.

Communication, in the case of a single-multicore processor or on a system on a chip (SoC), is very efficient because it involves only hardware busses or networks on a chip and router-based packet switching networks between SoC modules. Large-scale systems such as supercomputers are built around complex and expensive interconnection networks controlled mostly by hardware thus support lower-latency and higher-bandwidth communication than cloud interconnects.

Distributed systems and computer clouds communicate over communication networks or the Internet. The Internet is a network of networks interconnected by switches operating under the control of routing algorithms and it is based on intricate software-based communication protocols that increase the communication latency and limit the bandwidth.

The designers of a cloud computing infrastructure are acutely aware that the farther data travels from the CPU the lower the bandwidth and the larger the communication latency. The limits of cloud interconnection networks bandwidth and latency are tested by the demands of systems with millions of servers; some cloud workloads are more affected by these limits than others.

Cloud workloads fall into four broad categories based on their dominant resource needs: CPU-intensive, memory-intensive, I/O-intensive, and storage-intensive. While the first two benefit from, but do not require, high-performing networking, the latter two do. Networking performance directly impacts the performance of I/O- and storage-intensive workloads.

The costs of the networking infrastructure continue to rise at a time when the cost of other components of cloud infrastructure continue to decrease. Moreover, applications in science and engineering are data- and network-intensive and require more expensive networks with higher bandwidth and lower latency than today's cloud interconnects.

This chapter starts with an overview of the Internet and the World Wide Web in Sections 6.1, 6.2, and 6.3, covering basic concepts on network architecture. The Internet is in fact *content-centric*, i.e., it is used primarily to access content, and Section 6.4 presents research ideas for a content-centric Internet. Software-defined networks, covered in Section 6.5, support an alternative to network management that enables dynamic, programmatically efficient network configuration to improve network performance and monitoring.

The focus then changes to the communication fabric used by the cloud infrastructure and to analysis of interconnection networks architecture and algorithms. After an overview of the interconnection networks in Section 6.6, the chapter covers multistage networks, Infiniband, and storage area networks in Sections 6.7, 6.8, and 6.9, respectively.

A scalable data-center architecture and network resource management are the topics of Sections 6.10 and 6.11. Then, the focus changes again, this time to content-delivery networks and vehicular networks in Sections 6.12 and 6.13. Further readings, historical notes, and a set of exercises and problems conclude the chapter.

6.1 Packet-switched networks and the Internet

Computer clouds are accessed through the Internet, a network of packet-switched networks. A packet-switched network transports data units called *packets* through a maze of *switches* where packets are queued and routed towards their destination. Packets are subject to random delays, loss, and may arrive at their final destination out of order.

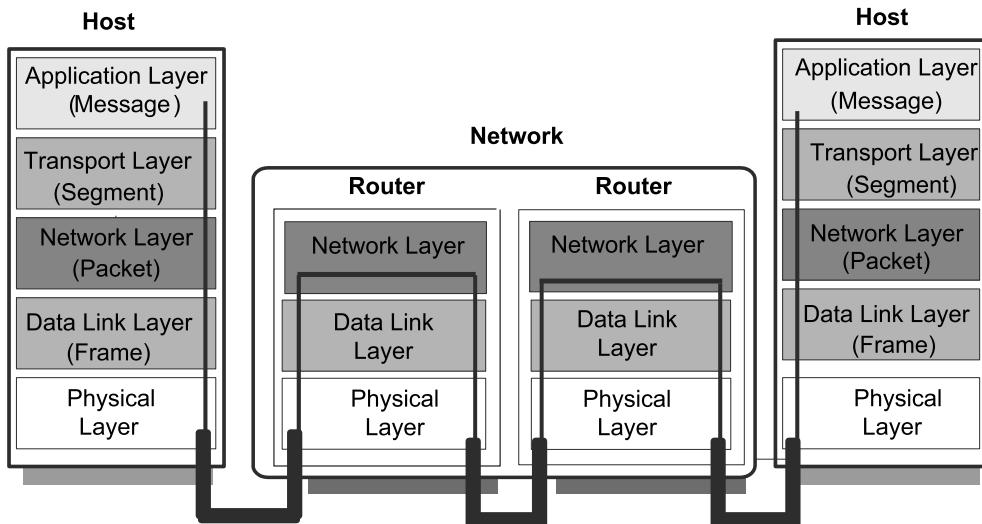
A few basic concepts are defined next. A *datagram* is the transfer unit in a packet-switched network. In addition to its *payload*, a datagram has a *header* containing control information necessary for its transport through the network. A *network architecture* describes the protocol stack used for communication. A *protocol* is a set of rules on how to communicate, specifying the actions taken by the sender and the receiver of a data unit. A *network host* identifies a system located at the network edge capable to initiate and to receive communication, a computer, a mobile device, e.g., a phone, or a sensor.

Network architecture and protocols. A packet-switched network has a *network core* consisting of routers and control systems interconnected by very-high-bandwidth communication channels and a *network edge* where the end-user systems reside.

A packet-switched network is a complex system consisting of a large number of autonomous components subject to complex and, sometimes contradictory requirements. Basic strategies for implementing a complex system are *layering* and *modularization*. Layering means decomposing a complex function into elements interacting through well-defined channels; a layer can only communicate with its adjacent layers. *Modularization* means dividing a system into interchangeable modules to create a flexible system while reducing the number of unique building blocks.

The Internet protocol stack, based on the TCP/IP network architecture, is shown in Fig. 6.1. Data flows down the protocol stack of the sending host from the application layer to the transport layer, then to the network layer, and to the data link layer. The physical layer pushes the streams of bits through a physical communication link encoded either as electrical, optical, or electromagnetic signals. At the receiving host, the packets flow up from the physical, to the data link, then to the network and transport layers and are finally delivered to the application layer. The corresponding data units for the five-layer architecture are: messages, segments, packets, frames, and encoded bits, respectively.

The *transport layer* is responsible for end-to-end communication, from an application running on the sending host to its peer running on the destination host using either TCP or UDP protocols. The network layer decides where the packet should be sent, either to another router, or to a destination host connected to a local area network connected to the router. IP, the *network layer* protocol, guides packets through the packet-switched network from the point of entry to the place where a packet exits

**FIGURE 6.1**

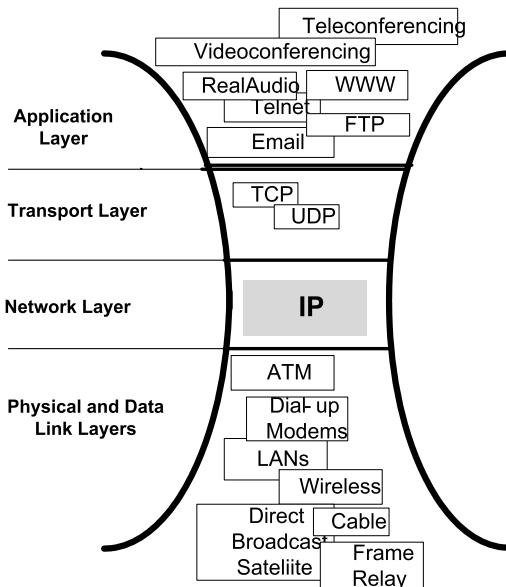
Internet protocol stack. Applications running on hosts at the edge of the network communicate using application-layer protocols. The transport layer deals with end-to-end delivery. The network layer is responsible for routing a packet through the network. The data link layer ensures reliable communication between adjacent nodes of the network, and the physical layer transports streams of bits encoded as electrical, optical, or electromagnetic signals (the thick lines represent such bit pipes).

the network. The *data link layer* encapsulates the packet for the communication link to the next hop. Once a packet reaches a router, the bits are passed to the data link and then to the network layer.

A protocol on one system communicates with its *peer* on another system. For example, the transport protocol on the sender, host A, communicates with the transport protocol on the receiver, host B. On the sending side, A, the transport protocol encapsulates the data from the application layer and adds control information as headers that can only be understood by its peer, the transport layer on host B. When the peer receives the data unit, it carries out a decapsulation, retrieves the control information, removes the headers, and then passes the payload to the next layer up, the application layer on host B.

The payload for the data link layer at the sending site includes the network header and the payload at the network layer. In turn, the network layer payload includes transport layer header and its payload consisting of the application layer header and application data.

The Internet. The Internet is a network of networks, a collection of separate, autonomous, and distinct networks. All networks adhere to a common framework and use: (i) globally unique IP addresses; (ii) the IP (Internet Protocol) routing protocol; and (iii) the Border Gateway Routing (BGP) protocol. BGP is a path vector reachability protocol that makes core routing decisions. BGP maintains a table of IP networks designating network reachability among autonomous systems. BGP makes routing decisions based on path, network policies, and/or rule sets.

**FIGURE 6.2**

The hourglass network architecture of the Internet. Regardless of the application, the transport protocol, and the physical network, all packets are routed from the source to the destination using the IP protocol and the IP address of the destination.

An *IP address* is a string of integers uniquely identifying every host connected to the Internet. An IP address allows the network to identify first the destination network and then the host in that network where a datagram should be delivered. A host may have multiple IP addresses, and it may be connected to more than one network. A host could be a supercomputer, a workstation, a laptop, a mobile phone, a network printer, or any other physical device with a network interface.

The Internet is based on a hourglass network architecture, as shown in Fig. 6.2. *The hourglass architecture is partially responsible for the explosive growth of the Internet*, it allowed the lower layers of the architecture to evolve independently from the upper layers. The communication technology drives dramatic changes of the lower layers of the Internet architecture, including the increase of the communication bandwidth and the widespread use of wireless networks and satellite communication. The software and the applications are the engines of progress for the upper layers of the architecture.

The hourglass model reflects the *end-to-end* architectural design principle. The model captures the fact that all packets transported through the Internet use IP to reach their destination. IP provides only best-effort delivery because any router along the path from the source to the destination may drop a packet when it is overloaded.

Another important architectural design principle of the Internet is the *separation between the routing and the forwarding planes*. This separation has allowed the forwarding plane to function while the routing evolved. The *forwarding plane* decides what to do with packets arriving on an inbound interface of a router. This plane uses a table to lookup the destination address of an incoming packet, then

retrieves the information to determine the path from the receiving interface to the proper outgoing interface(s) through the internal forwarding fabric of the router. The *routing plane* is responsible for building the routing table in each router.

In addition to the IP or logical address, each network interface, the hardware connecting a host with a network, has a unique *physical* or *MAC address*. While the MAC address is permanently assigned to a network interface of the device, the IP address may be dynamically assigned. The IP address of a mobile device changes depending on the device location and the network it is connected to.

The Dynamic Host Configuration Protocol (DHCP) is an automatic configuration protocol. DHCP assigns an IP address to a client system. A DHCP server has three methods of allocating IP addresses:

1. Dynamic allocation—a network administrator assigns a range of IP addresses to DHCP. During network initialization, each client computer on the LAN is configured to request an IP address from the DHCP server. The request-and-grant process uses a lease concept with a controllable time period, allowing the DHCP server to reclaim (and then reallocate) IP addresses that are not renewed.
2. Automatic allocation—the DHCP server permanently assigns a free IP address to a client from the range defined by the administrator.
3. Static allocation—the DHCP server allocates an IP address based on a manually filled in table with (MAC address–IP address) pairs. Only a client with a MAC address listed in this table is allocated an IP address.

Once a packet reaches the destination host, it is delivered to the proper transport protocol daemon which, in turn, delivers it to the application that listens to a *port*, an abstraction of the end-point of a logical communication channel, Fig. 6.3. The processes or threads running an application use an abstraction called *socket* to send and receive data through the network. A socket manages a queue of incoming messages and one for outgoing messages.

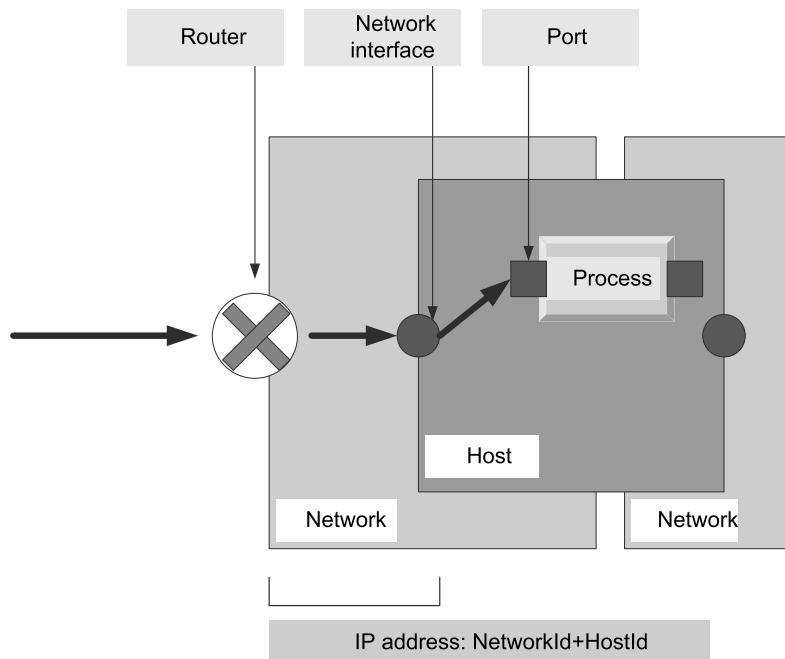
Internet transport protocols. Internet uses two transport protocols, a connectionless datagram protocol, UDP (User Datagram Protocol), and a connection-oriented protocol, TCP (Transport Control Protocol). The header of a datagram contains information sufficient for routing through the network from the source to the destination. The arrival time and the datagram delivery order are not guaranteed.

To ensure efficient communication, UDP assumes that error checking and error correction are either not necessary or performed by the application. Datagrams may arrive out of order, duplicated, or may not arrive at all. Applications using UDP include: DNS (Domain Name System), VoIP (Voice over IP), TFTP (Trivial File Transfer Protocol), streaming media applications such as IPTV, and online games.

TCP provides reliable, ordered delivery of a stream of bytes from an application on one system to its peer on the destination system; an application sends/receives data units called *segments* to/from a specific port, an abstraction of an endpoint of a logical communication link. TCP is the transport protocol used by the World Wide Web, email, file transfer, remote administration, and many other important applications.

TCP uses an end-to-end *flow-control mechanism* based on a sliding window, a range of packets the sender can send before receiving an acknowledgment from the receiver. These mechanisms enable the receiver to control the rate of segments sent and to process them reliably.

A network has a finite capacity to transport data, and when its load approaches capacity, we witness undesirable effects, such as routers start dropping packets and the delays and the jitter increase. An obvious analogy is a highway where the time to travel from point A to point B increases dramatically

**FIGURE 6.3**

Packet delivery to processes and threads; a packet is first routed by the IP protocol to the destination network and then to the host specified by the IP address. Applications listen to *ports*, abstractions of the endpoint of a communication channel.

in the case of congestion. A solution for traffic management is to introduce traffic lights limiting the rate at which new traffic is allowed to enter the highway, and this is precisely what the TCP emulates.

TCP uses several mechanisms for *congestion control* discussed in Section 6.3. These mechanisms control the rate of the data entering the network, keeping the data flow below a rate that would lead to a network collapse, and enforcing a fair allocation among flows. Acknowledgments coupled with timers are used to infer network conditions between the sender and receiver. TCP congestion control policies are based on four algorithms, *slow-start*, *congestion avoidance*, *fast retransmit*, and *fast recovery*. These algorithms use local information, such as the RTO (retransmission timeout) based on the estimated RTT (round-trip time) between the sender and receiver, as well as the variance in this round trip time to implement the congestion control policies. UDP is a connectionless protocol, thus there are no means to control the UDP traffic.

The review of basic networking concepts in this section shows why process-to-process communication incurs a significant overhead. While the raw speed of fiber-optic channels can reach Tbps, the actual transmission rate for end-to-end communication over a wide area network can only be on the order of Gbps, and the latency is on the order of milliseconds. The term “speed” is used informally to describe the maximum data-transmission rate, or the capacity of a communication channel; this ca-

pacity is determined by the physical bandwidth of the channel, and this explains why the term channel “bandwidth” is also used to measure the channel capacity, or the maximum data rate.

6.2 Internet evolution

The Internet is continually evolving under the pressure of its own success and the need to accommodate new applications and a larger number of users. Initially conceived as a data network, a network designed to transport data files, the Internet has morphed into today’s network supporting data streaming and applications with real-time constraints such as the Lambda service offered by the AWS. The discussion in this section is restricted to the aspects of the Internet evolution relevant to cloud computing.

Tier 1, 2, and 3 networks. To understand the architectural consequences of Internet evolution, we discuss first the relationship between two networks. *Peering* means that two networks exchange traffic between each other’s customers freely. *Transit* requires a network to pay another one for accessing the Internet. The term *customer* means that a network is receiving money to allow Internet access.

Based on these relationships, the networks are commonly classified as Tier 1, 2, and 3. A *Tier 1 network* can reach every other network on the Internet without purchasing IP transit or paying settlements; examples of Tier 1 networks are Verizon, ATT, NTT, and Deutsche Telecom; see Fig. 6.4.

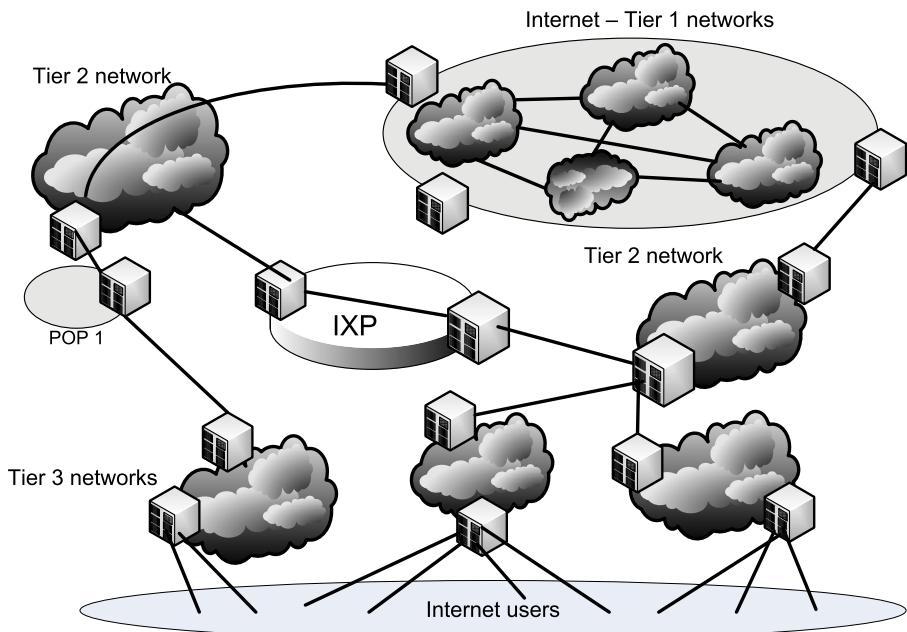
A *Tier 2 network* is an Internet service provider that engages in the practice of peering with other networks but still purchases IP transit to reach some portion of the Internet; Tier 2 providers are the most common providers on the Internet. A *Tier 3 network* purchases transit rights from other networks (typically Tier 2 networks) to reach the Internet. A *point-of-presence (POP)* is an access point from one place to the rest of the Internet.

An *Internet exchange point (IXP)* is a physical infrastructure enabling *Internet Service Providers (ISPs)* to exchange Internet traffic. IXPs interconnect networks directly via the exchange, rather than through one or more third-party networks. The advantages of the direct interconnection are numerous, but the primary reasons to implement an IXP are cost, latency, and bandwidth. Traffic passing through an exchange is typically not billed by any party, whereas traffic to an ISP’s upstream provider is.

IXPs reduce the portion of an ISP’s traffic that must be delivered via their upstream transit providers, thereby reducing the average per-bit delivery cost of their service. Furthermore, the increased number of paths found through the IXP improves routing efficiency and fault tolerance. A typical IXP consists of one or more network switches, to which each of the participating ISPs connects.

New technologies, such as web applications, cloud computing, and content-delivery networks, are reshaping the definition of a network. The World Wide Web, gaming, and entertainment are merging, and more computer applications are moving to the cloud. Data streaming consumes an increasingly larger fraction of the available bandwidth as high definition TV sets become less expensive and content providers such as Netflix and Hulu offer customers services that require a significant increase in the network bandwidth.

Does the network infrastructure keep up with the demand for bandwidth? A natural question to ask is: Where is the actual bottleneck limiting the bandwidth available to a typical Internet broadband user? The answer is: the “last mile,” the link connecting home to the ISP network. Recognizing that the broadband access infrastructure ensures continual growth of the economy and enables people to

**FIGURE 6.4**

There are three classes of networks, Tier 1, 2, and 3; an IXP is a physical infrastructure allowing ISPs to exchange Internet traffic.

work from any site, Google initiated in early 2010 a project to provide 1 Gbps access to individual households through FTTH.¹ ATT FTTH expects to have some seven million subscribers by 2022.

Migration to IPv6. Internet Protocol, Version 4 (IPv4), provides an addressing capability of 2^{32} , or approximately 4.3 billion IP addresses, a number that proved to be insufficient. Indeed, the Internet Assigned Numbers Authority (IANA) assigned the last batch of *five* address blocks to the Regional Internet Registries in February 2011, officially depleting the global pool of completely fresh blocks of addresses; each of the address blocks represents approximately 16.7 million possible addresses. Network Address Translation (NAT) provides a partial solution to this problem because it allows a single public IP address to support hundreds or even thousands of private IP addresses.

Internet Protocol Version 6 (IPv6) provides an addressing capability of 2^{128} , or 3.4×10^{38} addresses. There are other major differences between IPv4 and IPv6:

¹ The fiber-to-the-home (FTTH) is a broadband network architecture that uses optical fiber to replace the copper-based local loop used for the last-mile network access to homes.

1. *Multicasting.* IPv6 does not implement traditional IP broadcast, i.e., the transmission of a packet to all hosts on the attached link using a special broadcast address and, therefore, does not define broadcast addresses. IPv6 supports new multicast solutions, including embedding rendezvous point addresses in an IPv6 multicast group address. This solution simplifies the deployment of inter-domain solutions.
2. *Stateless address auto-configuration (SLAAC).* IPv6 hosts can configure themselves automatically using router-discovery messages when connected to a routed IPv6 network using ICMPv6 (Internet Control Message Protocol v6). When first connected to a network, a host sends a link-local router solicitation multicast request for its configuration parameters. If suitably configured, routers respond to such a request with a router advertisement packet that contains network-layer configuration parameters.
3. *Mandatory support for network security.* Internet Network Security (IPsec) is an integral part of the base protocol suite in IPv6, while it is optional for IPv4. IPsec is a protocol suite² operating at the IP layer. Each IP packet is authenticated and encrypted. Other security protocols, e.g., Secure Sockets Layer (SSL), Transport Layer Security (TLS), and Secure Shell (SSH) operate at the upper layers of the TCP/IP suite.

Migration to IPv6 is a very challenging and costly proposition because it involves upgrading applications, hosts, routers, and DNS infrastructure. Moving to IPv6 requires backward compatibility, and any organization migrating to IPv6 should maintain a complete IPv4 infrastructure [113].

6.3 TCP congestion control

Network congestion occurs when a router or a communication channel is forced to carry out more traffic than it can handle. Congestion leads to increased delays, packet loss, and retransmissions, blocking new connections, and ultimately to lower network capacity. Increased round-trip time of a *connection-oriented protocol* signals network congestion and could trigger congestion control mechanisms. *Connectionless protocols* like UDP cannot detect network congestion and cannot contribute to congestion control or avoidance.

TCP is a connection-oriented Internet transport protocol. TCP uses two mechanisms, flow control and congestion control, to achieve high channel utilization, avoid congestion, and, at the same time, ensure a fair sharing of the network bandwidth. The *flow-control* mechanism throttles the sender feedback from the receiver and forces the sender to transmit only the amount of data the receiver is able to buffer and then process. TCP uses a sliding-window flow-control protocol. If W is the window size, then the sender data rate S is:

$$S = \frac{W \times MSS}{RTT} \text{ bps} = \frac{W}{RTT} \text{ packets/second} \quad (6.1)$$

² IPsec uses several protocols: (1) Authentication Headers (AH) supports connectionless integrity, data origin authentication for IP datagrams, and protection against replay attacks; (2) Encapsulating Security Payloads (ESP) supports confidentiality, data-origin authentication, connectionless integrity, an anti-replay service, and limited traffic-flow confidentiality; (3) Security Associations (SA) provide the parameters necessary to operate the AH and/or ESP operations.

where MSS and RTT denote the Maximum Segment Size and the Round-Trip Time, respectively, assuming that MSS is one packet. If S is too small, the transmission rate is smaller than the channel capacity, while a large S leads to congestion. Channel capacity depends on the network load because the physical channels along the path of a flow are shared with many other Internet flows.

The actual window size W is affected by two factors: (a) the ability of the receiver to accept new data and (b) the sender's estimation of the available network capacity. The receiver specifies the amount of additional data it is willing to accept in the *receive window* field of every frame. The receiver's window shifts when the receiver receives and acknowledges a new segment of data. When a receiver advertises a window size of zero, the sender stops sending data and starts the *persist timer*. This timer is used to avoid deadlock when a subsequent window-size update from the receiver is lost.

When the persist timer expires, the sender sends a small packet, and the receiver responds by sending another acknowledgment containing the new window size. In addition to the flow control provided by the receiver, the sender attempts to infer the available network capacity and to avoid overloading the network. The source uses the losses and the delay to determine the level of congestion. If $awnd$ denotes the receiver window and $cwnd$ the congestion window set by the sender, the actual window should be:

$$W = \min(cwnd, awnd). \quad (6.2)$$

Several algorithms are used to calculate $cwnd$, including Tahoe and Reno, developed by Van Jacobson in 1988 and 1990. Tahoe was based on slow start (SS), congestion avoidance (CA), and fast retransmit (FR). The sender probes the network for spare capacity and detects congestion based on loss. *Slow start* means that the sender starts with a window of two times MSS, $init_cwnd = 1$. For every packet acknowledged, the congestion window increases by one MSS so that the congestion window effectively doubles for every RTT.

When the congestion window exceeds the threshold, *i.e.*, $cwnd \geq ssthresh$, the algorithm enters the CA state. In the CA state, on each successful acknowledgment $cwnd \leftarrow cwnd + 1/cwnd$ and on each RTT $cwnd \leftarrow cwnd + 1$. *Fast retransmit* is motivated by the fact that the timeout is too long thus, a sender retransmits immediately after three duplicate acknowledgments without waiting for a timeout. Two adjustments are made in this case:

$$flightsize = \min(awnd, cwnd) \quad \text{and} \quad ssthresh \leftarrow \max(flightsize/2, 2), \quad (6.3)$$

and the system enters in the slow start state, $cwnd = 1$. The pseudocode describing the Tahoe algorithm is:

```

for every ACK {
    if (W < ssthresh) then W++      (SS)
    else      W += 1/W                (CA)
}
for every loss {
    ssthresh = W/2
    W   = 1
}

```

While the majority of the Internet traffic is due to long-lived, bulk-data transfers, *e.g.*, video and audio streaming, the majority of transactions, *e.g.*, web requests, are short-lived. So, a major challenge is to ensure some fairness for short-lived transactions.

To overcome the limitations of the slow start, application strategies have been developed to reduce the time needed to download data over the Internet. For example, two browsers, Firefox 3 and Google Chrome, open up to six TCP connections per domain to increase the parallelism and to boost start-up performance when downloading a web page. Internet Explorer 8 opens 180 connections. Clearly, these strategies circumvent the mechanisms for congestion control and incur a considerable overhead. A better solution is to increase the initial congestion window of TCP for several reasons [155]:

1. TCP latency is dominated by the number of RTT's during the slow start phase. Increasing the *init_cwnd* parameter enables the data transfer to be completed with fewer RTTs.
2. A single TCP connection requires multiple RTTs to download a single page because the average page size is 384 KB.
3. This will ensure fairness between short-lived transactions, which are a majority of Internet transfers, and long-lived transactions, which transfer very large amounts of data.
4. This will allow faster recovery after losses through Fast Retransmission.

It can be shown that the latency of a transfer completing during the slow start without losses is given by the expression:

$$\left\lceil \log_{\gamma} \left(\frac{L(\gamma - 1)}{init_cwnd} + 1 \right) \right\rceil \times RTT + \frac{L}{C}, \quad (6.4)$$

with L being the transfer size, C is the bottleneck-link rate, and γ is a constant equal to 1.5 or 2 depending on whether the acknowledgments are delayed or not; $L/init_cwnd \geq 1$. In the experiments reported in [155], the TCP latency was reduced from about 490 msec when $init_cwnd = 3$ to about 466 msec for $init_cwnd = 16$.

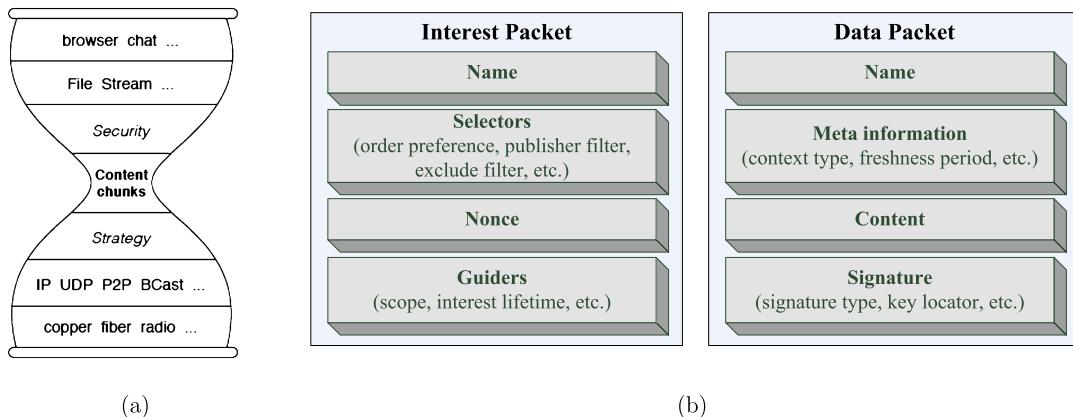
6.4 Content-centric networks; named data networks (R)

Internet data is tied to a particular host, and packets only name the communication end points; this makes data replication and migration difficult. The Internet is mostly used by applications in digital media, electronic commerce, Big Data analytics, etc. as a *distribution network*. In a distribution network, communication is *content-centric*, named objects are requested by an agent, and, once a site holding the object is located, the network transports it to the site that requested the object. The end user in a distribution network is oblivious to the location of the data and is only interested in the content.

The idea of a content-centric network has been around for some time. In 1999, the TRIAD project at Stanford³ [210,479] proposed to use an object name for routing towards the object replica. In this proposal, the Internet Relay Protocol performs name-to-address conversion using routing information maintained by relay nodes. The Name-Based Routing Protocol performs a function similar to the BGP protocol, supporting a mechanism for updating routing information in relay nodes.

In 2006 the Data Oriented Network Architecture (DONA) project at UC Berkeley [149] extended TRIAD incorporating security and persistence as primitives of the new network architecture. DONA architecture exposes two primitives: FIND—allows a client to request a particular piece of data by

³ TRIAD stands for Translating Relaying Internet Architecture Integrating Active Directories.

**FIGURE 6.5**

Named Data Networks. (a) NDN hourglass architecture parallels the one of the Internet; it separates the lower layers of the protocol stack from the upper ones, thus data naming can evolve independently from networking. (b) NDN networking service semantics is to fetch a data chunk identified by name, while Internet semantics is to deliver a packet to a given network address through an end-to-end channel identified by a source and a destination IP address. Interest packets and data packets are used for NDN routing and forwarding.

its name and REGISTER—enables content providers to indicate their intent to serve a particular data object. In 2006 Van Jacobson from UCLA argued that Named Data Networks (NDNs) should be the architecture of the future Internet.

In 2012, the Internet Research Task Force (IRTF) established an information-centric networking (ICN) research working group for investigating architecture designs for NDN. A 2014 survey of ICN research and a succinct presentation of NDNs can be found in [522] and [539], respectively. The important features of the NDN architecture addressed by the current research efforts include namespaces, trust models, in-network storage, data synchronization, and, last but not least, rendezvous, discovery, and bootstrapping.

The hourglass model can be extended, and packets can name objects rather than communication endpoints. NDN's hourglass model separates the lower layers of the protocol stack from the upper ones, thus data naming can evolve independently from networking, as in Fig. 6.5(a). NDN packet delivery is driven by data consumers; communication is initiated by an agent generating an *Interest* packet containing the name of the data as in Fig. 6.5(b).

Once the *Interest* packet reaches a network host that has a copy of the data item, a *Data* packet containing the name, the contents of the data, and the signature is generated. The signature consists of producer's key. The *data* packet follows the route traced by the *interest* packet, and it is delivered to the data consumer agent. The information necessary to forward a packet is maintained by an NDN router in several places:

1. *Content Store*—local cache for *Data* packets previously crossing the router. When an *interest* packet arrives, a search of the content store determines if a matching data exists, and, if so, the data is forwarded on the same router interface the *interest* packet was received. If not, the router uses

a data structure, the Forwarding Information Base (FIB), to forward the packet. More recently, NDN supports more persistent and larger-volume in-network storage, called Repositories, providing services similar to that of the Content Delivery Networks.

2. *Forwarding Information Base*—FIB entries are populated by a name-prefix-based procedure. The Forwarding Strategy retrieves the longest prefix matched entry from the forwarding information base for an interest packet.
3. *Pending Interest Table* (PIT)—stores all interest packets the router has forwarded but have not yet been satisfied. A PIT entry records the data name carried in the Internet, together with its incoming and outgoing router interface(s). When a data packet arrives, the router finds the matching PIT entry and forwards the data to all downstream interfaces listed in that PIT entry; then, it removes the PIT entry and caches the data packet in the Content Store.

Names are essential in NDN, though namespace management is not part of the NDN architecture. The scope and contexts of NDN names are different. Globally accessible data must have globally unique names, while local names require only local routing and must only be locally unique.

There are significant differences between NDN and TCP/IP. NDN namespace cannot be exhausted, while IPv4 address limitation forced migration to IPv6. NDN supports data-centric security; each *Data* packet is cryptographically signed, while TCP/IP security is left to the communication endpoints. An NDN router announces name prefixes for data it is willing to serve, while an IP router announces IP prefixes.

NDN is a universal overlay.⁴ NDN can run over any datagram network and, conversely, any datagram network, including IP, can run over NDN. Classical algorithms such as Open Shortest Path First (OSPF) route data chunks using component-wise, longest-prefix match of the name in an interest packet with the FIB entries of a router.

Wide-area applications can operate over IP tunnels, and islands of NDN nodes could be interconnected by tunneling over non-NDN clouds. Scalable forwarding, along with robust and effective trust management, are critical challenges for the future of NDN networks. Namespace design is at the heart of NDNs because it involves application data, communication, storage models, routing, and security.

NDN could be useful for cloud data centers and, in particular, for those supporting the IaaS cloud delivery model. Data is replicated, and multiple instances could access data concurrently from their nearest storage servers. This would be useful for some MapReduce applications, but it would require major changes in frameworks supporting this paradigm. On the other hand, many applications caching data at the server side would only marginally benefit from NDN support.

6.5 Software-defined networks; SD-WAN

Software-defined Networks (SDNs) allow programmatic control of communication networks extending the basic principles of resource virtualization to networking. SDNs introduce an abstraction layer separating network configuration from the physical communication resources. A network operating system running inside the control layer, sandwiched between the application layer and the infrastructure layer,

⁴ An overlay network is a network built on top of another physical network; its nodes are connected by virtual links, each of which corresponds to a path, possibly through many physical links, in the underlying network.

enables applications to reconfigure dynamically the communication substrate to adapt to their security, scalability, and manageability needs.

Though enthusiastically embraced by networking vendors, there is little agreement either on the architecture, the APIs, or the overlay networks among vendors. A 2012 white paper of the Open Network Foundation (ONF) (<https://www.opennetworking.org/sdn-resources>) promoted OpenFlow-based SDNs. According to ONF, a new network architecture, supporting traffic patterns changes due to extensive use of cloud services and the need for more bandwidth demanded by Big Data applications is necessary. ONF architecture includes several components:

- a. Applications—generate network requirements to controllers.
- b. Controllers—translate application requirements to SDN datapaths and provide an abstract view of the network to applications.
- c. Datapaths—logical network devices that expose control to the physical substrate, consisting of agents and forwarding engines.
- d. Control-to-data-plane—interfaces between SDN controllers and SDN datapaths.
- e. Northbound interfaces—interfaces between applications and controllers.
- f. Interface drivers and agents—driver—agent pairs.
- g. Management and administration.

OpenFlow is an API for programming data plane switches. The datapath and the control paths of an OpenFlow switch consist of a flow table, including an action associated with each flow entry and a controller that programs the flow entry. The controller configures and manages the switch and receives events from the switch.

Software-Defined Wide Area Network (SD-WAN). Informally, SD-WAN is a network where there is a separation between the data service and the rest of the protocol stack. SD-WAN changes WAN from being a hardware-centric network to a software-defined service. An SD-WAN creates a virtual overlay across the underlying data services, separating the upper protocol stack from the network. This separation allows organizations to use the best transport for each application. Transport services used by SD-WANs are:

1. Multiprotocol Label Switching (MPLS)—network routing that directs data from one node to the next based on short path labels rather than network addresses. MPLS avoids complex routing table lookups and speeds traffic flows.
2. Long-Term Evolution/4G/5G. LTE is associated with 4G and 5G wireless networks.
3. Broadband Internet connections.
4. DSL/Cable.

SD-WAN expands network through software rather than hardware and provides a simple interface for WAN management. SD-WAN is a cloud-scale architecture, open and scalable. SD-WAN uses a centralized control function to direct traffic and allows for more consistent, predictable app performance, and improved security. SD-WAN deployment can include existing router/switches, or even virtualized customer-premises equipment.

SD-WAN main advantages: (i) optimized user experience and efficiency for cloud applications; (ii) reduced cost for IT management due to automation and reduced connectivity costs by moving from expensive MPLS to lower-cost connections; (iii) improved organization traffic management; (iv) includes next generation firewalls with malware protection; (v) increased network security with encrypted

network traffic and network segmentation to limit and manage attacks; (vi) automated common tasks and configuration; (vii) redundancy to improve availability and reduce risk/failure; and (viii) easy deployment.

Citrix, Riverbed <https://www.riverbed.com/faq/what-is-sd-wan.html>, CISCO, and Cato Cloud <https://www.catonetworks.com/sd-wan/> are SD-WAN vendors. CISCO offers Meraki <https://meraki.cisco.com/products/security-sd-wan/> as well as Viptela SD-WANs. Citrix <https://www.citrix.com> SD-WAN optimize delivery of virtual, cloud, and SaaS applications on a secure, flexible WAN.

6.6 Interconnection networks for computer clouds

Computing and communication are deeply intertwined, as we have seen in Chapter 3. Interconnection networks discussed in the next sections are critical for the performance of computer clouds and supercomputers. Concepts important for understanding interconnection networks are introduced next.

A network consists of nodes and links or communication channels. The *degree* of a node is the number of links the node is connected with. The *nodes* of an interconnection network could be processors, memory units, or servers. The *network interface* of a node is the hardware connecting it to the network. *Switches* and *communication channels* are the elements of the interconnection fabric. Switches receive data packets, examine each packet to identify the destination IP addresses, and then use the routing tables to forward it to the next hop towards its final destination. An *n-way switch* has n ports connected to n communication links. An interconnection network can be *non-blocking* if it is possible to connect any permutation of sources and destinations at any time or *blocking* if this requirement is not satisfied.

While processor and memory technology have followed Moore's law, interconnection networks have evolved at a slower pace and have become a major factor in determining the overall performance and cost of large-scale systems. For example, from 1997 to 2010, the speed of Ethernet networks has increased from one to 100 Gbps. This increase is slightly slower than the Moore's law for traffic [359] which predicted a 1-Tbps Ethernet by 2013.

Interconnection networks are distinguished by their topology, routing, and flow control. *Network topology* is determined by the way nodes are interconnected, routing decides how a message gets from its source to destination, and the flow control negotiates how the buffer space is allocated. There are two basic types of network topologies:

- Static networks where there are direct connections between servers;
- Switched networks where switches are used to interconnect the servers.

The topology of an interconnection network determines the *network diameter*, the average distance between all pairs of nodes, the *bisection width*, the minimum number of links cut to partition the network into two halves, and the bisection bandwidth, as well as the cost and power consumption [273]. When a network is partitioned into two networks of the same size, the bisection bandwidth measures the communication bandwidth between the two.

The *full bisection bandwidth* allows one half of the network nodes to communicate simultaneously with the other half of the nodes. Assume that half of the nodes inject data into the network at a rate B Mbps. When the bisection bandwidth is B , then the network has full bisection bandwidth. The most popular topologies with a static interconnect are:

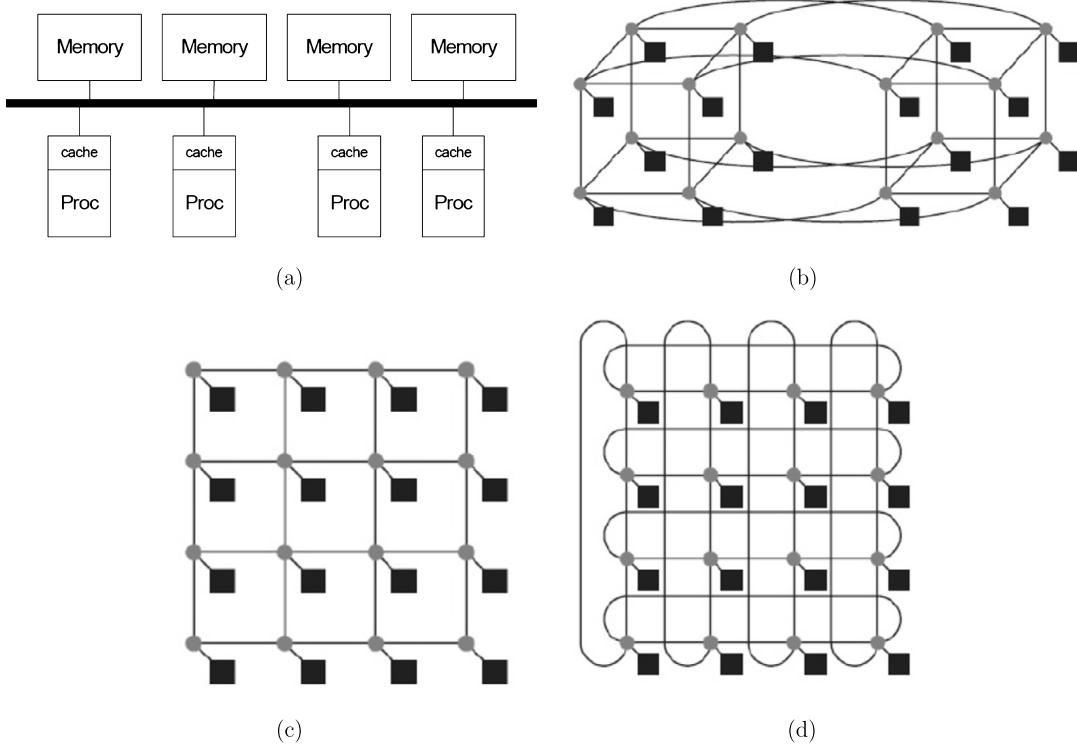


FIGURE 6.6

Static networks: (a) Bus; (b) Hypercube; (c) 2D-mesh; (d) 2D-torus.

- Bus**, a simple and cost-effective network, see Fig. 6.6(a). It does not scale, but it is easy to implement cache coherence through snooping for distributed memory systems. A bus is often used in shared-memory multiprocessor systems.
 - Hypercube** of order n ; see Fig. 6.6(b). A hypercube has a good bisection bandwidth; the number of nodes is $N = 2^n$, the degree is $n = \log N$, and the average distance between nodes is $\mathcal{O}(N)$ hops. Example of use: SGI Origin 2000.
 - 2D-mesh**; see Fig. 6.6(c). An $n \times n$ 2D-mesh has many paths to connect nodes, a cost of $\mathcal{O}(n)$, and an average latency of $\mathcal{O}(\sqrt{n})$. A mesh is not symmetric on edges, thus its performance is sensitive to the placement of communicating nodes on edges versus the middle. Example of use: Intel Paragon supercomputer of the 1990s.
 - Torus**, avoids mesh asymmetry but has a higher cost due to a larger number of components. A torus is effective for applications using nearest-neighbor communication; see Fig. 6.6(d). It is prevalent for proprietary interconnects. Example of use: Six-D Mesh/torus of Fujitsu K supercomputer.

Switched networks have multiple layers of switches connecting the nodes: (i) a crossbar switch has N^2 crosspoint switches; see Fig. 6.7(a); (ii) Omega (Butterfly, Benes, Banyan, etc.) have $(N \log N)/2$

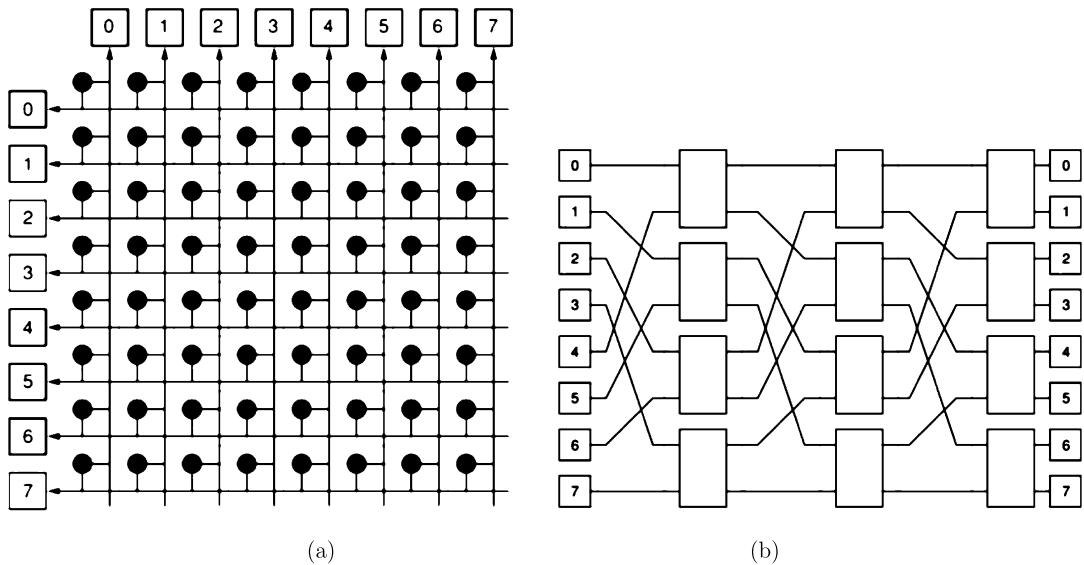


FIGURE 6.7

Switched networks: (a) An 8×8 crossbar switch. Sixteen nodes are interconnected by 49 switches represented by the dark circles; (b) An 8×8 Omega switch. Sixteen nodes are interconnected by 12 switches represented by white rectangles.

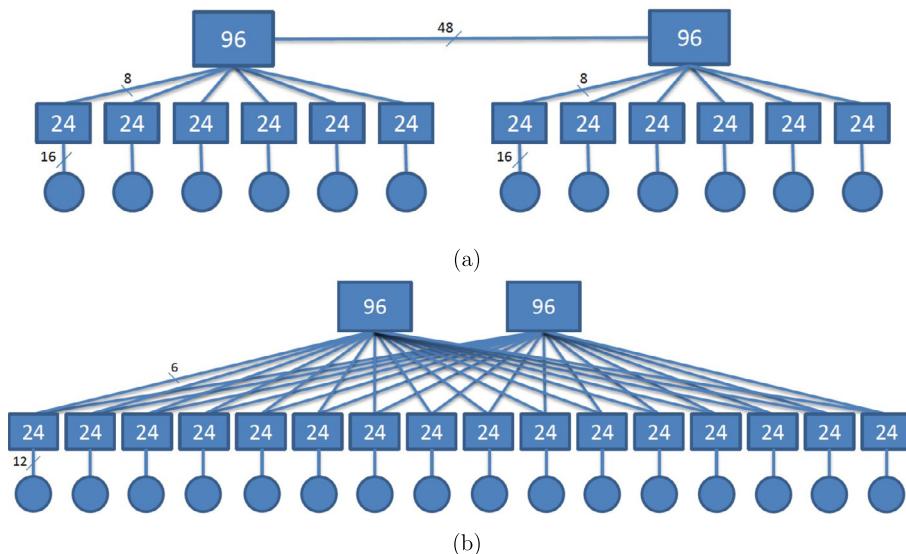
switches; see Fig. 6.7(b). The cost is $\mathcal{O}(N \log N)$ and the latency is $\mathcal{O}(\log N)$. Omega networks are small diameter networks.

Cloud interconnection networks. The cloud infrastructure consists of one or more warehouse-scale computers discussed in Section 4.2. WSCs have an infrastructure with a very large number of servers interconnected by high-speed networks.

The networking infrastructure of a cloud must satisfy several requirements, including scalability, cost, and performance. The network should allow low-latency, high-speed communication, and, at the same time, provide *location transparent communication* between servers. In other words, each server should be able to communicate with every other server with similar speed and latency. This requirement ensures that *applications need not be location aware*, and, at the same time, it reduces the complexity of the system management.

Typically, the networking infrastructure is organized hierarchically. The servers of a WSC are packed into racks and interconnected by a rack router. Then, rack routers are connected to cluster routers, which in turn are interconnected by a local communication fabric. Finally, inter-data center networks connect multiple WSCs [279]. Clearly, in a hierarchical organization, true location transparency is not feasible. Cost considerations ultimately decide the actual organization and performance of the communication fabric.

Oversubscription, defined as the ratio of the worst-case achievable aggregate bandwidth among the servers to the total bisection bandwidth of the interconnect, is a useful measure of the fitness of an

**FIGURE 6.8**

One hundred ninety-two node fat-tree interconnection network with two 96-way and 12 24-way switches in a computer cloud. (a) The two 96-way switches at the root are connected with one another and with six 24-way switches; (b) Each 96-way switch is connected to each one of the 12 24-way switch.

interconnection network for a large cluster. An oversubscription of one-to-one indicates that any host may communicate with an arbitrary hosts at the full bandwidth of the interconnect. An oversubscription value of four to one means that only 25% of the bandwidth available to servers can be attained for some communication patterns. Typical oversubscription figures are in the 2.5 to 1 and 8 to 1 range.

The cost of routers and of cables interconnecting the routers are major components of the overall cost of an interconnection network. Wire density has scaled up at a slower rate than processor speed, and wire delay has remained constant over time, thus better performance and lower costs can only be achieved with innovative router architecture. This motivates us to take a closer look at routers design.

Fat-trees are optimal interconnects for large-scale clusters and, by extension, for WSCs. Servers are placed at the leaves, and switches populate the root and the internal nodes of the tree. Fat-trees have additional links to increase the bandwidth near the root of the tree. A fat-tree interconnect can be built with inexpensive commodity parts since all switching elements of a fat-tree are identical. A fat-tree is a particular instance of Clos topology discussed in Section 6.7.

Two 192 node fat-trees built with two 96-way switches and 12 24-way switches are shown in Fig. 6.8. The two 96-way switches at the root are connected via 48 links and each 24-way switch has 6×8 uplink connections to the root and 6×16 down connections to 16 servers in Fig. 6.8(a), while in Fig. 6.8(b) each 96-way switch is connected via six links to each one of the 12 24-way switches connected via 12 links with servers.

Table 6.1 from [232] summarizes the performance/cost for a 2D-mesh, a 2D-torus, a hypercube of order seven, a fat-tree, and a fully connected network. Two dimensions of interconnection network

Table 6.1 A side-by-side comparison of performance and cost figures of several interconnection network topologies for 64 nodes.

Property	2D mesh	2D torus	Hypercube	Fat-tree	Fully connected
BW in # of links	8	16	32	32	1 024
Max/Avg hop count	14/7	8/4	6/3	11/9	1/1
I/O ports per switch	5	5	7	4	64
Number of switches	64	64	64	192	64
Total number of links	176	192	256	384	2 080

performance, the bisection bandwidth and the average and maximum number of hops, along with three elements affecting the cost, the number of ports per switch, the number of switches, and the total number of links, are shown. Fat-tree interconnects have the largest bisection bandwidth with the smallest number of I/O ports per switch, while a fully connected interconnect has a prohibitively large number of links.

6.7 Multistage interconnection networks

Low-radix routers have a small number of ports. *High-radix* routers have a large number of ports, a considerably smaller number of stages and enjoy lower latency and reduced power consumption. High-radix chips divide the bandwidth into larger number of narrow ports, while low-radix chips divide the bandwidth into a smaller number of wide ports. Every five years during the past two decades, the pin bandwidth of the chips used for switching has increased by an order of magnitude as a result of signaling rate increase and larger number of signals.

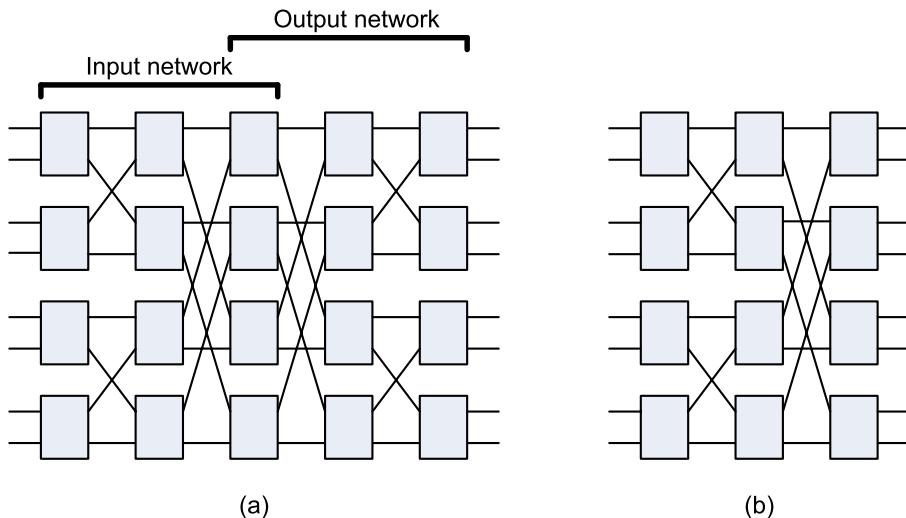
Clos networks were invented in early 1950s by Charles Clos from Bell Labs [110]. A Clos network is a multistage nonblocking network with an odd number of stages; see Fig. 6.9(a). The network consists of two butterfly networks where the last stage of the input is fused with the first stage of the output.

The name *butterfly* comes from the pattern of inverted triangles created by the interconnections, which look like butterfly wings. A butterfly network transfers data packets using the most efficient route, but it is blocking, it cannot handle a conflict between two packets attempting to reach the same port at the same time. In a Clos network, all packets overshoot their destination and then hop back to it. Most of the time, the overshoot is not necessary and increases the latency, so a packet takes twice as many hops as it really needs.

In a *folded Clos* topology the input and the output networks share switch modules. Such networks are sometimes called *fat-tree*; many commercial high-performance interconnects, such as Myrinet, InfiniBand, and Quadrics, implement a fat-tree topology. Some folded Clos networks use low-radix routers, e.g., the Cray XD1 uses radix-24 routers. The latency and the cost of the network can be lowered using high-radix routers.

The *Black Widow* topology extends the folded Clos topology and has a lower cost and latency; it adds side links, and this permits a statical partitioning of the global bandwidth among peer subtrees [443]. The Black Widow topology is used in Cray computers.

Flattened butterfly network topology [272] is similar to the generalized hypercube proposed in early 1980s and it is able to exploit high-radix routers. In a flattened butterfly, we start with a conventional

**FIGURE 6.9**

(a) A five-stage Clos network with radix-2 routers and unidirectional channels equivalent to two back-to-back butterfly networks. (b) The corresponding folded-Clos network with bidirectional channels; input and output networks share switch modules.

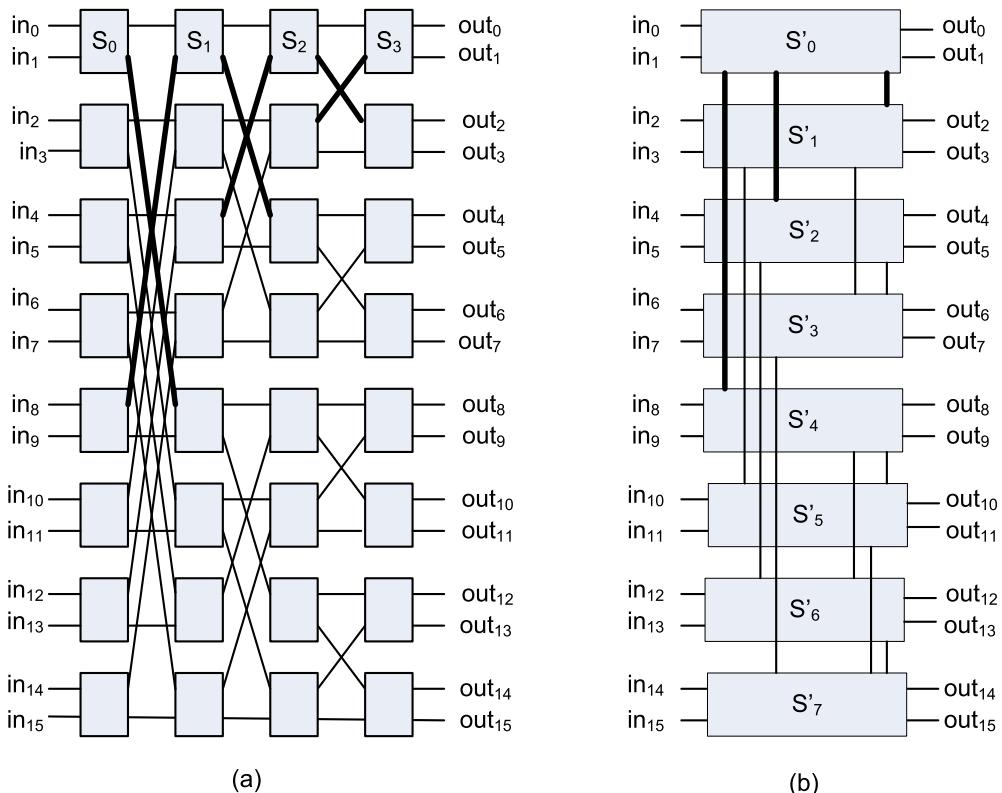
butterfly and combine the switches in each row into a single, higher-radix one. Each router is linked to more processors, and this halves the number of router-to-router connections.

Data sent by one processor can reach its destination with fewer hops and latency is reduced, though the physical path may be longer. For example, in Fig. 6.10(a), we see a two-ary four-fly butterfly; we combine the four switches S_0 , S_1 , S_2 and S_3 in the first row into a single switch, S'_0 . The flattened butterfly adaptively senses congestion and overshoots only when it needs to. On adversarial traffic patterns, the flattened butterfly has a similar performance as the folded Clos but provides over an order of magnitude increase in performance compared to the conventional butterfly.

The network cost for computer clusters can be reduced by a factor of two when high-radix routers (radix-64 or higher), and the flattened butterfly topology are used [273]. Flattened butterfly topology does not reduce the number of local cables, e.g., backplane wires from the processors to routers, but it reduces the number of global cables. The cost of the cables represents as much as 80% of the total network cost, e.g, for a 4K system, the cost savings of the flattened butterfly exceed 50%.

6.8 InfiniBand and Myrinet

InfiniBand is an interconnection network used by supercomputers and computer clouds backed by top companies in the industry, including Dell, HP, IBM, Intel, and Microsoft. Intel manufactures InfiniBand host bus adapters and network switches. InfiniBand uses a switched fabric rather than a shared communication channel and has a high throughput and low latency. Every link of the fabric has exactly one

**FIGURE 6.10**

(a) A two-ary four-fly butterfly with unidirectional links. (b) The corresponding two-ary four-flat flattened butterfly is obtained by combining the four switches S₀, S₁, S₂ and S₃ in the first row of the traditional butterfly into a single switch S'₀ and by adding additional connections between switches [272].

device connected at each end of the link, thus the worst case is the same as the typical case. A switched fabric architecture is fault tolerant and scalable, while in a shared channel architecture, all connected devices share the same bandwidth and the larger the number of devices, the less bandwidth is available to each one of them.

InfiniBand architecture implements the “bandwidth-out-of-the-box” concept and delivers bandwidth traditionally trapped inside a server across the interconnect fabric. InfiniBand supports several signaling rates, and the energy consumption depends on the throughput. Links can be bonded together for additional throughput, as shown in Table 6.2. InfiniBand specifications define multiple operational modes: single data rate (SDR), double data rate (DDR), quad data rate (QDR), 14 data rate (FDR), and enhanced data rated (EDR).

The signaling rates are: 2.5 Gbps in each direction per connection for an SDR connection; 5 Gbps for DDR; 10 Gbps for QDR; 14.0625 Gbps for FDR; 25.78125 Gbps for EDR per lane. The SDR, DDR,

Table 6.2 Evolution of high-speed interconnects. Several networking equipment suppliers offered proprietary solutions for 200 GE and 400 GE in 2016. Ethernet Alliance's 2020 technology roadmap expects 800 Gbps and 1.6 Tbps to become IEEE standard after 2020.

Network	Year	Speed
Gigabit Ethernet (GE)	1995	1 Gbps
10-GE	2001	10 Gbps
40-GE	2010	40 Gbps
100-GE	2010	100 Gbps
Myrinet	1993	1 Gbps
Fiber Channel	1994	1 Gbps
InfiniBand (IB)	2001	2 Gbps (1X SDR)
	2003	8 Gbps (4X SDR)
	2005	16 Gbps (4X DDR) & 24 Gbps (12X SDR)
	2007	32 Gbps (4X QDR)
	2011	56 Gbps (4X FDR)
	2012	100 Gbps (4X EDR)

and QDR link encoding is 8 B/10 B, every 10 bits sent carry 8 bits of data. Thus single, double, and quad data rates carry 2, 4, or 8 Gbit/s useful data, respectively, because the effective data transmission rate is four-fifths the raw rate.

InfiniBand allows links to be configured for a specified speed and width; the reactivation time of the link can vary from several nanoseconds to several microseconds. Exadata and Exalogic systems from Oracle implement the InfiniBand QDR with 40 Gbit/s (32 Gbit/s effective). InfiniBand fabric is used to connect compute nodes, compute nodes with storage servers, and Exadata and Exalogic systems.

In addition to high throughput and low latency, InfiniBand supports quality of service guarantees and failover—the capability to switch to a redundant or standby system. It offers point-to-point bidirectional serial links intended for the connection of processors with high-speed peripherals, such as disks, as well as multicast operations.

Initially, InfiniBand deployment pushed the limitations of personal computers (PCs) buses. Introduced as a standard PC architecture in the early 90s, the PCI bus has evolved from 32 bit/33 MHz to 64 bit/66 MHz, while PCI-X has doubled the clock rate to 133 MHz. PCI Express 2.0/x16 has a 16 Gbps bandwidth. The number of pins of a parallel bus is quite large; the number of pins per connection is large, e.g., a 64 bit PCI bus requires 90 pins.

Myrinet is an interconnect for massively parallel systems developed at Caltech. It is no longer in use, though it has several remarkable features [65]:

- a. Robustness ensured by communication channels with flow control, packet framing, and error control.
- b. Self-initializing, low-latency, cut-through switches.
- c. Host interfaces that can map the network, select routes, and translate from network addresses to routes, as well as handle packet traffic.

- d. Streamlined host software that allows direct communication between user processes and the network.

Myrinet design benefits from the experience gathered by a project to construct a high-speed local-area network at USC/ISI. A Myrinet is composed of point-to-point, full-duplex links connecting hosts and switches, an architecture similar to the one of the ATOMIC project at USC. Multiple-port switches may be connected to other switches and to the single-port host interfaces in any topology. Transmission is synchronous at the sending end, while the receiver circuits are asynchronous. The receiver injects control symbols in the reverse channel of the link for flow control. Myrinet switches use blocking-cut-through routing similar to the one in Intel Paragon and Cray T3D.

Myrinet supports high data rates; a Myrinet link consists of a full-duplex pair of 640 Mbps channels and has regular topology with elementary routing circuits in each node. The network is scalable, and its aggregate capacity grows with the number of nodes. Algorithmic routing enables multiple packets to flow concurrently and avoids deadlocks.

Store and forward and *cut-through or wormhole* networks are very different. In the former, an entire packet is buffered, and its checksum is verified in each node along the path from the source to the destination. In wormhole networks the packet is forwarded to its next hop as soon as the header is received and decoded. This decreases the latency, but a packet can still experience blocking if the outgoing channel expected to carry it to the next node is in use. In this case the packet has to wait until the channel becomes free.

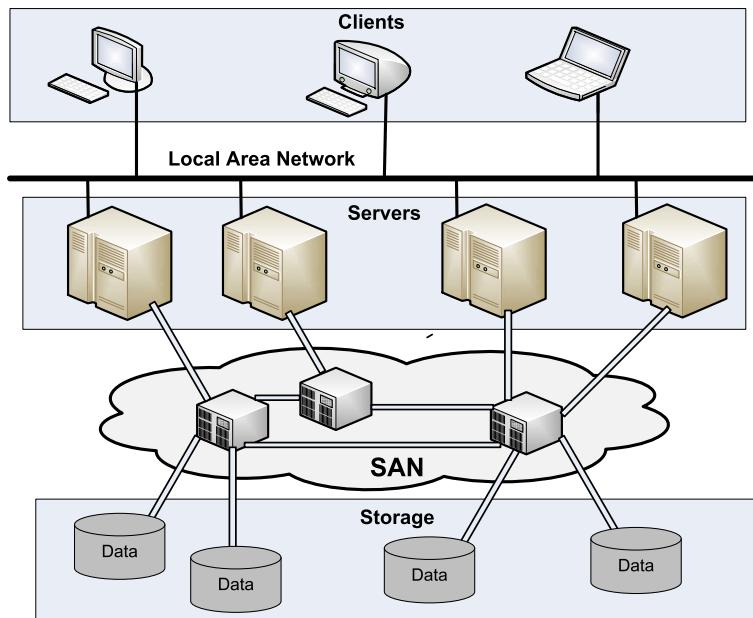
A comparison of Myrinet and Ethernet performance as a communication substrate for MPI libraries in [323] shows that the MG NAS benchmark,⁵ over Myrinet, only achieves 5% higher performance than TCP over Ethernet. MPI library implementations for Ethernet have a higher message latency and lower message bandwidth because they use the OS network protocol stack.

6.9 Storage area networks and the Fibre Channel

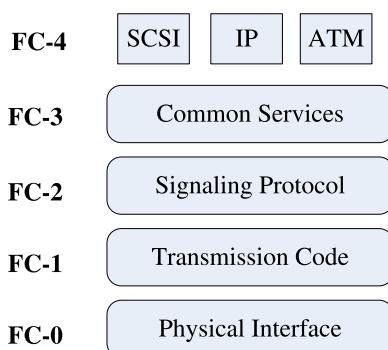
A storage area network (SAN) is a specialized, high-speed network for data block transfers. SANs interconnect servers to servers, servers to storage devices, and storage devices to storage devices; see Fig. 6.11. A SAN consists of a communication infrastructure and a management layer. The Fibre Channel (FC) is the dominant architecture of SANs. FC is a layered protocol with several layers, as depicted in Fig. 6.12, discussed next.

- A. Lower layer protocols: FC-0 physical interface; FC-1 transmission protocol responsible for encoding/decoding; and FC-2 signaling protocol responsible for framing and flow control. FC-0 uses laser diodes as the optical source and manages point-to-point fiber connections; when a fiber connection is broken, the ports send a series of pulses until the physical connection is reestablished and the necessary handshake procedures are followed. FC-1 controls the serial transmission and integrates

⁵ NASA Parallel Benchmarks, NAS, are used to evaluate the performance of parallel supercomputers. The original benchmark included five kernels: IS—Integer Sort, random memory access; EP—Embarrassingly Parallel; CG—Conjugate Gradient; MG—Multi-Grid on a sequence of meshes, long- and short-distance communication, memory intensive; and FT—discrete 3D Fast Fourier Transform, all-to-all communication.

**FIGURE 6.11**

A storage area network interconnects servers to servers, servers to storage devices, and storage devices to storage devices. Typically, it uses fiber optics and the FC protocol.

**FIGURE 6.12**

FC protocol layers. FC-4 supports communication with Small Computer System Interface (SCSI), IP, and Asynchronous Transfer Mode (ATM) network interfaces.

data with clock information. It ensures encoding to the maximum length of the code, maintains DC-balance, and provides word alignment. FC-2 provides the transport methods for data transmitted in 4-byte ordered sets containing data and control characters. It handles the topologies, based on the presence or absence of a fabric, the communication models, the classes of service provided by the fabric and the nodes, sequence and exchange identifiers, and segmentation and reassembly.

- B.** Upper layer protocols: FC-3 common services layer; and FC-4 protocol mapping layer. FC-3 supports multiple ports on a single-node or fabric using:

- Hunt groups—sets of associated ports assigned an alias identifier that allows any frame containing that alias to be routed to any available port within the set.
- Striping to multiply bandwidth, using multiple ports in parallel to transmit a single information unit across multiple links.
- Multicast and broadcast to deliver a single transmission to multiple destination ports or to all ports.

FC supports several classes of service to accommodate various application needs:

Class 1: rarely used blocking connection-oriented service; acknowledgments ensure that the frames are received in the same order in which they are sent, and reserves full bandwidth for the connection between the two devices.

Class 2: acknowledgments ensure that the frames are received; enables the fabric to multiplex several messages on a frame-by-frame basis; since frames can take different routes, it does not guarantee in-order delivery, rather it relies on upper layer protocols to take care of frame sequence.

Class 3: datagram connection; no acknowledgments.

Class 4: connection-oriented service. Virtual circuits (VCs) established between ports, guarantees in-order delivery and acknowledgment of delivered frames. The fabric is responsible for multiplexing frames of different VCs. Supports Guaranteed Quality of Service (QoS), including bandwidth and latency. This layer is intended for multimedia applications.

Class 5: isochronous service for applications requiring immediate delivery, without buffering.

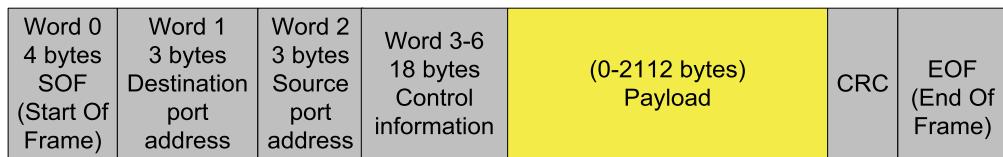
Class 6: supports dedicated connections for a reliable multicast.

Class 7: similar to Class 2 but used for the control and management of the fabric; a connectionless service with notification of non-delivery.

While every device connected to a LAN has a unique MAC address,⁶ each FC device has a unique Id called the WWN (World Wide Name), a 64 bit address. Every port in the switched fabric has its own unique 24-bit address consisting of: the domain (bits 23–16), the area (bits 15–08), and the port physical address (bits 07–00).

The switch of a switched fabric environment assigns dynamically and maintains the port addresses. When a device with a WWN address logs into a given port, the switch maintains the correlation between that port address and the WWN address of the device using the Name Server. The Name Server is a component of the fabric operating system, which runs inside the switch.

⁶ Typically, a virtualization platform generates a new, random MAC address for each virtual network interface at creation. Moreover, in Linux, it is possible to spoof the MAC address.

**FIGURE 6.13**

The format of an FC frame; the payload can be at most 2 112 bytes; and larger data units are carried by multiple frames.

The format of an FC frame is shown in Fig. 6.13. Zoning permits finer segmentation of the switched fabric; only the members of the same zone can communicate within that zone. It can be used to separate different environments, e.g., a Microsoft Windows NT from a UNIX environment.

SANs use several other protocols, for example, Fibre Channel over IP (FCIP) allows transmission of Fibre Channel data through IP networks using tunneling. *Tunneling* allows a network protocol to carry a payload over an incompatible delivery-network, or to provide a secure path through an untrusted network. A protocol normally blocked by a firewall is wrapped inside a protocol that the firewall does not block. For example, an HTTP tunnel can be used for communication from network locations with restricted connectivity, e.g., behind NATs, firewalls, or proxy servers. Restricted connectivity is a commonly used method to lock down a network to secure it against internal and external threats.

Internet Fibre Channel Protocol (iFCP) is a gateway-to-gateway protocol supporting communication among FC storage devices in a SAN, or on the Internet, using TCP/IP; iFCP replaces lower-layer Fibre Channel transport with TCP/IP and Gigabit Ethernet. Fibre Channel devices connect to an iFCP gateway or switch, and each Fibre Channel session is terminated at the local gateway and converted to a TCP/IP session via iFCP.

6.10 Scalable data center communication architectures

The question addressed now is: How does one organize the communication infrastructure of large cloud data centers for best performance at the lowest cost? *Data center networks* (DCNs) architectures providing an answer to this question face major challenges: (i) aggregate cluster bandwidth scales poorly with the cluster size; (ii) the bandwidth needed by many cloud applications comes at a high price, and the interconnect cost increases dramatically with cluster size; and (iii) the communication bandwidth of DCNs may become oversubscribed by a significant factor depending on communication patterns.

We only mention two DCN architectural styles, *three-tier* and *fat-tree* DCNs. The former has a multiple-rooted tree topology with three layers: core, aggregate, and access. Servers are connected to switches at tree leaves at the *access layer*. Enterprise switches at the root of the tree form the *core layer* and connect switches together at the *aggregate layer* connecting the data center to the Internet. The uplinks of the *aggregate layer* switches connect them to the core layer, and their download links connect to the access layer. Three-tier DCN architecture is not suitable for computer clouds, it is not scalable, the bisection bandwidth is not optimal, and core layer switches are expensive and power-hungry.

Fat-tree topology is optimal for computer clouds, the bandwidth is not severely affected for messages crossing multiple switches, and the interconnection network can be built with commodity rather than enterprise switches. The implementation of the fat-tree topology proposed in [20] is based on the following principles:

- i. The network should scale to a very large number of nodes.
- ii. The fat-tree should have multiple core switches.
- iii. The network should support multipath routing. Equal-cost multi-path (ECMP) routing algorithm [244] performing static load splitting among flows should be used.
- iv. The network should use switches with optimal cost/performance ratios.

Hierarchical and fat-tree interconnection networks cost per GigE⁷ has decreased by one order of magnitude and the cost/performance indicator for fat-tree built with commodity switches is almost an order of magnitude lower than that of hierarchical networks. The choice of multirouted fat-tree topology and multipath routing is justified because in 2008 the largest cluster that could be supported with a single rooted core 128-port router with 1:1 oversubscription would be limited to 1 280 nodes.

A WSC interconnect can be organized as a fat-tree with k -port switches, $k = 48$, but the organization can be supported for any k . The network consists of k pods⁸ each pod has two layers and $k/2$ switches at each layer. Each switch at the lower layer is connected directly to $k/2$ servers. The other $k/2$ ports are connected to $k/2$ of k ports in the aggregation layer. The total number of switches is $k(k + 1)$ and the total number of servers connected to the system is k^2 ; $(k/2)^2$ paths connect each pair of servers.

A WSC with 16 384 servers can be built with 128-port switches, and one with 262 144 servers requires 512-port switches. Fig. 6.14 shows a fat-tree interconnection network for $k = 4$. The core, the aggregation, and the edge layers are populated with *four*-port switches. Each core switch is connected to one of the switches at the aggregation layer of each pod. The network has four pods, four switches at each pod, two at aggregation layer and two at the edge. Four servers are connected to each pod.

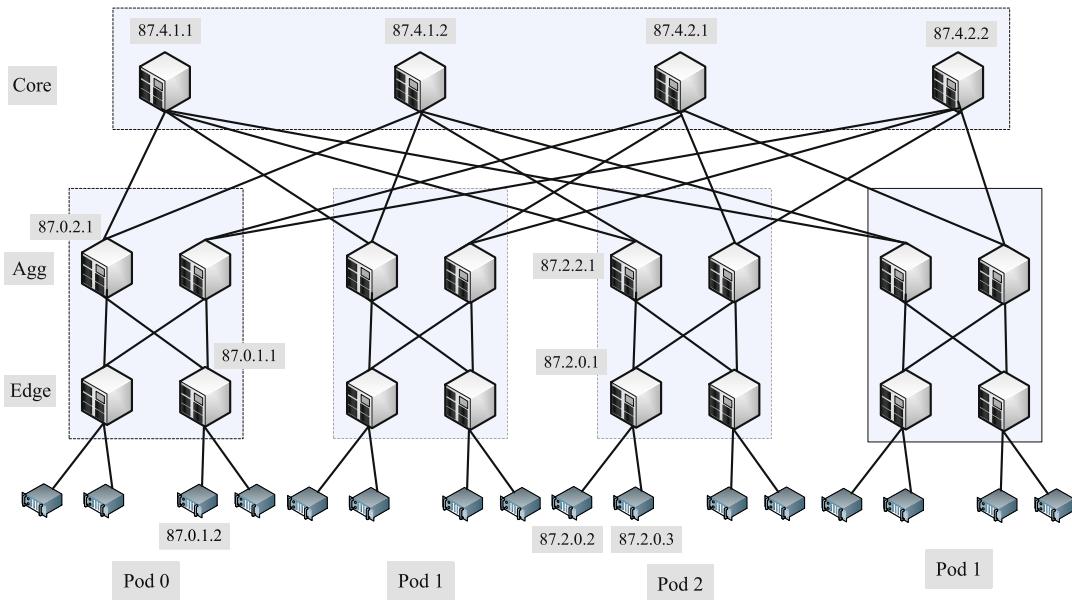
Switch IP addresses are $87.pod.switch.1$; switches are numbered left to right, and bottom to top. Core switches IP addresses are $10.k.j.i$, where j and i denote switch coordinates in the $(k/2)^2$ core switch grid starting from top-left. For example, the four switches of pod 2 have IP addresses 87.2.0.1, 87.2.1.1, 87.2.2.1, and 87.2.3.1. Server IP addresses are $87.pod.switch.serverID$, where $serverID$ is server position in the subnet of the edge router starting from left to right. For example, the IP addresses of the two servers connected to the switch with IP address 87.2.0.1 are 87.2.0.2 and 87.2.0.3.

There are multiple path between any pair of servers. For example, packets sent by server with IP address 87.0.1.2 to server with IP address 87.2.0.3 can follow the following routes:

$$\begin{aligned}
 & 87.0.1.1 \mapsto 87.0.2.1 \mapsto 87.4.1.1 \mapsto 87.2.2.1 \mapsto 87.2.0.1 \\
 & 87.0.1.1 \mapsto 87.0.2.1 \mapsto 87.4.1.2 \mapsto 87.2.2.1 \mapsto 87.2.0.1 \\
 & 87.0.1.1 \mapsto 87.0.1.1 \mapsto 87.4.2.1 \mapsto 87.2.2.1 \mapsto 87.2.0.1 \\
 & 87.0.1.1 \mapsto 87.0.1.1 \mapsto 87.4.2.2 \mapsto 87.2.2.2 \mapsto 87.2.0.1
 \end{aligned} \tag{6.5}$$

⁷ The IEEE 802.3-2008 standard defines GigE as a technology for transmitting Ethernet frames. One GigE corresponds to a rate of one gigabit per second, i.e., 10^9 bits per second.

⁸ A pod is a repeatable design pattern to maximize the modularity, scalability, and manageability of data center networks.

**FIGURE 6.14**

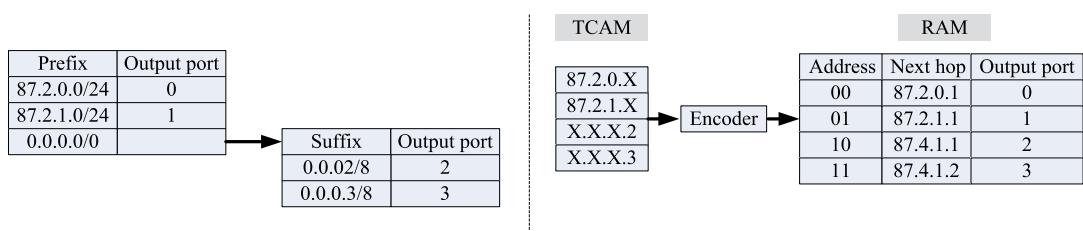
A fat-tree network with k -port switches and $k = 4$. Four core four-way switches are at the root, and there are four pods, two switches at the aggregation layer, and two at the edge layer of each pod. Each switch at the edge of a pod is connected to two servers.

Packet routing uses two-level routing tables and supports two-level prefix lookup. This strategy may increase the lookup latency, but the prefix search can be done in parallel and compensate for the latency increase. The main table entries are of the form $(prefix, outputport)$ and could have an additional pointer to a secondary table or could be terminating if none of its entries point to secondary table. A secondary table consists of $(suffix, outputport)$ entries and may be pointed to by more than one first-level entry.

Fig. 6.15 (Left) shows the two-level routing tables for switch 87.2.2.1 and routing for two incoming packets for servers with the IP addresses 87.2.1.2 and 87.3.0.3; the incoming packets are forwarded on ports 1 and 3, respectively. Lookup engines use a ternary version of content-addressable memory (CAM), called TCAM. Fig. 6.15 (Right) shows that TCAM stores address prefixes and suffixes, which index a RAM that stores the IP address of the next hop and the output port.

The prefix entries are stored with numerically smaller addresses first and the right-handed (suffix) entries in larger addresses. The output of the CAM is encoded so that the entry with the numerically smallest matching address is the output. When the destination IP address of a packet matches both a left-handed and a right-handed entry, then the left-handed entry is chosen.

The k switches in a pod have terminating prefixes to the subnets in that pod. When two servers in the same pod, but on a different subnet communicate, all upper-level switches of the pod will have a terminating prefix pointing to the destination subnet's switch.

**FIGURE 6.15**

(Left) The two-level routing tables for switch 87.2.2.1. Two incoming packets for IP addresses 87.2.1.2 and 87.3.0.3 are forwarded on ports 1 and 3, respectively. (Right) The RAM implementation of a two-level TCAM routing table.

For all outgoing pod traffic, the pod switches have a default /0 prefix with a secondary table matching the least-significant byte of the destination IP address, the server ID. Traffic diffusion occurs only in the first half of a packet's journey. Once a packet reaches a core switch, there is exactly one link to its destination pod, and that switch will include a terminating /16 prefix for the pod of that packet ($87.\text{pod}.0.0/16, \text{port}$). Once a packet reaches its destination pod, the receiving upper-level pod switch will also include a $(10.\text{pod}.\text{switch}.0/24, \text{port})$ prefix to direct the packet to its destination subnet switch, where it is finally switched to its destination server.

Each pod switch assigns terminating prefixes for subnets in the same pod and adds a /0 prefix with a secondary table matching the *serverIDs* for inter-pod traffic. The routing tables for the upper pod switches are generated with the pseudocode of Algorithm 6.1.

Algorithm 6.1: Generates aggregation switch routing tables

```

1 foreach pod x  $\in [0, k - 1]$  do
2     foreach switch z  $\in [k/2, k - 1]$  do
3         foreach subnet i  $\in [0, k/2 - 1]$  do
4             addPrefix(10.x.z.1, 10.x.i.0/24, i);
5         end
6         addPrefix(10.x.z.1, 0.0.0.0/0, 0);
7         foreach hostID i  $\in [2, (k/2) + 1]$  do
8             addSuffix(10.x.z.1, 0.0.0.i/8, (i - 2 + z)  $\bmod (k/2) + (k/2)$ );
9         end
10    end
11 end

```

For the lower pod switches, the /24 subnet prefix in line 3 is omitted since that subnet's own traffic is switched, and intra- and interpod traffic should be evenly split among the upper switches. Core switches contain only terminating /16 prefixes pointing to their destination pods, as shown in Algorithm 6.2. The maximum numbers of first-level prefixes and second-level suffixes are k and $k/2$, respectively.

Flow classification with dynamic port-reassignment in pod switches overcomes cases of local congestion when two flows compete for the same output port, even though another port that has the same

Algorithm 6.2: Generates routing tables for core switches

```

1 foreach  $j \in [0, k - 1]$  do
2     foreach  $i \in [1, k/2]$  do
3         foreach destination pod  $\in [0, k/2 - 1]$  do
4             addPrefix(10.k.j.i, 10.x.0.0/16, x);
5         end
6     end
7 end

```

cost to the destination is underused. Power consumption and heat dissipation are major concerns for the cloud data centers. Switches at the higher tiers of the interconnect in data centers consume several kW, and the entire interconnection infrastructure consumes hundreds to thousands of kW.

6.11 Network resource management algorithms (R)

A critical aspect of resource management in cloud computing is to guarantee the communication bandwidth required by an application as specified by an SLA. The solutions to this problem are based on the strategies used for some time on the Internet to support the data streaming QoS requirements.

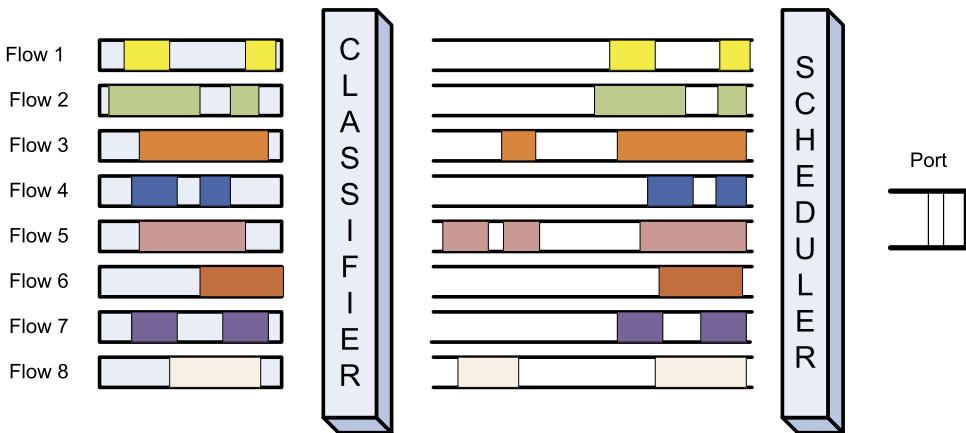
Cloud interconnects consist of communication links of limited bandwidth and switches of limited capacity. When the load exceeds its capacity, a switch starts dropping packets because it has limited input buffers for the switching fabric and for the outgoing links, as well as limited CPU cycles. A switch must handle multiple flows, pairs of source-destination endpoints of the traffic, thus a scheduling algorithm has to manage several quantities at the same time: the *bandwidth*, the amount of data each flow is allowed to transport, the *timing* when the packets of individual flows are transmitted, and the *buffer space* allocated to each flow.

Communication and computing require scheduling, therefore it should be no surprise that the first algorithm we discuss can be used for scheduling packets transmission, as well as threads. A first strategy to avoid network congestion is to use a FCFS scheduling algorithm. The advantage of FCFS algorithm is a simple management of the three quantities: bandwidth, timing, and buffer space. Nevertheless, the FCFS algorithm does not guarantee fairness; greedy flow sources can transmit at a higher rate and benefit from a larger share of the bandwidth.

Fair Queuing (FQ). The algorithm ensures that a high-data-rate flow cannot use more than its fair share of the link capacity. Packets are first classified into flows by the system and then assigned to a queue dedicated to the flow. Packet queues are serviced one packet at a time in round-robin (RR) order as depicted in Fig. 6.16. FQ's objective is *max-min* fairness. This means that it maximizes first the minimum data rate and then the second minimum data rate of any data flow. Starvation of expensive flows is avoided, but the throughput is low.

The FQ algorithm guarantees buffer space management fairness but does not guarantee fairness of bandwidth allocation; a flow transporting large packets will benefit from a larger bandwidth [360].

The FQ algorithm in [138] proposes a solution to this problem. First, it introduces a *Bit-by-bit Round-robin (BR)* strategy; as the name implies, in this rather impractical scheme, a single bit from

**FIGURE 6.16**

Fair Queuing—packets are first classified into flows by the system and then assigned to a queue dedicated to the flow; queues are serviced one packet at a time in round-robin order and empty queues are skipped.

each queue is transmitted, and the queues are visited in a round-robin fashion. Let $R(t)$ be the number of rounds of the BR algorithm up to time t and $N_{active}(t)$ be the number of active flows through the switch. Call t_i^a the time when the packet i of flow a , of size P_i^a bits arrives, and call S_i^a and F_i^a the values of $R(t)$ when the first and the last bit, respectively, of the packet i of flow a are transmitted. Then,

$$F_i^a = S_i^a + P_i^a \quad \text{and} \quad S_i^a = \max[F_{i-1}^a, R(t_i^a)]. \quad (6.6)$$

$R(t)$, $N_{active}(t)$, S_i^a , and F_i^a in Fig. 6.17 depend only on the arrival time of the packets, t_i^a , and not on their transmission time, provided that a flow a is active as long as

$$R(t) \leq F_i^a \quad \text{when } i = \max(j | t_i^a \leq t). \quad (6.7)$$

The authors of [138] use for the packet-by-packet transmission time the following nonpreemptive scheduling rule which emulates the BR strategy: *the next packet to be transmitted is the one with the smallest F_i^a .* A preemptive version of the algorithm requires that the transmission of the current packet be interrupted as soon as one with a shorter finishing time, F_i^a , arrives.

A fair allocation of the bandwidth does not have an effect on the timing of the transmission. A possible strategy is to allow less delay for the flows using less than their fair share of the bandwidth. The same paper [138] proposes the introduction of a quantity called the *bid*, B_i^a , and scheduling the packet transmission based on its value. The bid is defined as

$$B_i^a = P_i^a + \max[F_{i-1}^a, (R(t_i^a) - \delta)], \quad (6.8)$$

with δ being a nonnegative parameter. The properties of the FQ algorithm, as well as the implementation of a nonpreemptive version of the algorithm, are analyzed in [138].

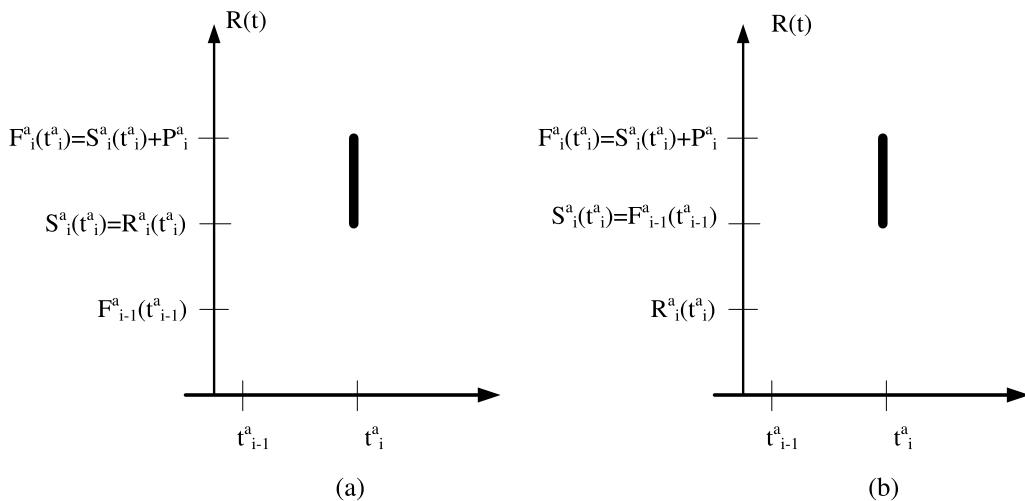


FIGURE 6.17

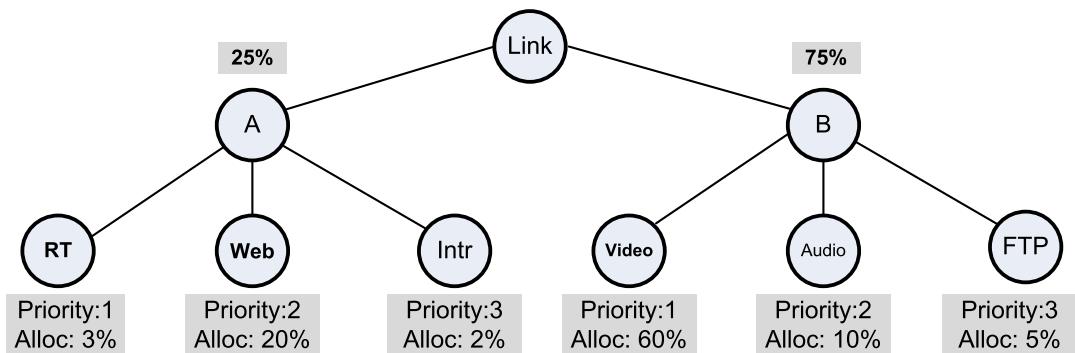
Transmission of packet i of flow a arriving at time t_i^a of size P_i^a bits. Transmission starts at time $S_i^a = \max[F_{i-1}^a, R(t_i^a)]$ and ends at time $F_i^a = S_i^a + P_i^a$ with $R(t)$ the number of rounds of the algorithm. (a) The case $F_{i-1}^a < R(t_i^a)$. (b) The case $F_{i-1}^a \geq R(t_i^a)$.

Stochastic Fairness Queueing (SFQ) requires less calculations; it is a simpler and less accurate implementation of FQ algorithms. SFQ ensures that each flow has the opportunity to transmit an equal amount of data and takes into account data packet sizes [343].

Class-Based Queuing (CBQ). This algorithm is a widely used strategy for link sharing proposed by Sally Floyd and Van Jacobson in 1995 [181]. The objective of CBQ is to support flexible link sharing for applications that require bandwidth guarantees such as VoIP, video streaming, and audio streaming. At the same time, CBQ supports some balance between short-lived network flows, such as web searches, and long-lived ones, such as video streaming or file transfers.

CBQ aggregates connections and constructs a hierarchy of classes with different priorities and throughput allocations. CBQ, uses several functional units for link sharing: (i) a *classifier* uses the information in the packet header to assign arriving packets to classes; (ii) an *estimator* of the short-term bandwidth for the class; (iii) a *selector*, or scheduler, identifies the highest priority class to send next and, if multiple classes have the same priority, to schedule them on a round-robin base; and (iv) a *delayer* to compute the next time when a class that has exceeded its link allocation is allowed to send.

Classes are organized in a tree-like hierarchy; for example, in Fig. 6.18, group *A* corresponds to short-lived traffic and group *B* to long-lived traffic. The leaves of the tree are considered Level 1, and this example includes six classes of traffic: real-time, web, interactive, video streaming, audio streaming, and file transfer. At Level 2, there are the two classes of traffic, *A* and *B*. The root at Level 3 is the link itself. The link-sharing policy aims to ensure that, if sufficient demand exists, then, after some time intervals, each interior or leaf class receives its allocated bandwidth. The distribution of the “excess” bandwidth follows a set of guidelines but does not support mechanisms for congestion avoidance.

**FIGURE 6.18**

CBQ link sharing for two groups A, of short-lived traffic, and B, of long-lived traffic, allocated 25% and 75% of the link capacity, respectively. There are six classes of traffic with priorities 1, 2, and 3. The RT (real-time) and the video streaming have priority 1 and are allocated 3% and 60%, respectively, of the link capacity. Web transactions and audio streaming have priority 2 and are allocated 20% and 10%, respectively of the link capacity. Intr (interactive applications) and FTP (file transfer protocols) have priority 3 and are allocated 2% and 5%, respectively, of the link capacity.

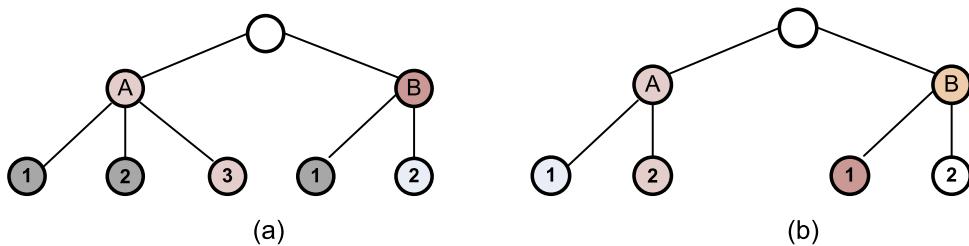
A class is *overlimit* if, over a certain recent period, it has used more than its bandwidth allocation (in bytes per second), *underlimit* if it has used less, and *atlimit* if it has used exactly its allocation. A leaf class is *satisfied* if it is underlimit, has a persistent backlog, and is *unsatisfied* otherwise; a non-leaf class is unsatisfied if it is underlimit and has some descendent class with a persistent backlog. A precise definition of the term “persistent backlog” is part of a local policy. A class does not need to be *regulated* if it is underlimit or if there are no unsatisfied classes; the class should be regulated if it is overlimit and if some other class is unsatisfied, and this regulation should continue until the class is no longer overlimit or until there are no unsatisfied classes; see Fig. 6.19 for two examples.

Linux kernel implements a link sharing algorithm called *Hierarchical Token Buckets* (HTB) inspired by CBQ. In CBQ every class has an *assured rate* (AR); in addition to AR every class in HTB has also a *ceil rate* (CR)⁹; see Fig. 6.20. The main advantage of HTB over CBQ is that it allows *borrowing*. If a class C needs a rate above its AR, it tries to borrow from its parent; then the parent examines its children and, if there are classes running at a rate lower than their AR, the parent can borrow from them and reallocate it to class C.

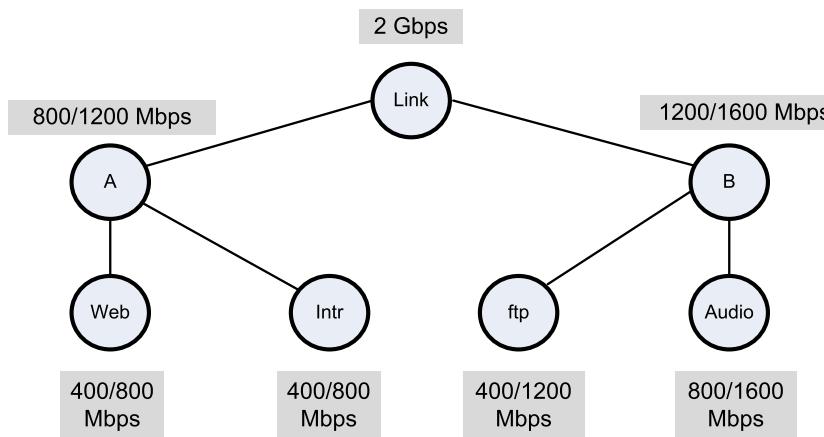
6.12 Content delivery networks

Computer clouds support network-centric computing and network-centric content. The vast amount of data stored on the cloud has to be delivered efficiently to a large user population and Content Delivery

⁹ The ceil rate is the rate a class is allowed to borrow from other classes.

**FIGURE 6.19**

There are two groups *A* and *B* and three types of traffic, e.g., web, real-time, and interactive, denoted as 1, 2, and 3.
 (a) Group *A* and class *A.3* traffic are underlimit and unsatisfied; classes *A.1*, *A.2*, and *B.1* are overlimit, unsatisfied and with persistent backlog, and have to be regulated; type *A.3* is underlimit and unsatisfied; group *B* is overlimit.
 (b) Group *A* is underlimit and unsatisfied; Group *B* is overlimit and needs to be regulated; class *A.1* traffic is underlimit; class *A.2* is overlimit and with persistent backlog; class *B.1* traffic is overlimit and with persistent backlog and needs to be regulated.

**FIGURE 6.20**

HTB packet scheduling adds a ceil rate to the allowed rate of every node.

Networks (CDNs) offer fast and reliable content delivery and reduce communication bandwidth by caching and replication. A CDN receives the content from an *Origin* server, then replicates it to its *Edge* cache servers; the content is delivered to an end-user from the “closest” Edge server.

CDNs support scalability, increase reliability and performance, and provide better security. The volume of transactions and data transported by the Internet increases dramatically every year; additional resources are necessary to accommodate the additional load placed on the communication and storage systems and to improve the end-user experience. CDNs place additional resources pro-

visioned to absorb the traffic caused by *flash crowds*¹⁰ and, in general, to provide capacity on demand.

The additional resources are placed strategically throughout the Internet to ensure scalability. The resources provided by a CDN are replicated, and when one replica fails, the content is available from another one; replicas are “close” to the consumers of the content, and this placement reduces the start-up time and the communication bandwidth. A CDN uses two types of servers: the *origin* server updated by the content provider and *replica* servers, which cache the content and serve as authoritative reference for client requests. Security is a critical aspect of the services provided by a CDN; the replicated content should be protected from the increased risk of cyber fraud and unauthorized access.

CDNs deliver static content and/or live or on-demand streaming media. *Static content* refers to media stored using traditional caching technologies because changes are infrequent. Examples of static content are: HTML pages, images, documents, software patches, and audio and/or video files. *Live media* refers to live events when the content is delivered in real time from the encoder to the media server. On-demand delivery of audio and/or video streams, movie files, and music clips provided to the end-users is content-encoded and then stored on media servers. Virtually all CDN providers support static content delivery, while live or on-demand streaming media is considerably more challenging.

CDN providers and protocols. The first CDN was setup by *Akamai*, a company that evolved from an MIT project to optimize network traffic. Since its inception Akamai has placed some 20 000 servers in 1 000 networks in 71 countries; in 2009, it controlled some 85% of the market [391]. Akamai mirrors the contents of clients on multiple systems placed strategically through the Internet. Though the domain name (but not subdomain) is the same, the IP address of the resource requested by a user points to an Akamai server rather than the customer’s server. Then the Akamai server is automatically picked depending on the type of content and the network location of the end user.

There are several other active commercial CDNs including EdgeStream providing video streaming and Limelight Networks providing on-demand and live delivery of video, music, and games. There are several academic CDNs: Coral is a freely available network designed to mirror web content, hosted on PlanetLab; Globule is an open-source collaborative CDN developed at Vrije Universiteit in Amsterdam.

The communication infrastructure among different CDN components uses a fair number of protocols including: Network Element Control Protocol (NECP), Web Cache Coordination Protocol (WCCP), SOCKS, Cache Array Routing Protocol (CARP), Internet Cache Protocol (ICP), Hypertext Caching Protocol (HTCP), and Cache Digest described succinctly in [391]. For example, caches exchange ICP queries and reply to locate the best sites to retrieve an object; HTCP is used for discovering HTTP caches, caching data, managing sets of HTTP caches, and monitoring cache activity.

CDN organization, design decisions, and performance. There are two strategies for CDN organization; in the so-called *overlay*, the network core does not play an active role in the content delivery. On the other hand, the *network* approach requires the routers and the switches to use dedicated software to identify specific application types and to forward user’s requests based on predefined policies.

The first strategy is based exclusively on content replication on multiple caches and redirection based on proximity to the end-user. In the second approach the network core elements redirect con-

¹⁰ The term *flash crowd* refers to an event disrupting the life of a significant segment of the population, such as an earthquake in a populated area; it causes the Internet traffic to increase dramatically.

tent requests to local caches or redirect data center's incoming traffic to servers optimized for specific content type access. Some CDNs including Akamai use both strategies.

The most important design and policy decisions for a CDN are: (i) the placement of the edge servers; (ii) the content selection and delivery; (iii) content management; and (iv) request routing policies. The placement problem is often solved with suboptimal heuristics using as input the workload patterns and the network topology. The simplest, but costly, approach for content selection and delivery is the *full-site* replication suitable for static content; the edge servers replicate the entire content of the origin server. On the other hand, the *partial-site* selection and delivery retrieves the base HTML page from the origin server and the objects referenced by this page from the edge caches. The objects can be replicated based on their popularity, or on some heuristics.

Content management depends on the caching techniques, the cache maintenance, and the cache update policies. CDNs use several strategies to manage the consistency of content at replicas: periodic updates, updates triggered by the content change, and on-demand updates. Request-routing in a CDN directs users to the closest edge server that can best serve the request; metrics, such as network proximity, client perceived latency, distance, and replica server load, are taken into account when routing a request. Round-robin is a nonadaptive request-routing that aims to balance the load; it assumes that all edge servers have similar characteristics and can deliver the content.

Adaptive algorithms perform considerably better but are more complex and require some knowledge of state of the system. The algorithm used by *Akamai* takes into consideration metrics such as: the load of the edge server, the bandwidth currently available to a replica server, and the reliability of client connection to edge servers. CDN routing can exploit an organization where several edge servers are connected to a service node aware of the load and the information about each edge server connected to it and attempts to implement a *global load balancing policy*.

An alternative is *DNS-based routing*, when a domain name has multiple IP addresses associated to it and the service provider's DNS server returns the IP addresses of the edge servers holding the replica of the requested object; then, the client's DNS server chooses one of them. Another alternative is the *HTTP redirection*; in this case a web server includes in the HTTP header of a response to a client the address of another edge server. Finally, *IP anycasting* requires that the same IP address is assigned to several hosts and the routing table of a router contains the address of the host closest to it.

Critical CDN performance metrics are: (i) cache hit ratio—the ratio of the number of cached objects versus total number of objects requested; (ii) reserved bandwidth for the origin server; (iii) latency—based on the perceived response time by the end users; (iv) edge server utilization; and, last but not least, (v), reliability—based on packet loss measurements.

CDNs face considerable challenges due to increased appeal of data streaming and to the proliferation of mobile devices such as smartphones and tablets. On-demand video streaming requires enormous bandwidth and storage space, as well as powerful servers; CDNs for mobile networks must be able to dynamically reconfigure the system in response to spatial and temporal demand variations.

Content-Centric Networks (CCNs) are related to information-centric networking architectures, such as Named Data Networks (NDNs) discussed in Section 6.4, where content is named and transferred throughout the network. The request for a named content is routed to the producer or to any entity that can deliver the expected content object.

CCN content may be served from any router's cache. Content is signed by its producer and the consumer is able to verify the signature before actually using the content. CCNs supports opportunistically caching. CCNs offer a number of advantages according to <http://chris-wood.github.io/2015/06/>

Table 6.3 VANETs applications, contents, local interest, local validity, and lifetime.

Application	Contents	Local interest	Local validity	Lifetime
Safety warnings	Dangerous Road	All	100 m	10 s
Safety warnings	Accident	All	500 m	30 s
Safety warnings	Work zone	All	1 Km	Construction
Public service	Emergency vehicle	All	500 m	10 min
Public service	Highway information	All	5 Km	All day
Driving	Road congestion	All	5 Km	30 min
Driving	Navigation Map	Subscribers	5 Km	30 min

[16/CCN-vs-CDN.html](#) i.e., content popularity need not be predicted beforehand. CCNs offer benefits not available from by IP-based CDNs: (i) active and intelligent forwarding strategies for routers; (ii) publisher mobility is easily supported via CCN routing protocols; (iii) congestion control can be enforced within the network; (iv) existing/problematic IP stack can be completely replaced with a new set of layers; (v) existing APIs can be completely reworked to focus on content, not on addresses; and (vi) content security is not tied to the channel but to the content itself.

Content service networks (CSN) introduced in [320] are overlay networks built around CDNs to provide an infrastructure service for processing and transcoding.

6.13 Vehicular ad hoc networks

A vehicular ad hoc network (VANET) consists of groups of moving or stationary vehicles connected by a wireless network. Until recently, the main use of VANETs was to provide safety and comfort to drivers in vehicular environments. This view is changing: Vehicular ad hoc networks are seen now as an infrastructure for an intelligent transportation system with an increasing number of autonomous vehicles and for any activity requiring Internet connectivity in a smart city. Also, VANETs allow on-board computers of mostly stationary vehicles, e.g., vehicles at an airport parking, to serve as resources of a mobile computer cloud with minimum help from the Internet infrastructure.

The contents produced and consumed by vehicles has *local relevance* in terms of time, space, and agents involved, the producer and the consumer. Vehicle-generated information has *local validity*, a limited spatial scope, an *explicit lifetime*, a limited temporal scope, and *local interest*, it is relevant to agents in a limited area around the vehicle. For example, the information that a car is approaching a congested area of a highway is relevant only for that particular segment of the road, at a particular time, and for vehicles nearby. The attributes of vehicular contents are summarized in Table 6.3 from [300].

One of the distinguishing characteristics of VANETs is the *content-centric distribution*, meaning that the content is important but the source is not. This is in marked contrast to the Internet, where an agent demands information from a specific source. For example, traffic information floods a specific area, and a vehicle retrieves it without concern for its source, while an Internet request for highway traffic information is directed to a specific site. Vehicle applications collect sensor data, and vehicles collaborate sharing sensory data. Sensory data is collected by vehicle-installed cameras, i.e., by on-board instruments. For example, *CarSpeak* allows a vehicle to access sensors on neighboring vehicles in

the same manner in which it can access its own [286]. *Waze* is a community-based traffic and navigation application allowing drivers real-time traffic and road information.

VANET communication protocols are similar to the ones used by wired networks; each host has an IP address. Assigning IP addresses to moving vehicles is far from trivial and often requires a Dynamic Host Configuration Protocol (DHCP) server, a heresy for ad hoc networks that operate without any infrastructure, using self-organization protocols. Vehicles frequently join and leave the network, and content of interest cannot be consistently bound to a unique IP address. A router typically relays and then deletes content.

6.14 Further readings

“Brief History of the Internet” [289] was written by the Internet pioneers Barry Leiner, Vinton Cerf, David Clark, Robert Kahn, Leonard Kleinrock, Daniel Lynch, Jon Postel, Larry Roberts, and Stephen Wolff. A widely used text of Kurose and Ross [285] is an excellent introduction to basic networking concepts. The book by Bertsekas and Gallagher [61] gives insights into performance evaluation of computer networks. The classic texts on queuing theory of Kleinrock [276] are required reading for those interested in network analysis.

A survey of information-centric networking research is given in [522], while [539] is a succinct presentation of NDNs. Several publications related to software-defined networks are available on the site of the Open Network Foundation, <https://www.opennetworking.org/sdn-resources>. Moore’s law for traffic is discussed in [359]. Class-based queuing algorithm was introduced by Floyd and Van Jacobson in [181]. The *Black Widow* topology for system interconnects is analyzed in [443]. Storage Area Networks are analyzed in [472].

Scale-free networks and their applications are described by Barabási and Albert in [17–19,46]. The small-worlds networks were introduced by Watts and Strogatz in [507]. Epidemic algorithms for the dissemination of topological information are presented in [215,258,259]. Erdős-Rényi random graphs are analyzed in [67,165]. Energy-efficient protocols for cooperative networks are discussed in [162].

The future of fiber networks is covered in [290]. Vehicle ad hoc networks and their applications to vehicular cloud computing are discussed in [194,300,511]. An analysis of peer-to-peer networks is reported in [195]. Network management for private clouds, P2P networks, and virtual networks are presented in [356], [357], and [364], respectively.

6.15 Exercises and problems

Problem 1. Four ground rules for an open-architecture principles are cited in the “Brief History of the Internet” [289]. Analyze the implication of each one of these rules.

Problem 2. Key issues for network design in [289] are: (1) algorithms to prevent lost packets from permanently disabling communications; (2) host-to-host “pipelining” so that multiple packets could be en-route from source to destination, if intermediate networks allow it; (3) end-end checksums, reassembly of packets from fragments, and detection of duplicates, if any; (4) the need for global addressing; and (5) techniques for host-to-host flow control. Discuss how these issues were addressed by the TCP/IP network architecture.

- Problem 3.** Analyze the challenges of transition to IPv6. What will be in your view the effect of this transition on cloud computing?
- Problem 4.** Discuss the algorithms used to compute the TCP window size.
- Problem 5.** Creating a virtual machine (VM) reduces ultimately to copying a file, therefore the explosion of the number of VMs cannot be prevented; see Section 8.10. Virtualization could drastically lead to an exhaustion of the IPv4 address space. Analyze the solution to this potential problem adopted by the IaaS cloud service delivery model.
- Problem 6.** Read the paper describing the stochastic fair queuing algorithm [181]. Analyze the similarities and dissimilarities of this algorithm with the start-time fair queuing algorithm discussed in Section 9.6.
- Problem 7.** The small-worlds networks were introduced by D. Watts and S. H. Strogatz. They have two desirable features, high clustering and small path length. Read [507] and design an algorithm to construct a small-worlds network.
- Problem 8.** Scale-free networks are analyzed in [17–19,46]. Discuss the important features of systems interconnected by scale-free networks discussed in these papers.
- Problem* 9.** Consider the two 192 node fat-trees interconnect with two 96-way and twelve 24-way switches in Fig. 6.8. Compute the bisection bandwidth of the two interconnects.

Cloud data storage

7

The volume of data generated by human activities is growing at a breathtaking rate. More data was generated in the last two years than in the entire human history before that. Computer clouds provide vast amounts of storage demanded by many applications using these data. Several data sources contribute to the massive amounts of data stored on a cloud. A variety of sensors feed streams of data to cloud applications or simply generate content. An ever-increasing number of cloud-based services collect detailed data about their services and information about the users of these services.

Big Data, discussed in depth in Chapter 12, reflects the reality that many applications use data sets so large that local computers, or even small-to-medium scale data centers, do not have the capacity to store and process such data. Big Data growth can be viewed as a three-dimensional phenomenon as it: (i) implies an increased volume of data; (ii) requires increased processing speed to process more data and produce more results; and (iii) involves a diversity of data sources and data types.

A network-centric data storage model is particularly useful for mobile devices with limited power reserves and local storage, now able to save and access large audio and video files on computer clouds. Billions of Internet-connected mobile and stationary devices access data stored on computer clouds.

Cloud data storage and processing are intimately tied to one another. Data analytics uses large volumes of data collected by an organization to optimize its businesses. In-depth data analysis allows an organization to reach a larger population of customers, identify the strengths of the products or shortcomings in the organization, save energy, and, last but not least, protect the environment.

Applications in many areas of science, including genomics, structural biology, high energy physics, astronomy, meteorology, and the study of the environment, carry out complex analysis of data sets often on the order of terabytes.¹ As a result, file systems, such as Btrfs, XFS, ZFS, exFAT, NTFS, HFS Plus, and ReFS, support disk formats with theoretical volume sizes of several exabytes.

While we emphasize the advantages of a concentration of resources, we have to be acutely aware that a cloud is a large-scale distributed system with a very large number of components that must work in concert. Management of a large collection of storage systems poses significant challenges and requires novel approaches to storage system design. Effective data replication and storage management strategies are critical for cloud applications.

Sophisticated strategies to reduce the access time and to support multimedia access are necessary to satisfy the timing requirements of data streaming and content delivery. Data replication allows concurrent access to data from multiple processors and decreases the chances of data loss. Maintaining consistency among multiple copies of data records increases the complexity of data management software and could negatively affect the storage system performance if data is frequently updated.

¹ Terabyte, 1 TB = 10^{12} bytes; Petabyte, 1 PB = 10^{15} bytes; Exabyte, 1 EB = 10^{18} bytes; Zettabyte, 1 ZB = 10^{21} bytes.

Nowadays, large-scale systems are built with off-the-shelf components, while the distributed file systems of the past used custom-designed reliable components. The storage system design philosophy has shifted from performance-at-any-cost to reliability-at-the-lowest-possible-cost. This shift is evident in the evolution of ideas from the file systems of the 1980s, such as the Network File System (NFS), the Andrew File System (AFS), and the Sprite File System (SFS), to the Google File System (GFS), the Megastore, and the Colossus [173], developed during the last two decades.

The discussion of cloud storage starts with a review of storage technology in Sections 7.1 and 7.2, followed by an overview of storage models in Section 7.3. Evolution of file systems from distributed file system to parallel file systems, then to the file systems capable of handling massive amounts of data, is presented in Sections 7.4 and 7.5, where we discuss distributed file systems, General Parallel File Systems, and Google File System, respectively.

A locking service, Chubby, based on the Paxos algorithm, is presented in Section 7.7. The mismatch between relational database systems and cloud requirements in Section 7.8 is followed by a discussion of NoSQL databases in Section 7.9 and of transaction processing systems in Section 7.10. Sections 7.11 and 7.12 analyze the BigTable and the Megastore systems, respectively. Storage reliability at scale, data center disk locality, and database provenance are discussed in Sections 7.13, 7.14, and 7.15, respectively.

7.1 Dynamic random access memories and hard disk drives

During the last decades, the storage technology has evolved at an accelerated pace, and the volume of data stored every year has constantly increased [236]. Though it pales in comparison to the evolution of processor technology, the evolution of storage technology is astounding. A 2003 study [359] shows that during the 1980–2003 period the storage density of hard disk drives (HDD) has increased by four orders of magnitude from about 0.01 Gb/in² to about 100 Gb/in². During the same period the prices have fallen by five orders of magnitude to about 1 cent/Mbyte. HDD densities were projected to climb to 1 800 Gb/in² by 2016, up from 744 Gb/in² in 2011. Solid-state drives (SSDs) can support hundreds of thousands of I/O operations per second (IOPS).

Dynamic Random Access Memory (DRAM). The density of DRAM (Dynamic Random Access Memory) devices increased from about 1 Gb/in² in 1990 to 100 Gb/in² in 2003. DRAM cost tumbled from \$80/MB to less than \$1/MB during the same period. In April 2020, TSMC (Taiwan Semiconductor Manufacturing Company) started mass production of integrated circuits using 5 nm lithographic process.

Recent advancements in storage technology have a broad impact on the storage systems used for cloud computing. The capacity of NAND flash-based devices outpaced DRAM capacity growth, and the cost per gigabyte has significantly declined. Manufacturers of storage devices are investing in competing solid-state technologies such as Phase-Change Memory. While solid-state memories are based on electron charge, another fundamental property of an electron, its spin, is used to store information. *Spintronics*, an acronym for spin transport electronics, promises storage media based on antiferromagnetic materials insensitive to perturbations by stray fields and with much shorter switching times [518].

While the density of storage devices has increased and the cost has decreased dramatically, the access time has improved only slightly. Performance of I/O subsystems has not kept pace with the

performance of processors, and this performance gap affects multimedia, scientific and engineering, and other applications that process increasingly large volumes of data.

Storage systems face substantial pressure because the volume of data generated has increased exponentially during the last decades. In the 1980s and 1990s, data was primarily generated by humans, nowadays, machines generate data at an unprecedented rate. Mobile devices, such as smartphones and tablets recording static images and movies have limited local storage capacity and rely on transferring the data to cloud storage. Sensors, surveillance cameras, and digital medical imaging devices generate data at a high rate and dump it on storage systems accessible via the Internet.

Online digital libraries, eBooks, and digital media, along with reference data, add to the demand for massive amounts of storage. The term reference data is used for infrequently used data, such as archived copies of medical or financial records, customer account statements, etc.

As the volume of data increases, new methods and algorithms for data mining that require powerful computing systems are being developed. Only a concentration of resources could provide the CPU cycles, along with the vast storage capacity, necessary when performing such intensive computations and when accessing the very large volume of data.

The rapid technological advancements have changed the balance between the initial investment in the storage devices and the system management costs. Now, the cost of storage management is the dominant element of the total cost of a storage system. This effect favors the centralized storage strategy supported by a cloud; indeed, a centralized approach can automate some of the storage management functions such as replication and backup and, thus, reduce substantially the storage management cost.

Hard disk drives (HDDs) are ubiquitous secondary storage media for general-purpose computers. An HDD is a nonvolatile random-access data storage device consisting of one or more rotating platters coated with magnetic material. Magnetic heads mounted on a moving actuator arm read and write data to the surface of the platters.

A typical HDD has a spindle motor and an actuator that positions the read/write head assembly across the spinning disks. Rotation speed of platters in today's HDDs ranges from 4 200 rpm for energy-efficient portable devices, to 15 000 rpm for high-performance servers. HDDs for desktop computers and laptops are 3.5-inch and 2.5-inch, respectively.

HDDs are characterized by capacity and performance. The capacity is measured in Megabytes (MB), Terabytes (TB), or Gigabytes (GB). The average access time is the most relevant HDD performance indicator. The access time includes the *seek time*, the time for the arm to reach to the cylinder/track, and the *search time*, the time to locate the record on a track. HDD technology has improved dramatically since the disk first introduced by IBM in 1956, as shown in Table 7.1.

7.2 Solid-state disks

Solid-state disks are solid-state persistent storage devices using integrated circuit assemblies as memory. SSD interfaces are compatible with the block I/O of HDDs, thus they can replace the traditional disks. SSDs do not have moving parts, are typically more resistant to physical shock, run silently, and have shorter access time and lower latency than HDDs.

Lower-priced SSDs use triple-level or multi-level cell (MLC) flash memory, slower and less reliable than single-level cell (SLC) flash memory. MLC to SLC ratios of persistence, sequential write, sequen-

Table 7.1 Evolution of hard disk drive technologies from 1956 to 2016. The improvement ranges from astounding ratios such as 650×10^6 to one, 300×10^6 to one, and 2.7×10^6 to one for density, price, and capacity, respectively, to a modest 200 to one for average access time and 11 to one for MTBF, the mean time between failures.

Parameter	1956	2016
Capacity	3.75 MB	10 TB
Average access time	≈ 600 msec	2.5–10 ms
Density	200 bits/sq. inch	1.3 TB sq. inch
Average life span	$\approx 2\,000$ hours/MTBF	$\approx 22\,500$ hours/MTBF
Price	\$9 200/MB	\$0.032/GB
Weight	910 Kg	62 g
Physical volume	1.9 m^3	34 cm^3

tial read, and price are 1 : 10, 1 : 3, 1 : 1, and 1 : 1.3, respectively. Most SSDs use MLC NAND-based flash memory, a nonvolatile memory that retains data when power is lost.

SLC NAND I/O operation latency is: 25 μ sec to fetch a 4 KB page from the array to the I/O buffer on a read, 250 μ sec to commit a 4 KB page from the I/O buffer to the array on a write, and 2 msec to erase a 256 KB block. When multiple NAND devices operate in parallel, SSD bandwidth scales up and high latencies can be hidden, provided that the load is evenly distributed between NAND devices and sufficient outstanding operations are pending. Most SSD manufacturers use nonvolatile NAND due to lower cost compared with DRAM and to the ability to retain data without a constant power supply.

Solid-state hybrid disks (SSHDs) combine the features of SSDs and HDDs; they include a large HDD and an SSD cache to improve performance of frequently accessed data. The SSD cost-per-byte is reduced by about 50% every year, but it must be reduced by up to three orders of magnitude to be competitive with HDD.

Table 7.2 shows that SSD capacity has increased by a factor of five million from 1991 to 2018, and the price per GB has decreased by a factor of 555 000, while the access times has decreased 11 times for read and 38 times for write.

Enterprise flash drives (EFDs) are SSDs with a higher set of specifications, designed for applications requiring high I/O performance (IOPS), reliability, energy efficiency, and consistent performance.

An EFD includes a controller, an embedded processor critical for EFD performance that executes firmware-level code. Some of the functions performed by the controller are: (i) bad block mapping; (ii) read and write caching; (iii) encryption; (iv) crypto-shredding; (v) error detection and correction; (vi) garbage collection; (vi) read scrubbing and read disturb management; and (vii) wear leveling.

The performance scales with the number of parallel NAND flash chips. A single NAND chip is relatively slow, due to the narrow (8/16 bit) asynchronous I/O interface, and additional high latency of basic I/O operations. Multiple NAND devices operate in parallel and the bandwidth scales. High latencies can be hidden, as long as enough outstanding operations are pending and the load is evenly distributed between devices.

NVM Express (NVMe) is an optimized, high-performance scalable host controller interface for non-volatile memory (NVM) technologies; see Fig. 7.1. Typical SAS (Serial Attached SCSI) devices

Table 7.2 SSD Evolution. Leftmost column—SSD parameters, capacity, bandwidth, latency, and cost. Middle column—characteristics, high price, pick performer of enterprise SSDs. Right column—the same parameters for consumer SSDs.

Parameter	Enterprise	Available for Consumers
Capacity	100 TB in 2018 DC100 (20 MB in 1991)	8 TB in 2020
Sequential read speed	15 GB/sec in 2019 (49.3 MB/sec in 2007)	6.795 GB/sec in 2020
Sequential write speed	15.2 GB/sec in 2019 (80.0 MB/sec in 2007)	4.397 GB/sec in 2020
IOPS (I/O operations)	2 500 000 in 2019 (79 in 2007)	736 000 read & 700 000 write in 2020
Access time	45 nsec read & 13 nsec write in 2020 (500 nsec in 2007)	45 nsec read & 13 nsec write in 2020
Price	\$ 0.10/GB in 2020 (\$ 50 000/GB in 1991)	\$ 0.10/GB in 2020

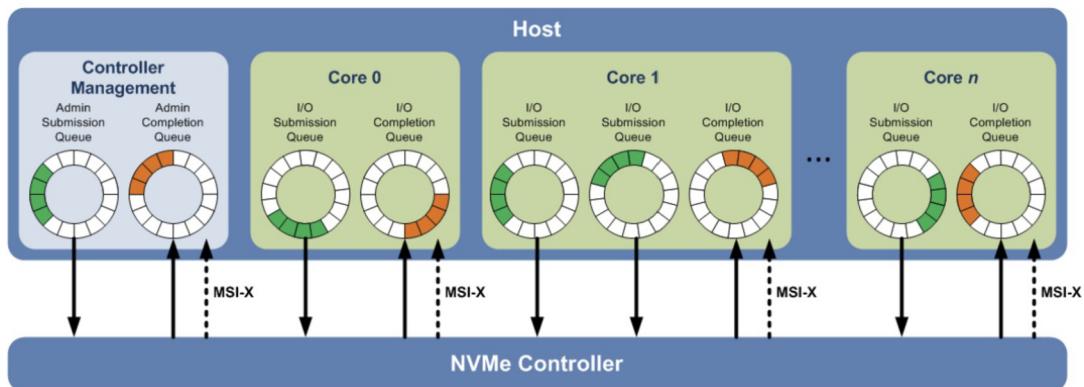
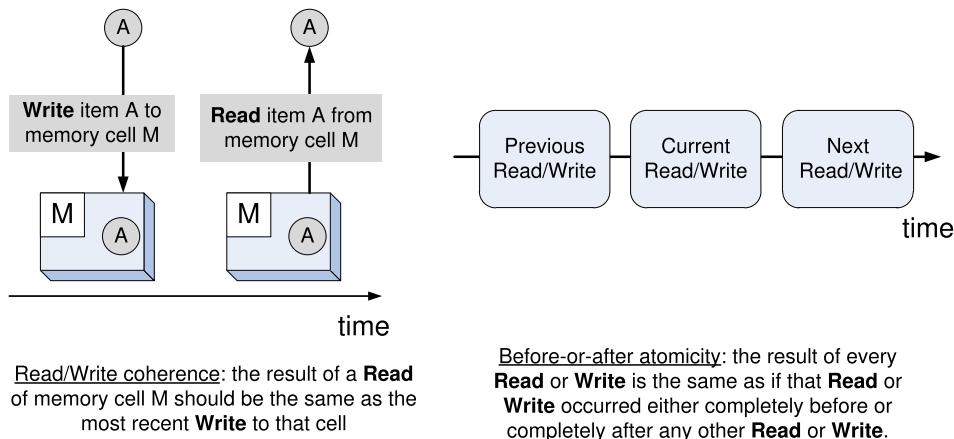


FIGURE 7.1

NVMe controller allows each core to have multiple queues for outstanding operations and to hide the latency of individual operations.

support up to 256 commands; and SATA (Serial AT Attachment) devices support up to 32 commands in a single queue.

NVMe supports 64K commands per queue and up to 64K queues. These queues are designed such that I/O commands and responses to those commands operate on the same processor core and can take advantage of the parallel processing capabilities of multi-core processors. Each application or thread can have its own independent queue, so no I/O locking is required (https://nvmeexpress.org/wp-content/uploads/NVMe_Overview.pdf).

**FIGURE 7.2**

Semantics of read/write coherence and before-or-after atomicity.

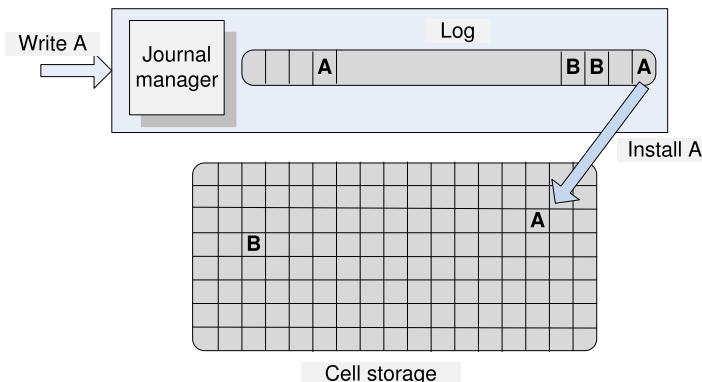
The standard form factor for an SSD is 2.5-inch, which fits inside the drive bay of most laptop or desktop computers. M.2 is a form factor specification for internally mounted SSDs, more expensive per GB than other devices. Samsung 970 EVO Plus M2 NVMe SS has impressive performance: 3 430 GB/sec and 2 813 GB/sec for sequential read and write, respectively, and 278 077 and 190 185 for random IOPS read and write, respectively, according to <https://www.samsung.com/semiconductor/minisite/ssd>.

7.3 Storage models, file systems, and databases

Storage models describe the layout of a data structure in a physical storage; a data model captures the most important logical aspects of a data structure in a database. Physical storage can be a local disk, a removable media, or storage accessible via the network.

Two abstract models of storage are commonly used: cell storage and journal storage. *Cell storage* assumes that the storage consists of cells of the same size and that each object fits in one cell. This model reflects the physical organization of several storage media; the primary memory of a computer is organized as an array of memory cells and a secondary storage device, e.g., a disk, is organized in sectors or blocks read and written as a unit. Read/write coherence and before-or-after atomicity are two highly desirable properties of any storage model and, in particular, of cell storage; see Fig. 7.2.

Journal storage is a fairly elaborate organization for storing composite objects such as records consisting of multiple fields. Journal storage consists of a manager and a cell storage where the entire history of a variable is maintained, rather than just the current value. The user does not have direct access to the cell storage, instead it can request the journal manager to: (i) start a new action; (ii) read the value of a cell; (iii) write the value of a cell; (iv) commit an action; and (v) abort an action. The

**FIGURE 7.3**

A log contains the entire history of all variables; the log is stored on a nonvolatile media of a journal storage. If the system fails after the new value of a variable is stored in the log, but before the value is stored in the cell memory, the value can be recovered from the log. If the system fails while writing the log, the cell memory is not updated. This guarantees that all actions are all-or-nothing. Two variables **A** and **B** in the log and the cell storage are shown. A new value of **A** is written first to the log and then installed on cell memory at the unique address assigned to **A**.

journal manager translates user requests into commands sent to the cell storage: (i) read a cell; (ii) write a cell; (iii) allocate a cell; and (iv) deallocate a cell.

The log of a storage system contains a history of all variables in a cell storage. Information about the updates of each data item forms a record appended to the end of the log. A log provides authoritative information about the outcome of an action involving the cell storage; the cell storage can be reconstructed using the log which can be easily accessed, we only need a pointer to the last record.

An all-or-nothing action first records the action in a log in journal storage and then installs the change in the cell storage by overwriting the previous version of a data item; see Fig. 7.3. The log is always kept on non-volatile storage, e.g., disk, and the considerably larger cell storage resides typically on nonvolatile memory, but can be held in memory for real-time access or using a write-through cache.

A file system consists of a collection of directories and each directory provides information about a set of files. High-performance systems can choose among three classes of file systems: Network File Systems (NFS), Storage Area Networks (SAN), and Parallel File Systems (PFS). Network file systems are very popular and have been used for some time, but do not scale well and have reliability problems; an NFS server could be a single point of failure.

Advances in networking technology allow separation of storage systems from computational servers; the two can be connected by an SAN. SANs offer additional flexibility and allow cloud servers to deal with nondisruptive changes in the storage configuration. Moreover, the storage in an SAN can be *pooled* and then allocated based on the needs of the servers; pooling requires additional software and hardware support and represents another advantage of a centralized storage system. An SAN-based implementation of a file system can be expensive because each node must have a Fibre Channel adapter to connect to the network.

Parallel file systems are scalable, are capable of distributing files across a large number of nodes, and provide a global naming space. In a parallel data system, several I/O nodes serve data to all computational nodes; the system includes also a metadata server that contains information about the data stored in the I/O nodes. The interconnection network of a parallel file system could be a SAN.

Databases and database management systems. Most cloud applications do not interact directly with the file systems but through an application layer that manages a database. A database is a collection of logically related records. The software that controls the access to the database is called a Data Base Management System (DBMS). The main functions of a DBMS are: enforce data integrity, manage data access and concurrency control, and support recovery after a failure.

A DBMS supports a query language, a dedicated programming language used to develop database applications. Several database models, including the navigational model of the 1960s, the relational model of the 1970s, the object-oriented model of the 1980s, and the NoSQL model of the first decade of the 2000s, reflect the limitations of the hardware available at the time and the requirements of the most popular applications of each period.

Cloud databases. Most cloud applications are data-intensive, test the limitations of existing cloud storage infrastructure, and demand database management systems capable of supporting rapid application development and a short time-to-market. Cloud applications require low latency, scalability, high availability, and demand a consistent view of data. These requirements cannot be satisfied simultaneously by existing database models; for example, relational databases are easy to use for application development but do not scale well.

As its name implies, the NoSQL model does not support SQL as a query language and may not guarantee the ACID, Atomicity, Consistency, Isolation, and Durability properties of traditional databases. It usually guarantees an eventual consistency for transactions limited to a single data item. The NoSQL model is useful when the structure of the data does not require a relational model and the amount of data is very large. Several types of NoSQL databases have emerged in the last few years. Based on the manner the NoSQL databases store the data, we recognize several types such as key-value stores, BigTable implementations, document store databases, and graph databases.

Replication, used to ensure fault tolerance of large-scale systems built with commodity components, requires mechanisms to guarantee that replicas are consistent with one another. This is yet another example of increased complexity of modern computing and communication systems when the software has to support desirable properties of the physical systems. Section 7.7 contains an in-depth analysis of a service implementing a consensus algorithm to guarantee that replicated objects are consistent.

Many cloud applications support online transaction processing and have to guarantee the correctness of the transactions. Transactions consist of multiple actions; for example, the transfer of funds from one account to another requires withdrawing funds from one account and crediting it to another. The system may fail during or after each one of the actions, and steps to ensure correctness must be taken. Correctness of a transaction means that the result should be guaranteed to be the same as if the actions were applied one after another regardless of the order. More stringent conditions must sometimes be observed; for example, banking transactions must be processed in the order they are issued, the so-called *external time consistency*. To guarantee correctness, a transaction processing system supports *all-or-nothing atomicity*, discussed in Section 10.12.

7.4 Distributed file systems; the precursors

The first distributed file systems were developed in the 1980s by software companies and universities. The systems covered are: NFS, developed by Sun Microsystems in 1984; AFS, developed at Carnegie Mellon University as part of the Andrew project; and SFS, developed by John Osterhout's group at U.C. Berkeley as a component of the Unix-like distributed operating system called Sprite. Other systems developed at about the same time are Locus [500], Apollo [299], and Remote File System (RFS) [41]. Main concerns in the design of these systems are scalability, performance, and security; see Table 7.3.

In 1980s, many organizations, including research centers, universities, financial institutions, and design centers, considered that networks of workstations are an ideal environment for their operations. Diskless workstations were appealing due to reduced hardware costs and lower maintenance and system administration costs. Soon, it became obvious that a distributed file system could be very useful for the management of a large number of workstations, and Sun Microsystems, one of the main promoters of a distributed systems based on workstations, proceeded to develop the NFS in early 1980s.

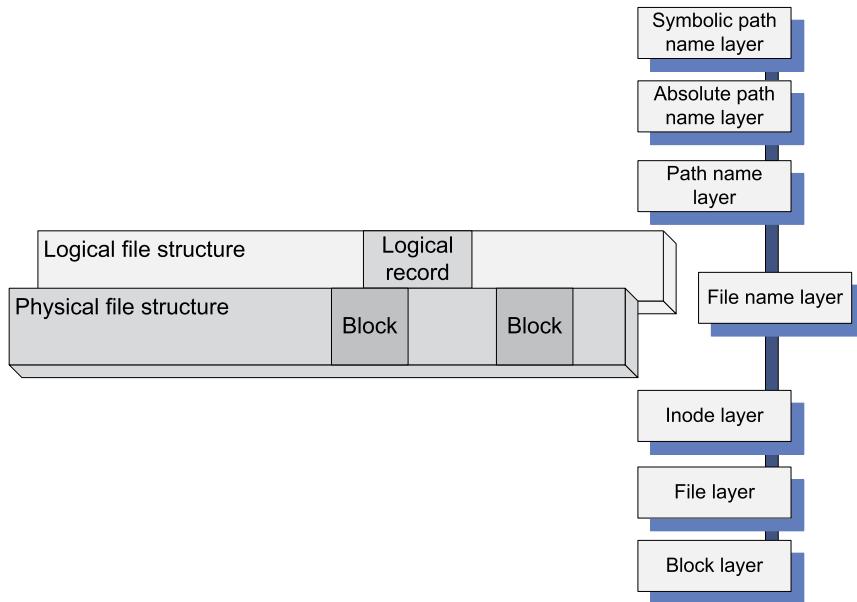
Network File System. NFS was the first widely used distributed file system; the development of this application based on the client-server model was motivated by the need to share a file system among a number of clients interconnected by a local area network.

A majority of workstations were running under UNIX; thus, many design decisions for the NFS were influenced by the design philosophy of the UNIX File System (UFS). It is not surprising that the NFS designers aimed to: (i) provide the same semantics as a local UFS to ensure compatibility with existing applications and facilitate easy integration into existing UFS; (ii) ensure that the system will be widely used, thus support clients running on different operating systems; and (iii) accept a modest performance degradation due to remote access over a network with a bandwidth of several Mbps.

UFS has three important characteristics enabling extension from local to remote file management:

1. Layered design provides the necessary flexibility of the file system. Layering allows separation of concerns and minimization of the interaction among the modules necessary to implement the system. The addition of the vnode layer allowed UNIX file system to treat uniformly local and remote file access.
2. Hierarchical design supports file system scalability; it allows grouping of files into special files called directories, supports multiple levels of directories and collections of directories and files, the so-called file systems. The hierarchical file structure is reflected by the file-naming convention.
3. Metadata supports a systematic rather than an ad hoc design philosophy of the file system. The inodes contain information about individual files and directories and are kept on persistent media together with the data. Metadata includes the file owner, the access rights, the creation time or the time of the last modification of the file, and the file size, as well as information about the structure of the file and the persistent storage device cells where the data is stored. Metadata also supports device independence, a very important objective due to the very rapid pace of storage technology development.

The *logical organization* of a file reflects the data model, the view of the data from the perspective of the application. The *physical organization* reflects the storage model and describes the manner the file is stored on a given storage media. The layered design allows UFS to separate the concerns for the physical file structure from those for the logical one.

**FIGURE 7.4**

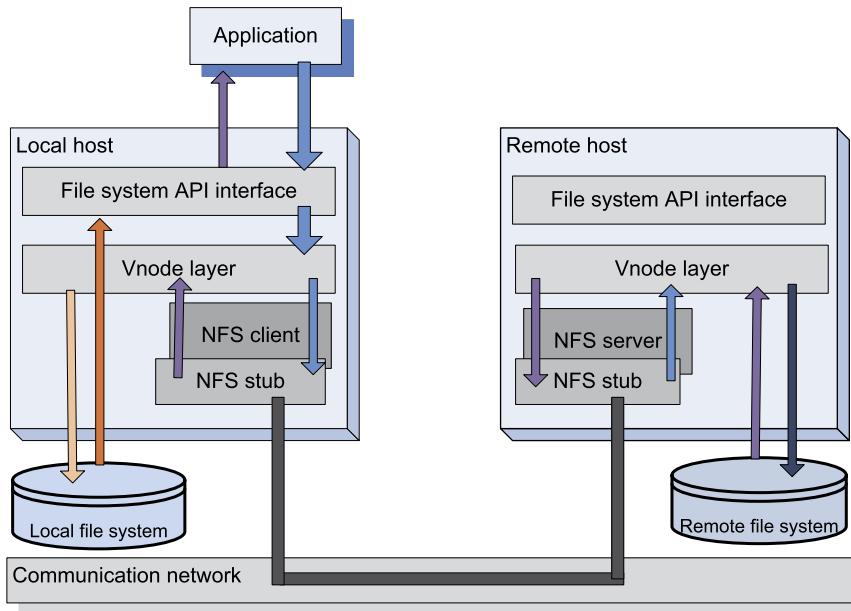
UFS layered design separates the physical file structure from the logical one. The lower three layers, block, file, and inode, are related to the physical file structure, while the upper three layers, path name, absolute path name, and symbolic path name, reflect the logical organization. The file name layer mediates between the two groups.

Recall that a file is a linear array of cells stored on a persistent storage device; the file pointer identifies a cell used as a starting point for a read or write operation. This linear array is viewed by an application as a collection of logical records; the file is stored on a physical device as a set of physical records, or blocks, of size dictated by the physical media.

The lower three layers of the UFS hierarchy, the block, the file, and the inode layer, reflect the physical organization. The block layer allows the system to locate individual blocks on the physical device; the file layer reflects the organization of blocks into files; and the inode layer provides the metadata for the objects (files and directories). The upper three layers, the path name, the absolute path name, and symbolic path name layer, reflect the logical organization. The file name layer mediates between the machine-oriented and the user-oriented views of the file system; see Fig. 7.4.

Several control structures maintained by the kernel of the operating systems support the file handling by a running process; these structures are maintained in the user area of the process address space and can only be accessed in kernel mode. To access a file, a process must first establish a connection with the file system by opening the file; at that time, a new entry is added to the file description table, and the metainformation is brought in to another control structure, the open file table.

A *path* specifies the location in a file system of a file or directory; a *relative path* specifies this location relative to the current/working directory of the process, while a *full path*, also called an absolute

**FIGURE 7.5**

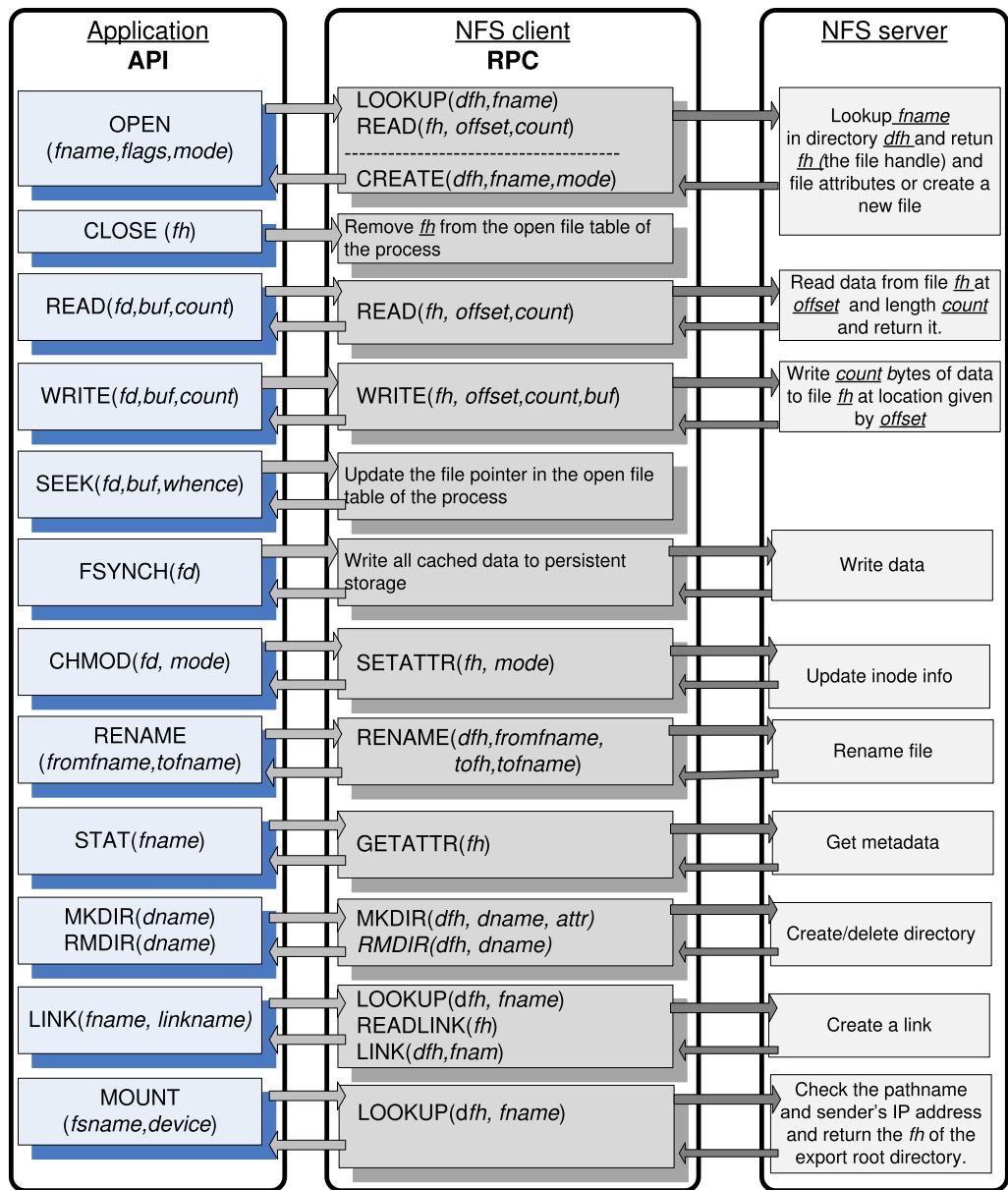
The NFS client-server interaction. The vnode layer implements file operation in a uniform manner, regardless of whether the file is local or remote. An operation targeting a local file is directed to the local file system, while one for a remote file involves NFS; an NFS client packages the relevant information about the target and the NFS server passes it to the vnode layer on the remote host, which, in turn, directs it to the remote file system.

path, specifies file location independently of the current directory, typically relative to the root directory. A local file is uniquely identified by a *file descriptor* (fd), generally, an index in the open file table.

NFS is based on the client–server paradigm. The client runs on the local host, while the server is at the site of the remote file system, and they interact by means of Remote Procedure Calls (RPCs). NFS uses a vnode layer to distinguish between operations on local and remote files; see Fig. 7.5. The API interface of the local file system distinguishes file operations on a local file from the ones on a remote file and, in the latter case, invokes the RPC client. Fig. 7.6 shows the API for a UNIX file system and the calls made by the RPC client in response to API calls issued by a user program for a remote file system, as well as some of the actions carried out by the NFS server in response to an RPC call.

A remote file is uniquely identified by a file handle rather than a file descriptor; a file handle is a 32-byte internal name, a combination of a file system identification, an inode number, and a generation number. A file handle allows a host to locate the remote file system and the file on that system. A generation number allows the system to reuse inode numbers and ensures a correct semantics when multiple clients operate on the same remote file.

While many RPC calls, such as *Read*, are idempotent, an action is idempotent if repeating it several times has the same effect as if the action was executed only once. Communication failures could sometimes lead to an unexpected behavior. Indeed, if the network fails to deliver the response to a *Read* RPC,

**FIGURE 7.6**

The API of UFS and the corresponding RPCs issued by an NFS client to the NFS server. The actions of the server in response to an RPC issued by the NFS client are too complex to be fully described here. *fd* stands for file descriptor, *fh* for file handle, *fname* for file name, *dname* for directory name, *dfh* for the directory where the file handle can be found, *count* for the number of bytes to be transferred, *buf* for the buffer to transfer the data to/from, and *device* for the device where the file system is located.

then the call can be repeated without any side effects. By contrast, when the network fails to deliver the response to the *Rmdir* RPC, the second call returns an error code to the user if the call was successful the first time; if the server fails to execute the first call, then the second call returns normally. Note also that there is no *Close* RPC because this action only makes changes to the open file data structure of the process and does not affect the remote file.

NFS has undergone significant transformations over the years; it has evolved from Version 2 [433], discussed in this section, to Version 3 [392] in 1994, and then to Version 4 [393] in 2000; see Section 7.16.

Andrew File System. AFS is a distributed file system developed in the late 1980s at Carnegie Mellon University in collaboration with IBM [358]. System designers envisioned a very large number of workstations interconnected with a relatively small number of servers. It was anticipated that each individual at CMU would have an Andrew workstation, thus the system would connect up to 10 000 workstations.

The set of trusted servers in AFS form a structure called Vice. The workstation OS, 4.2BSD UNIX, intercepts file system calls and forwards them to a user-level process called Venus which caches files from Vice and stores modified copies of files back on the servers they came from. Reading and writing operations are performed directly on the cached copy of the file and bypass Venus. Only when a file is opened or closed does Venus communicate with Vice.

The emphasis of the AFS design is on performance, security, and simple management of the file system [245]. The local disk of a workstation acts as a persistent cache ensuring scalability and reducing the response time. The master copy of a file residing on one of the servers is updated only when the file is modified. This strategy reduces the server load and improves the system performance.

Another major objective of the AFS design is improved security. The communications between clients and servers are encrypted, and all file operations require secure network connections. When a user signs in to a workstation, the password is used to obtain security tokens from an authentication server; these tokens are then used every time a file operation requires a secure network connection.

The AFS uses Access Control Lists (ACLs) to allow control sharing of the data. An ACL specifies the access rights of an individual user or of a group of users. A set of tools support the management of ACLs. Another facet of the effort to reduce the user involvement in the file management is *location transparency*. Files could be accessed from any location and could be moved automatically, or at the request of system administrators, without user's involvement and inconvenience. The relatively small number of servers reduces drastically the efforts related to system administration because operations, such as backups, affect only the servers, while workstations can be added, removed, or moved from one location to another without administrative intervention.

Sprite Network File System. SFS is a component of the Sprite network operating system [237]. SFS supports non-write-through caching of files on the client, as well as the server systems [363]. Processes running on all workstations enjoy the same semantics for file access as if they would run on a single system; this is possible due to a cache consistency mechanism that flushes portions of the cache and disables caching for shared files opened for read-write operations.

Caching not only hides the network latency but also reduces the server utilization and improves the performance by reducing the responder time. A file access request made by a client process could be satisfied at various levels. First, the request is directed to the local cache; if not satisfied there, it is passed to the local file system of the client. If it cannot be satisfied locally, then the request is sent to the remote server; if the request cannot be satisfied by the remote server's cache, then it is sent to the file system running on the server.

Sprite system design decisions were influenced by resources available when a typical workstation had a 1 to 2 MIPS processor and 4 to 14 Mbytes of physical memory. The main-memory caches allowed diskless workstations to be integrated in the system and enabled the development of unique caching mechanisms and policies for both clients and servers. The results of a file-intensive benchmark reported by [363] show that SFS was 30 to 35% faster than either NFS or AFS.

The file cache is organized as a collection of 4K blocks; a cache block has a virtual address consisting of a unique file identifier supplied by the server and a block number in the file. Virtual addressing allows the clients to create new blocks without the need to communicate with the server; file servers map virtual addresses to physical disk addresses. Note also that the page size of the virtual memory in Sprite is also 4K. The size of the cache available to an SFS client or a server system changes dynamically function of the needs. This is possible because the Sprite operating system ensures an optimal sharing of the physical memory between file caching by SFS and virtual memory management.

The file system and the virtual memory manage separate sets of physical memory pages and maintain a time-of-last-access for each block or page, respectively. Virtual memory uses a version of the clock algorithm [362] to implement a Least Recently Used (LRU) page replacement algorithm, and the file system implements a strict LRU order since it knows the time of each read and write operation. Whenever the file system or the virtual memory management experiences a file cache miss or a page fault, it compares the age of its oldest cache block or page, respectively, with the age of the oldest one of the other system; the oldest cache block or page is forced to release the real memory frame.

An important design decision of SFS was to delay write-backs; this means that a block is first written to cache, and the writing to the disk is delayed for a time on the order of tens of seconds. This strategy speeds up writing and also avoids writing when the data is discarded before the time to write it to the disk. The obvious drawback of this policy is that data can be lost in the case of a system failure. A write-through is the alternative to the delayed write-back; it guarantees reliability because the block is written to the disk as soon as it is available on the cache, but it increases the time for a write operation.

Most network file systems guarantee that, once a file is closed, the server will have the newest version on persistent storage. As far as concurrency is concerned, we distinguish sequential write-sharing, when a file cannot be opened simultaneously for reading and writing by several clients, from concurrent write-sharing, when multiple clients can modify the file at the same time. Sprite allows both concurrency modes and delegates cache consistency to the servers. In case of concurrent write-sharing, the client cashing for the file is disabled; all reads and writes are carried out through the server.

Table 7.3 presents a comparison of caching, writing strategy, and consistency of NFS, AFS [358], Sprite [237], Locus [500], Apollo [299], and RFS [41].

7.5 General parallel file system

Once the distributed file systems became ubiquitous, the natural next step in the file systems evolution was supporting parallel access. Parallel file systems allow multiple clients to read and write concurrently from the same file. Support for parallel I/O is essential for the performance of many applications [339]. Early supercomputers such as Intel Paragon took advantage of parallel file systems to support data-intensive applications.

Concurrency control is a critical issue for parallel file systems. Several semantics for handling shared and concurrent file access are possible. One option is to have a shared file pointer; in this case,

Table 7.3 A comparison of several network file systems [358].

File system	Cache size and location	Writing policy	Consistency guarantees	Cache validation
NFS	Fixed, memory	On close or 30 sec. delay	Sequential	On open, with server consent
AFS	Fixed, disk	On close	Sequential	When modified server asks client
SFS	Variable, memory	30 sec. delay	Sequential, concurrent	On open, with server consent
Locus	Fixed, memory	On close	Sequential, concurrent	On open, with server consent
Apollo	Variable, memory	Delayed or on unlock	Sequential	On open, with server consent
RFS	Fixed, memory	Write-through	Sequential, concurrent	On open, with server consent

successive reads issued by different clients advance the file pointer. Another semantics is to allow each client to have its own file pointer.

IBM developed the General Parallel File System [440] in early 2000s as a successor of TigerShark multimedia file system [230]. GPFS emulates closely the behavior of a general-purpose POSIX system running on a single system. This parallel file system was designed for optimal performance of large clusters and can support a file system of up to 4 petabytes consisting of up to 4 096 disks of 1 TB each.

The maximum file size is $(2^{63} - 1)$ bytes. A file consists of blocks of equal size, ranging from 16 KB to 1 MB, striped across several disks. The system could support not only very large files but also a very large number of files. GPFS directories use the *extensible hashing* techniques to access a file.

A hash function is applied to the file name; then, the n low-order bits of the hash value give the block number of the directory where the file information can be found with n , a function of the number of files in the directory. Extensible hashing is used to add a new directory block. The system maintains user data, file metadata, such as last modified time, and file system metadata, such as allocation maps. Metadata, such as file attributes and data block addresses, is stored in inodes and in indirect blocks.

Reliability is a major concern in a system with many physical components. To recover from system failures, GPFS records all metadata updates in a write-ahead log file. *Write-ahead* means that updates are written to persistent storage only after the log records have been written. For example, when a new file is created, a directory block must be updated, and an inode for the file must be created. These records are transferred from cache to disk after the log records have been written. When the system ends up in an inconsistent state, the directory block is written, and then, if the I/O node fails before writing the inode, the log file allows to recreate the inode record.

The log files are maintained by each I/O node for each file system it mounts; thus, any I/O node is able to initiate recovery on behalf of a failed node. Disk parallelism is used to reduce the access time; multiple I/O read requests are issued in parallel, and data is pre-fetched in a buffer pool.

Data striping allows concurrent access and improves performance but can have unpleasant side effects. Indeed, when a single disk fails, a large number of files are affected. To reduce the impact of such undesirable events, the system attempts to mask a single disk failure or the failure of the access path to a disk. The system uses RAID devices with the stripes equal to the block size and dual-attached

RAID controllers. To further improve the fault tolerance of the system, GPFS data files and metadata are replicated on two different physical disks.

Consistency and performance, critical for any distributed file system, are difficult to balance; support for concurrent access improves the performance but faces serious challenges for maintaining consistency. GPFS consistency and synchronization are ensured by a distributed locking mechanism; a *central lock manager* grants *lock tokens* to *local lock managers* running in each I/O node. Lock tokens are also used by the cache management system.

Lock granularity has important implications for the performance of a file system, and GPFS uses a variety of techniques for different types of data. *Byte-range tokens* are used for read and write operations to data files as follows: The first node attempting to write to a file acquires a token covering the entire file, $[0, \infty]$. This node is allowed to carry out all reads and writes to the file without any need for permission until a second node attempts to write to the same file; then, the range of the token given to the first node is restricted. More precisely, if the first node writes sequentially at offset fp_1 and the second one at offset $fp_2 > fp_1$, then the ranges of the two tokens are $[0, fp_2]$ and $[fp_2, \infty]$, respectively, and the two nodes can operate concurrently without the need for further negotiations. Byte-range tokens are rounded to block boundaries.

Byte-range token negotiations among nodes use the *required range* and the *desired range* for the offset and for the length of the current and the future operations, respectively. The *data-shipping*, an alternative to byte-range locking, allows fine-grain data sharing. In this mode, the file blocks are controlled by the I/O nodes in a round-robin manner. A node forwards a read or write operation to the node controlling the target block, the only one allowed to access the file.

A *token manager* maintains the state of all tokens; it creates and distributes tokens, collects tokens once a file is closed, and downgrades/upgrades tokens when additional nodes request access to a file. Token management protocols attempt to reduce token manager load; for example, when a node wants to revoke a token it sends messages to all other nodes holding the token and forwards the reply to the token manager.

Access to metadata is synchronized; for example, when multiple nodes write to the same file, the file size and the modification dates are updated using a *shared write lock* to access an inode. One of the nodes assumes the role of a *metanode* and all updates are channeled through it, the file size and the last update time are determined by the metanode after merging the individual requests. The same strategy is used for updates of the indirect blocks. GPFS global data, such as ACLs (Access Control Lists), quotas, and configuration data, are updated using the distributed locking mechanism.

GPFS uses *disk maps* for the management of the disk space. The GPFS block size can be as large as 1 MB, and a typical block size is 256 KB. A block is divided into 32 sub-blocks to reduce disk fragmentation for small files, thus the block map has 32 bits to indicate if a subblock is free or used. The system disk map is partitioned into n regions, and each disk map region is stored on a different I/O node; this strategy reduces the conflicts and allows multiple nodes to allocate disk space at the same time. An *allocation manager* running on one of the I/O nodes is responsible for actions involving multiple disk map regions. For example, it updates free space statistics and helps with deallocation by sending periodically hints of the regions used by individual nodes.

A detailed discussion of system utilities and of the lessons learned from the deployment of the file system at several installations in 2002 can be found in [440]; the documentation of the GPFS is available from [250].

7.6 Google file system

Google File System, developed in late 1990s, uses thousands of storage systems built from inexpensive commodity components to provide petabytes of storage to a large user community with diverse needs [197]. Thus, it should not be surprising that a main concern of GFS designers was reliability of a system exposed to hardware failures, system software errors, application errors, and, last but not least, human errors. The system was designed after a careful analysis of file characteristics and access models. Some of the most important aspects of this analysis reflected in the GFS design are:

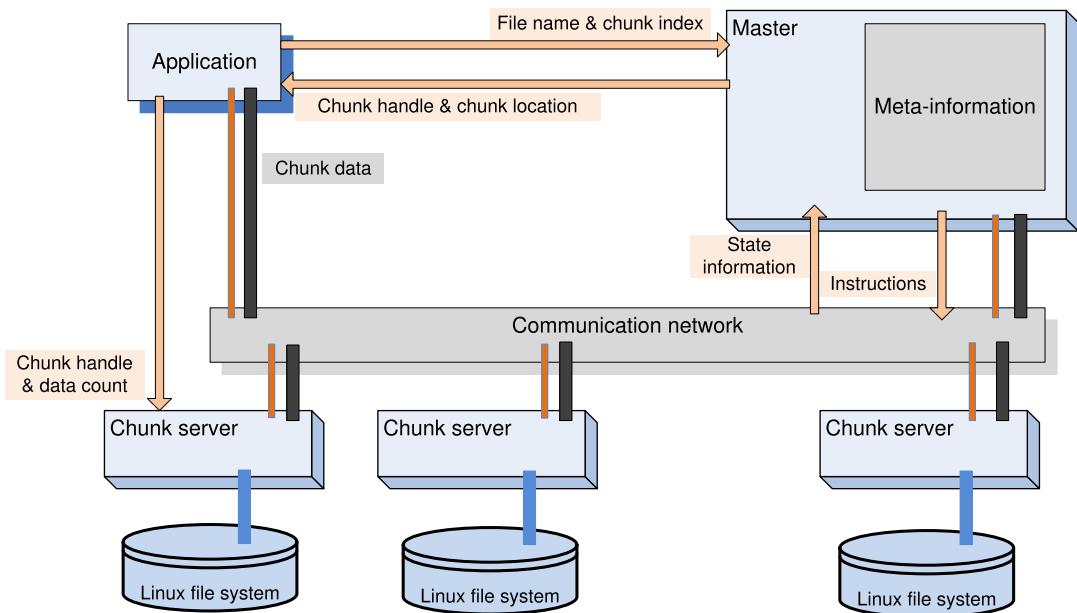
- a. Scalability and reliability are critical features of the system; they must be considered from the beginning, rather than at a later design stage.
- b. The vast majority of files range in size from a few GB to hundreds of TB.
- c. The most common operation is to append to an existing file; random write operations to a file are extremely infrequent.
- d. Sequential read operations are the norm.
- e. Users process the data in bulk and are less concerned with the response time.
- f. The consistency model should be relaxed without placing an additional burden on the application developers to simplify system implementation.

This analysis led to several major design decisions:

1. Segment a file in large chunks.
2. Implement an atomic file *append* operation allowing multiple applications operating concurrently to append to the same file.
3. Build the cluster around a high-bandwidth rather than low-latency interconnection network. Separate the flow of control from the data flow; schedule the high-bandwidth data flow by pipelining the data transfer over TCP connections to reduce the response time. Exploit network topology by sending data to the closest node in the network.
4. Eliminate client site caching; caching increases the overhead for maintaining consistency among cached copies at multiple client sites and is not likely to improve performance.
5. Ensure consistency by channeling critical file operations through a master controlling the entire system.
6. Minimize master's involvement in file access operations to avoid hotspot contention and to ensure scalability.
7. Support efficient checkpointing and fast recovery mechanisms.
8. Support efficient garbage-collection mechanisms.

GFS files are collections of fixed-size segments called *chunks*; at the time of file creation, each chunk is assigned a unique *chunk handle*. A chunk consists of 64 KB blocks and each block has a 32 bit checksum. Chunks are stored on Linux files systems and are replicated on multiple sites; a user may change the number of the replicas, from the standard value of three to any desired value. The chunk size is 64 MB; this choice is motivated by the desire to optimize the performance for large files and to reduce the amount of metadata maintained by the system.

A large chunk size increases the likelihood that multiple operations will be directed to the same chunk, thus it reduces the number of requests to locate the chunk, and, at the same time, it allows an application to maintain a persistent network connection with the server where the chunk is located.

**FIGURE 7.7**

GFS cluster architecture. The *master* maintains state information about all system components and controls a number of *chunk servers*. A chunk server runs under Linux and uses metadata provided by the master to communicate directly with an application. The data flow is decoupled from the control flow. The data and the control paths are shown separately, data paths with thick lines and the control paths with thin lines. Arrows show the flow of control between an application, the master, and the chunk servers.

Space fragmentation occurs infrequently because the chunk of a small file and the last chunk of a large file are only partially filled.

The architecture of a GFS cluster is illustrated in Fig. 7.7. The *master* controls a large number of *chunk servers*; it maintains metadata such as the file names, access control information, the location of all the replicas for every chunk of each file, and the state of individual chunk servers. Some of the metadata is stored in persistent storage, e.g., the *operation log* records the file namespace, as well as the file-to-chunk-mapping. Chunks location is stored only in the control structure of master's memory and is updated at the system start up, or when a new chunk server joins the cluster. This strategy allows the master to have up-to-date information about the location of the chunks.

System reliability is a major concern, and the operation log maintains a historical record of metadata changes enabling the master to recover in case of a failure. As a result, such changes are atomic and are not made visible to the clients until they have been recorded on multiple replicas on persistent storage. To recover from a failure, the master replays the operation log. To minimize the recovery time, the master periodically checkpoints its state and, at recovery time, it replays only the log records after the last checkpoint.

Each chunk server is a commodity Linux system. A chunk server receives instructions from the master and responds with status information. For file read or write operations an application sends to the master the file name, the chunk index, and the offset in the file. The master responds with the chunk handle and the location of the chunk. Then the application *communicates directly* with the chunk server to carry out the desired file operation.

The consistency model is very effective and scalable. Operations, such as file creation, are atomic and are handled by the master. To ensure scalability, the master has a minimal involvement in file mutations, operations such as *write* or *append*, that occur frequently. In such cases the master grants a lease for a particular chunk to one of the chunk servers called the *primary*; then, the primary creates a serial order for the updates of that chunk.

When a data for a *write* straddles the chunk boundary, two operations are carried out, one for each chunk. The steps for a *write* request illustrate a process that buffers data and decouples the control flow from the data flow for efficiency:

- i. The client contacts the master, which assigns a lease to one of the chunk servers for the a particular chunk, if no lease for that chunk exists; then, the master replies with the id of the primary as well as secondary chunk servers holding replicas of the chunk. The client caches this information.
- ii. The client sends the data to all chunk servers holding replicas of the chunk; each one of the chunk servers stores the data in an internal LRU buffer and then sends an acknowledgment to the client.
- iii. The client sends the *write* request to the primary chunk server once it has received the acknowledgments from all chunk servers holding replicas of the chunk. The primary chunk server identifies mutations by consecutive sequence numbers.
- iv. The primary chunk server sends the *write* requests to all secondaries.
- v. Each secondary chunk server applies the mutations in the order of the sequence number and then sends an acknowledgment to the primary one.
- vi. Finally, after receiving the acknowledgments from all secondaries, the primary informs the client.

The system supports an efficient checkpointing procedure based on *copy-on-write* to construct system snapshots. A lazy garbage-collection strategy is used to reclaim the space after a file deletion. As a first step, the file name is changed to a hidden name and this operation is time stamped. The master periodically scans the namespace and removes metadata for files with a hidden name older than a few days. This mechanism enables a user who deleted files by mistake to recover the files with little effort.

Periodically, chunk servers exchange with the master the list of chunks stored on each one of them; the master supplies them with the identity of orphaned chunks, whose metadata has been deleted, and such chunks are then deleted. Even when control messages are lost, a chunk server will carry out the house cleaning at the next *heartbeat* exchange with the master. Each chunk server maintains in core the checksums for the locally stored chunks to guarantee data integrity. *CloudStore* is an open source C++ implementation of GFS. CloudStore allows client access from C++, Java, and Python.

7.7 Locks; Chubby—a locking service

Locks support the implementation of reliable storage for loosely coupled distributed systems. Locks enable controlled access to shared storage and ensure atomicity of *read* and *write* operations. Consensus protocols are critically important for the design of reliable distributed storage systems. The election of

a leader or master from a group of data servers requires an effective consensus protocol because the master plays an important role in the management of a distributed storage system. For example, in GFS, the master maintains state information about all systems components.

Locking and the election of a master can be done using a version of the Paxos algorithm for asynchronous consensus. The algorithm guarantees safety without any timing assumptions, a necessary condition in a large-scale system when communication delays are unpredictable. Nevertheless, the algorithm must use clocks to ensure liveness and to avoid the impossibility of reaching consensus with a single faulty process [175]. Coordination and consensus using Paxos are discussed in depth in Sections 11.4 and 10.13, respectively.

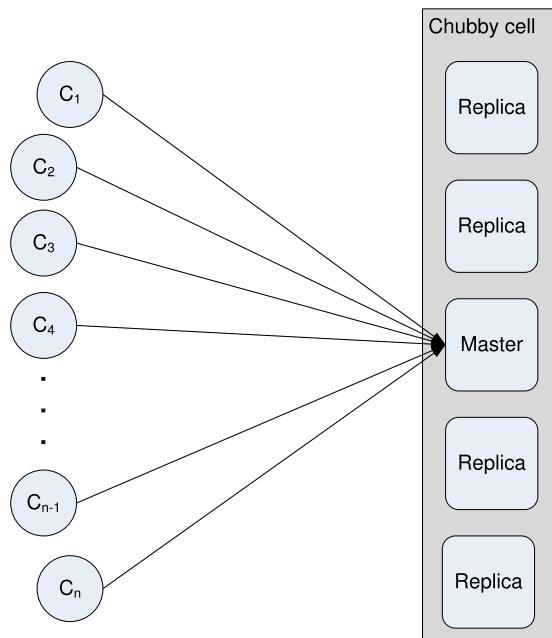
Distributed systems experience communication problems, such as lost messages, messages out of sequence, or corrupted messages. There are solutions for handling these undesirable phenomena; for example, one can use virtual time, i.e., sequence numbers, to ensure that messages are processed in an order consistent with the time they were sent by all processes involved, but this complicates the algorithms and increases the processing time.

Advisory locks are based on the assumption that all processes play by the rules; advisory locks do not have any effect on processes that circumvent the locking mechanisms and access the shared objects directly. *Mandatory locks* block access to the locked objects to all processes that do not hold the locks, regardless if they use locking primitives or not.

Locks held for a very short time are called *fine-grained*, while *coarse-grained* locks are held for a longer time. Some operations such as accessing the name of the lock, whether the lock is shared or held in exclusivity, the generation number of the lock require meta-information about the lock. The metainformation is sometimes aggregated into an opaque byte-string called a *sequencer*. The question how to most effectively support a locking and consensus component of a large-scale distributed system demands several design decisions.

A first decision is whether the locks should be mandatory or advisory. Mandatory locks have the obvious advantage of enforcing access control; a traffic analogy is that a mandatory lock is like a drawbridge: Once it is up, all traffic is forced to stop. An advisory lock is like a stop sign: Those who obey the traffic laws will stop, but some may not. The disadvantages of mandatory locks are added overhead and less flexibility. Once a data item is locked, even a high-priority task related to maintenance or recovery cannot access the data unless it forces the application holding the lock to terminate. This is a very significant problem in large-scale systems where partial system failures are likely.

A second design decision is whether the system should be based on fine-grained or coarse-grained locking. *Fine-grained locks* allow more application threads to access shared data, but generate a larger workload for the lock server. Moreover, when the lock server fails for a period of time, a larger number of applications are affected. Advisory locks and *coarse-grain locks* seem to be a better choice for a system expected to scale to a very large number of nodes distributed in data centers interconnected via wide area networks with a higher communication latency. A third design decision is how to support a systematic approach to locking. Two alternatives come to mind: (i) delegate the implementation of the consensus algorithm to the clients and provide a library of functions needed for this task; and (ii) create a locking service implementing a version of the asynchronous Paxos algorithm and provide a library to be linked with an application client to support service calls. Forcing application developers to invoke calls to a Paxos library is more cumbersome and more prone to errors than the service alternative. Of course, the lock service itself has to be scalable to support a potentially heavy load.

**FIGURE 7.8**

A Chubby cell consisting of five replicas, one of them is elected as the master. Clients c_1, c_2, \dots, c_n communicate with the master using RPCs.

Another consideration when making this choice is flexibility: the ability of the system to support a variety of applications. A name service comes to mind as many cloud applications *read* and *write* small files. To allow atomic file operations, the names of small files should be included in the namespace of the service. The choice should also consider the performance, if a service can be optimized and if clients can be allowed to cache control information. Lastly, the overhead and resources for reaching consensus should be considered. Again, the service alternative seems more advantageous because it needs fewer replicas for high availability.

In 2000 Google decided to use advisory locks and coarse-grained locks and started developing Chubby [81], a lock service. The service has been used by several Google systems including GFS discussed in Section 7.6 and BigTable presented in Section 7.11. A Chubby cell typically serves one data center. The cell server in Fig. 7.8 includes several *replicas*; the standard number of replicas is five. To reduce the probability of correlated failures, the servers hosting replicas are distributed across the campus of a data center.

Chubby replicas use a distributed consensus protocol to elect a new *master* when the current one fails. A master is elected by a majority, as required by the asynchronous Paxos algorithm, accompanied by the commitment that a new master will not be elected for a period of time, namely, the *master lease*. A session is a connection between a client and the cell server maintained over a period of time. Data cached by a client, locks acquired, and handles of all files locked by the client are only valid for the

duration of the session. Clients use RPCs to request services from the master. The master responds without consulting replicas when receiving a *read* request, propagates the request to all replicas, and waits for a reply from a majority of replicas before responding to a *write* request.

The client interface of the system is similar, yet simpler, than the one supported by the Unix file system; in addition, it includes notification for events related to file or system status. A client can subscribe to events such as: file contents modification, change or addition of a child node, master failure, lock acquired, conflicting lock requests, and invalid file handle. Files and directories of Chubby service are organized in a tree structure and use a naming scheme similar to Unix. Each file has a *file handle* similar to the file descriptor.

The master of a cell periodically writes a snapshot of its database to a GFS file server. Every file or directory can act as a lock. To *write* to a file, the client must be the only one holding the file handle, while multiple clients may hold the file handle to *read* from the file. Handles are created by a call to *open()* function and destroyed by a call to *close()*. Other calls supporting the service are *GetContentsAndStat()*, to get the file data and meta-information, and *SetContents*, *Delete*.

Several calls allow the client to acquire and release locks. Some applications may decide to create and manipulate a sequencer with calls to: *SetSequencer()* for associating a sequencer with a handle; *GetSequencer()* for obtaining the sequencer associated with a handle; or *CheckSequencer()* for checking the validity of a sequencer.

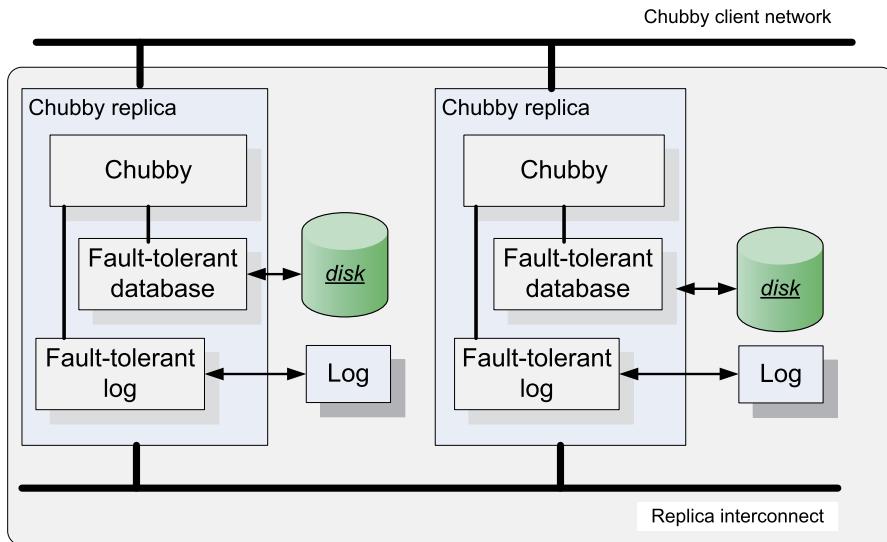
The sequence of calls *SetContents()*, *SetSequencer()*, *GetContentsAndStat()*, and *CheckSequencer()* can be used by an application for the election of a master. In this process all candidate threads attempt to open a lock file, call it *lfile*, in exclusive mode. The one that succeeds to acquire the lock for *lfile*, becomes the master, writes its identity in *lfile*, creates a sequencer for the lock of *lfile*, call it *lfseq*, and passes it to the server. The other threads read the *lfile* and discover that they are replicas. Periodically, they check the sequencer *lfseq* to determine if the lock is still valid. The example illustrates the use of Chubby as a name server; in fact, this is one of the most frequent uses of the system.

Chubby locks and Chubby files are stored in a replicated database. The architecture of these replicas shows that the stack consists of: (i) the Chubby component implementing the Chubby protocol for communication with the clients; and (ii) the active components writing log entries and files to the local storage of the replica; see Fig. 7.9.

An *atomicity log* for a transaction processing system allows a crash recovery procedure to undo all-or-nothing actions that did not complete or to finish all-or-nothing actions that committed but did not record all of their effects. Each replica maintains its own copy of the log; a new log entry is appended to the existing log, and the Paxos algorithm is executed repeatedly to ensure that all replicas have the same sequence of log entries.

The next element of the stack is responsible for the maintenance of a fault-tolerant database, in other words, that all local copies are consistent. Fault tolerance is a property enabling a system to continue operating properly in the event of the failure of one or more faults of some of its components. The database consists of the actual data, or the *local snapshot* in Chubby speak, and a *replay log* to allow recovery in case of failure. The state of the system is also recorded in the database.

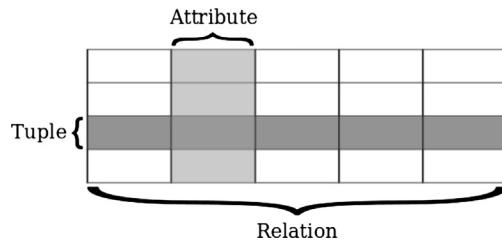
Paxos algorithm is used to reach consensus on sets of values, e.g., the sequence of entries in a replicated log. To ensure that the algorithm succeeds, in spite of the occasional failure of a replica, the following three phases of the algorithm are executed repeatedly:

**FIGURE 7.9**

Chubby replica architecture. Chubby implements the communication protocol with clients. The system includes a component to transfer files to a fault-tolerant database and a fault-tolerant log component to write log entries. The fault-tolerant log uses the Paxos algorithm to achieve consensus. Each replica has its own local file system; replicas communicate with one another using a dedicated interconnect and communicate with the clients through a client network.

- i. Elect a replica to be the master/coordinator. When a master fails, several replicas may decide to assume the role of a master. To ensure that the result of the election is unique, each replica generates a sequence number larger than any sequence number it has seen, in the range $(1, r)$ where r is the number of replicas. Then it broadcasts a *propose* message with this sequence number. The replicas that have not seen a higher sequence number broadcast a *promise* reply and declare that they will reject proposals from other candidate masters. The replica who sent the *propose* message is elected as the master if the number of respondents represent a majority of replicas.
- ii. The master broadcasts to all replicas an *accept* message, including the value it has selected and waits for replies, either *acknowledge* or *reject*.
- iii. Consensus is reached when the majority of the replicas send the *acknowledge* message; then, the master broadcasts the *commit* message.

Paxos algorithm implementation is far from trivial: While the algorithm can be expressed as a few tens of lines of pseudocode, its actual implementation could be several thousand lines of C++ code [92]. Moreover, the practical use of the algorithm cannot ignore the wide variety of failure modes, including algorithm errors and bugs in its implementation, and testing a software system of a few thousands lines of codes is challenging.

**FIGURE 7.10**

Data organization in relational databases.

7.8 RDBMS—cloud mismatch

Before the age of cloud computing, several data models were widely used: the *hierarchical* model for strictly hierarchical relations, the *network* model for many-to-many relationships, and, the most ubiquitous of all, the *relational* model (RDBMS) of Edgar F. Codd from IBM Almaden Research Center. Codd proposed a set of 12 rules for the model [111], but his rules are usually condensed into one sentence: *present data as relations and provide relational operators to manipulate data*.

The relational model organizes data into *tables/relations* representing one *entity type*. Tables consists of *columns/attributes* and rows, called records/tuples; a unique *key* identifies each row. Rows represent instances of a type of entity and columns represent values attributed to instance; see Fig. 7.10.

Structured Query Language (SQL) is a special-purpose language for managing structured data in a relational database system and the centerpiece of this storage technology. SQL has three components: a data definition language, a data manipulation language, and a data control language. Oracle, MySQL, IBM DB2, PostgreSQL, Hive, and Microsoft's SQL Server, Access, and Azure SQL Database are widely used relational database systems.

The set of possible values for a given attribute is described by a *domain*, and *constraints* can restrict the types and range of values of an attribute. Indices, created on any combination of attributes of a relation, usually implemented via B+ trees,² R-trees,³ and bitmaps, provide data access. Relational database queries are expressed using *relational algebra* (a set of algebraic structures with a well-defined semantics) and *relational calculus*.

RDBMS instances reference one another and form a graph, but relational schemas are based on relational algebras linking heterogeneous tuples. RDBMS applications are often written in object-oriented (OO) languages, which encapsulate objects hiding their internal representation. Type system differences between RDBMS and OO languages represent a major mismatch between relational semantics and OO languages. Operator semantics and scalar types are another facet of this mismatch. Further-

² A B+ tree is an m-arry tree (a rooted tree in which each node has no more than m children), with a variable, but often large, number of children per node. A B+ tree can be viewed as a B-tree in which each node contains only keys, and an additional level is added at the bottom with linked leaves. In a block-oriented storage context, e.g., filesystems, B+ trees, retrieve data efficiently.

³ R-trees are data structures used for spatial access methods, i.e., for indexing multidimensional information such as geographical coordinates or polygons. R-trees group nearby objects and represent them with a minimum bounding rectangle (thus, the R in the name) in the next higher level of the tree.

more, relational models prohibit pointers, while OO languages use extensively pointers to reference objects.

RDBMS transactions are much larger than operations by classes in OO languages. Moreover, RDBMS transactions are dynamically bounded sets of arbitrary data manipulations, whereas the granularity of transactions in an OO language is typically on the level of individual assignments to primitive-typed fields. All these facts contribute to the so-called *object-relational impedance mismatch* [427].

The values in a relational tuple cannot contain their own structure, like a nested record. In-memory data structures support much richer representations than relations, so this mismatch does not affect in-memory data representation. In-memory data structures must be translated to a relational representation before being stored on the disk. This mismatch is one of the main reasons for the transition to NoSQL cloud databases.

A critical requirement for cloud databases is *scalability*, the ability to store and process efficiently huge amounts of data distributed over a large set of storage devices. We distinguish two types of scalability, horizontal and vertical; in the former additional resources are provided by a larger number of storage servers, while in the latter additional resources are provided by more powerful servers.

RDBMS are not designed to run effectively on a large number of storage servers without a single point of failure. A case in point, the Microsoft SQL server runs on a cluster-aware file system storing state information on a highly available disk system, a single point of failure.

In Section 3.17, we have seen that the CAP theorem states that a distributed system cannot guarantee consistency, availability, and partition tolerance at the same time. RDBMS systems guarantee consistency, and sharding⁴ makes them tolerant to partitioning, therefore they cannot guarantee availability. That's why a standard RDBMS cannot scale very well: It is not able to guarantee availability.

Cloud computing brought along the demand for storing unstructured or semistructured data, scalability, and effectiveness, thus the need for a new database model, the NoSQL model discussed next.

7.9 NoSQL databases

Convenience prevailed when naming this new database model. An accurate description of this model, possibly *NO-Relational-database*, lacked appeal, so the name NoSQL was rapidly adopted by the community. But this name is misleading.

Michael Stonebreaker, the 2014 Turing Award winner for major contributions to the theory and practice of database systems, writes, “blinding performance depends on removing overhead. Such overhead has nothing to do with SQL, but instead revolves around traditional implementations of ACID transactions, multi-threading, and disk management” [458].

NoSQL refers to a set of databases developed in late 1990s that do not rely on the relational model. NoSQL databases differ from each other, run well on large clusters, and are open source. NoSQL schema is defined implicitly by the way application code utilizes the database. NoSQL does not reflect an advance in storing technology, rather a response to practical needs to efficiently access very large datasets stored on large computer clusters [75]. In the new model RDBMS names changed: a *partition* is a *shard*, a *table* is a *document root element*, a *row* is an *aggregate/record*,

⁴ A shard is a horizontal partitioning of a database, a row in a table structured data.

and a *column* is an *attribute/field/property*. There is no stand-alone query language for NoSQL databases.

NoSQL distribution models. Scalability is an essential attribute of a cloud database, therefore the ability to store very large volumes of data on a large number of storage servers is a major requirement for NoSQL databases. The mean time to failure in a large-scale storage systems built with off-the-shelf components is rather short, therefore the distribution model should support effective replication required for database availability and sharding. Replication stores data on multiple servers, while sharding partitions the data.

The distribution model is based on an aggregate data model that eliminates the relations and transactions of the relational model in favor of more complex data structures allowing lists and other record structures to be nested. The aggregate structures are naturally replicated and sharded; storing relevant data on the same storage server minimizes the number of servers accessed to respond to a query.

The NoSQL distribution models are: (i) sharding; (ii) master-slave replication; and (iii) peer-to-peer (P2P) replication. *Sharding* is based on horizontal scalability; the data is distributed to the servers aiming to balance the workload. The aggregates are stored together if accessed together; the aggregate distribution to storage servers may change over time in response to access pattern changes. Sharding optimizes the writes.

Master-slave replication designates one server as master and all other storage servers as slaves. The master is responsible for processing updates and is the authoritative data source. Servers are synchronized with the master. A drawback of this distribution model is inconsistency as updates take time to propagate from the master to the slaves. This distribution model is optimal for read-intensive database applications and horizontal scaling because new slaves can be added easily.

P2P replication. P2P systems are the most democratic; servers have the same functionality and accept write operations; and losing one or a few servers may affect performance but not system functionality. The price to pay is potential lack of consistency; a *write-write conflict* is unavoidable when a record is updated on two different servers at the same time. These models can be combined, for example, a P2P and sharding combination is appealing to column-family databases.

NoSQL database types. Four flavors can be distinguished:

1. *Key-value.* This type of NoSQL databases model data as an index key and a value, similar with hash tables. Access to database is via primary key. Modeling relationships among different sets of data, correlating data with different sets of keys, searching based on the value of the key-value pair, and operations on sets are not recommended for key-value databases.
2. *Document databases.* Similar to key-value, but the value associated with a key contains structured or semistructured data. These databases enjoy a number of useful features such as: (i) ability to query data in the document without having to retrieve the whole document by key; (ii) support nested documents within the value field. Document databases should not be used for complex transactions spanning different operations because atomicity is only guaranteed within a single document, not across documents.
3. *Column-family databases.* Such databases store large sparse tables with a very large number of rows and only a few columns. They store data in column families as rows. Several columns are associated with a row key for data often accessed together. Super columns consist of a map of columns; there are also super column families.

- 4. Graph databases.** In graph databases, the vertices represent entities, and the directional edges represent relationships among the entities. The graphs are meant to be stored with a single organizational structure and then interpreted in different ways based on the relationships. Scaling graph databases is different than for aggregate-oriented models, and sharding is difficult because the nodes are relationship-oriented rather than aggregate-oriented. Horizontal scaling is ineffective for graph databases, only vertical scaling is suitable.

Rating of several NoSQL databases. Recall from Section 3.17 that PACELC theorem states that, in case of (P), network partitioning, one has to choose between (A), availability, and (C), consistency. Else, (E), even when the system is running normally in absence of partitions, one has to choose between latency, and availability. According to [2] the choices made by several NoSQL databases are:

1. **PA/EL**—if a partition occurs, give up consistency for availability, and under normal operation give up consistency for lower latency. Examples: default versions of DynamoDB, Cassandra, Riak, and Cosmos DB.
2. **PC/EC**—to achieve consistency pay availability and latency costs. Examples: fully ACID systems such as VoltDB/H-Store, Megastore, and MySQL Cluster.
3. **PA/EC**—guarantees reads and writes to be consistent. Example: MongoDB.
4. **Tradeoffs between C/A during P, and L/C during E.** For example, Cosmos DB supports five tunable consistency levels and never violates the specified consistency level.

Summary. Cloud stores such as *document stores* and NoSQL databases are designed to scale well, do not exhibit a single point of failure, have built-in support for consensus-based decisions, and support partitioning and replication as basic primitives. Systems such as *SimpleDB* discussed in Section 2.2, *CouchDB* (see <http://couchdb.apache.org/>), or *Oracle NoSQL database* [379] are very popular, though they provide less functionality than traditional databases. The *key-value* data model is very popular. Several such systems including Voldemort, Redis, Scalaris, and Tokyo cabinet are discussed in [89].

The *soft-state* approach in the design of NoSQL allows data to be inconsistent and transfers the task of implementing only the subset of the ACID properties required by a specific application to the application developer. The NoSQL systems ensure that data will be *eventually consistent* at some future point in time, instead of enforcing consistency at the time when a transaction is “committed.”

It was suggested to associate NoSQL databases with the BASE acronym reflecting their relevant properties, Basically Available, Soft state, and Eventually consistent, whereas traditional databases are characterized by ACID properties; see Section 7.3. Data partitioning among multiple storage servers and data replication are also tenets of the NoSQL philosophy; they increase availability, reduce the response time, and enhance scalability.

7.10 Data storage for online transaction processing systems

Many cloud services are based on Online Transaction Processing (OLTP) and operate under tight latency constraints. Moreover, OLTP applications have to deal with extremely high data volumes and are expected to provide reliable services for very large communities of users. It did not take very long for organizations heavily involved in cloud computing, such as Google and Amazon, eCommerce companies, such as *eBay*, and social media networks such as Facebook, Twitter, or LinkedIn, to discover that traditional relational databases are not able to handle the massive amount of data and the real-time demands of online applications critical for their business model.

The search for alternate models to store the data on a cloud is motivated by the need to decrease the latency by caching frequently used data in memory on dedicated servers, rather than fetching it repeatedly. Distributing data to a large number of servers allows multiple transactions to occur at the same time and decreases the response time. The relational schema is of little use for OLTP applications, and conversion to key-value databases seems a much better approach. Of course, such systems do not store meaningful metadata information, unless they use extensions that cannot be exported easily.

Reducing the response time is a major concern of OLTP system designers. The term *memcaching* refers to a general purpose distributed memory system that caches objects in main memory. The system is based on a very large hash table distributed across many servers. A *memcached* system is based on a client–server architecture and runs under several operating systems, including Linux, Unix, Mac OS X, and Windows. The servers maintain a key–value associative array. The API allows clients to add entries to the array and to query it; a key can be up to 250 bytes long, and a value can be not larger than 1 MB. A *memcached* system uses the LRU cache replacement strategy. Scalability is the other major concern for cloud OLTP applications and implicitly for datastores. There is a distinction between *vertical scaling*, where the data and the workload are distributed to systems that share resources, such as cores/processors, disks, and possibly RAM, and *horizontal scaling*, where the systems do not share either the primary or secondary storage [89].

The overhead of OLTP systems is due to four sources with equal contributions: logging, locking, latching, and buffer management. Logging is expensive because traditional databases require transaction durability, thus every write to the database can only be completed after the log has been updated. To guarantee atomicity, transactions lock every record, and this requires access to a lock table.

Many operations require multithreading, and the access to shared data structures, such as lock tables, demands short-term latches⁵ for coordination. The breakdown of the instruction count for these operations in existing DBMS is: 34.6% buffer management, 14.2% latching, 16.3% locking, 11.9% logging, and 16.2% for hand-coded optimization [228].

Today, OLTP databases could exploit the vast amounts of resources of modern computing and communication systems to store the data in main memory rather than rely on disk-resident B-trees and heap files, locking-based concurrency control, and the support for multithreading optimized for the computer technology of past decades [228]. Logless, single threaded, and transactionless databases could replace the traditional ones for some cloud applications.

Data replication is critical not only for system reliability and availability but also for its performance. In an attempt to avoid catastrophic failures due to power blackouts, natural disasters, or other causes (see also Section 1.4), many companies have established multiple data centers located in various geographic regions. Thus, data replication must be done over a wide area network. This could be quite challenging especially for log data, metadata, and system configuration information due to larger communication delays and an increased probability of communication failures. Several strategies are possible, some based on master–slave configurations, other based on homogeneous replica groups.

Master–slave replication can be asynchronous or synchronous. In the first case, the master replicates write-ahead log entries to at least one slave, and each slave acknowledges appending the log record as soon as the operation is done, while in the second case, the master must wait for the acknowledgments from all slaves before proceeding. Homogeneous replica groups enjoy shorter latency and higher

⁵ A latch is a counter that triggers an event when it reaches zero; for example, a master thread initiates a counter with the number of worker threads and waits to be notified when all of them have finished.

availability than master-slave configurations; any member of the group can initiate mutations which propagate asynchronously.

In summary, the “one-size-fits-all” approach in the traditional storage system design is replaced by a flexible one, tailored to the specific requirements of the applications. Sometimes, the data management of a cloud computing environment integrates multiple databases. For example, Oracle integrates its NoSQL database with the HDFS discussed in Section 11.7, with the Oracle Database, and with the Oracle Exadata. Another approach, discussed in Section 7.12, partitions the data and guarantees full ACID semantics within a partition, while supporting eventual consistency among partitions.

7.11 BigTable

BigTable is a distributed storage system developed by Google to store massive amounts of data and to scale up to thousands of servers [94]. The system uses GFS discussed in Section 7.6 to store user data and system information. BigTable uses the Chubby distributed lock service to guarantee atomic *read* and *write* operations. Directories and files in Chubby’s namespace are used as locks. Client applications written in C++ can add/delete values, search for a subset of data, and lookup for data in a row.

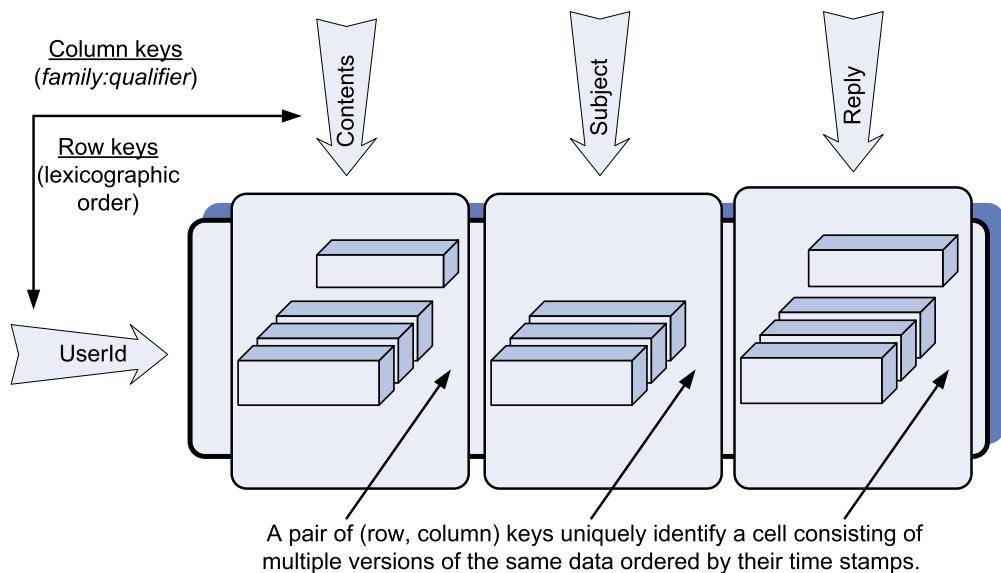
BigTable is based on a simple and flexible data model and allows an application developer to exercise control on the data format and layout. It also reveals data locality information to the application clients. Column keys identify units of access control called *column families*, including data of the same type. A column key consists of a string defining the family name, a set of printable characters, and an arbitrary string as qualifier. A row key is an arbitrary string of up to 64 KB, and a row range is partitioned into *tablets* serving as units for load balancing. Any *read* or *write* row operation is atomic, even when it affects more than one column. Time stamps used to index different versions of the data in a cell are 64-bit integers. The interpretation of time stamps can be defined by the application, while the default is the time of an event in microseconds.

The organization of a BigTable—see Fig. 7.11—shows a sparse, distributed, multidimensional map for an email application. The system consists of three major components: a library linked to application clients to access the system, a master server, and a large number of tablet servers. The master server controls the entire system, it assigns tablets to tablet servers and balances the load among them, manages garbage collection, and handles table and column family creation and deletion.

Internally, the space management is ensured by a three-level hierarchy: the *root tablet* whose location is stored in a Chubby file, points to entries in the second element, the *metadata tablet* which, in turn, points to *user tablets*, collections of locations of user’s tablets. An application client searches through this hierarchy to identify the location of its tablets and then caches addresses for further use.

BigTable performance reported in [94] is summarized in Table 7.4 which shows the number of random and sequential *read* and *write* and scan operations for 1 000 bytes, when the number of servers increases from 1 to 50, then to 250, and finally to 500. Locking prevents the system from achieving a linear speedup, but the BigTable performance is remarkable due to a fair number of optimizations. For example, the number of scans on 500 tablet servers is 7843×500 instead of $15\,385 \times 500$. It is reported that only 12 clusters use more than 500 tablet servers, while some 259 clusters use between 1 and 19 tablet servers.

BigTable is used by a variety of applications including Google Earth, Google Analytics, Google Finance, and web crawlers. For example, Google Earth uses two tables, one for preprocessing and

**FIGURE 7.11**

BigTable example; the organization of an email application as a sparse, distributed, multidimensional map. The slice of the BigTable shown consists of a row with the *UserId* key and three *family* columns; the *Contents* key identifies the cell holding the contents of emails received, the one with the *Subject* key identifies the subject of emails, and the one with the *Reply* key identifies the cell holding the replies; the version of records in each cell are ordered according to their time stamps. The row keys of this BigTable are ordered lexicographically; a column key is obtained by concatenating the *family* and the *qualifier* fields. Each value is an uninterpreted array of bytes.

Table 7.4 BigTable performance; the number of operations per tablet server.

Number of tablet servers	Random read	Sequential read	Random write	Sequential write	Scan
1	1 212	4 425	8 850	8 547	15 385
50	593	2 463	3 745	3 623	10 526
250	479	2 625	3 425	2 451	9 524
500	241	2 469	2 000	1 905	7 843

one for serving client data. The preprocessing table stores raw images; the table is stored on disk as it contains some 70 TB of data. Each row of data consists of a single imagery; adjacent geographic segments are stored in rows in close proximity to one another. The column family is very sparse; it contains a column for every raw image. The preprocessing stage relies heavily on MapReduce to clean and consolidate the data for the serving phase. The serving table is stored on GFS and is “only” 500 GB, and it is distributed across several hundred tablet servers that maintain in-memory column families; this organization enables the serving phase of Google Earth to provide a fast response time to tens of thousands of queries per second.

Google Analytics provides aggregate statistics such as the number of visitors of a web page per day. To use this service, web servers embed a JavaScript code in their web pages to record information every time a page is visited. The data is collected in a *raw click* BigTable of some 200 TB with a row for each end-user session. A *summary* table of some 20 TB contains predefined summaries for a website.

7.12 Megastore

Megastore is a scalable storage for online services. The system, distributed over several data centers, has a very large capacity and is highly available. Megastore is widely used internally at Google; in 2011, it had a capacity of 1PB, handled some 23 billion transactions daily, 3 billion *write*, and 20 billion *read* transactions [43].

The basic design philosophy of the system is to partition the data into *entity groups* and replicate each partition independently in data centers located in different geographic areas. The system supports full ACID semantics within each partition and provides limited consistency guarantees across partitions; see Fig. 7.12. Megastore supports only those traditional database features that allow the system to scale well and do not affect drastically the response time.

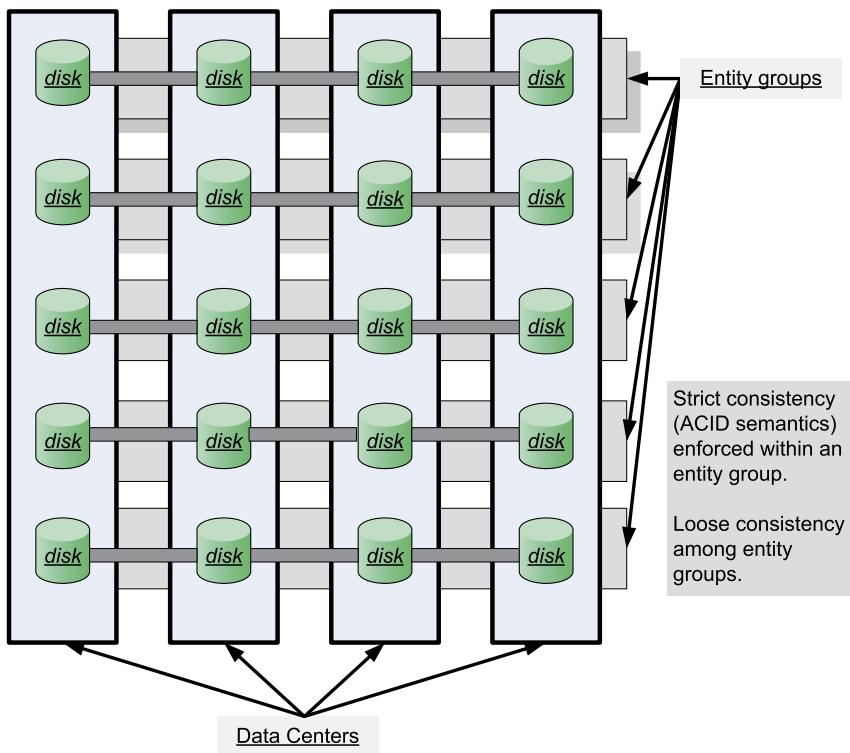
Another distinctive feature of the system is the use of the Paxos consensus algorithm discussed in Section 10.13 for replicating primary user data, metadata, and system configuration information across data centers and for locking. The version of the Paxos algorithm used by Megastore does not require a single master, but instead any node can initiate *read* and *write* operations to a write-ahead log replicated to a group of symmetric peers.

Entity groups are application-specific and store together logically related data; for example, an email account could be an entity group for an email application. Data should be carefully partitioned to avoid excessive communication between entity groups. Sometimes, it is desirable to form multiple entity groups as is the case of blogs [43].

This middle ground between traditional and NoSQL databases taken by the Megastore designers is also reflected by the data model. The data model is declared in a *schema* consisting of a set of *tables* that are composed of *entries*, each entry being a collection of named and typed *properties*. The unique primary key of an entity in a table is created as a composition of entry properties. A Megastore table can be a *root* or a *child* table; each *child entity* must reference a special entity, called *root entity* in its root table. An entity group consists of a primary entity and all the entities that reference it.

The system makes extensive use of BigTable. Entities from different Megastore tables can be mapped to the same BigTable row without collisions. This is possible because the BigTable column name is a concatenation of the Megastore table name and the name of a property. A BigTable row for the root entity stores the transaction and all metadata for the entity group. Multiple versions of the data with different time stamps can be stored in a cell as we have seen in Section 7.11.

Megastore takes advantage of this feature to implement *multiversion concurrency control*. When a mutation of a transaction occurs, this mutation is recorded along with its time stamp, rather than marking the old data as obsolete and adding the new version. This strategy has a couple of several advantages: *Read* and *write* operations can proceed concurrently, and a *read* always returns the last fully updated version.

**FIGURE 7.12**

Megastore organization. Data is partitioned into *entity groups*; full ACID semantics within each partition and limited consistency guarantees across partitions are supported. A partition is replicated across data centers in different geographic areas.

A *write* transaction involves several steps: (1) get the time stamp and the log position of the last committed transaction; (2) gather the *write* operations in a log entry; (3) use the consensus algorithm to append the log entry and then commit; (4) update the BigTable entries; and (5) cleanup.

7.13 Storage reliability at scale

Building reliable systems with unreliable components is a major challenge in system design identified and studied early on by John von Neumann [497]. This challenge is greatly amplified, on one hand, by the scale of the cloud computing infrastructure and by the use of off-the-shelf components that reduces the infrastructure cost and, on the other hand, by the latency constraints of many cloud applications.

Even though the mean time to failure of individual components can be on the order of month or years, it is unavoidable to witness a small, but significant, number of server and network components

that are failing at any given time. Data losses cannot be tolerated, thus the failure of storage devices is a major concern. It is left to the software to mask the failures of storage devices and avoid data loss.

Dynamo and DynamoDB. Amazon developed two database systems to support reliability at scale. Dynamo, a highly available key-value storage system, has been solely used by AWS core services for in-house applications since 2007 [133]. In 2012 DynamoDB, a NoSQL database service for latency-sensitive applications that need consistent access at any scale, was opened to AWS user community. Dynamo and DynamoDB use a similar data model, but Dynamo had a multimaster design requiring the client to resolve version conflicts, whereas DynamoDB uses synchronous replication across multiple data centers for high durability and availability.

DynamoDB is a fully managed database service designed to provide an “always-on” experience. It supports both document and key-value store models and has been used for mobile, web, gaming, IoT, advertising, real-time analytics, and other applications. DynamoDB stores data on SSDs to support latency-sensitive applications; typical requests take milliseconds to complete. DynamoDB allows developers to specify the throughput capacity required for specific tables within their database using the *provisioned throughput* feature to deliver predictable performance at any scale. The service is integrated with other AWS services, e.g., it offers integration with Hadoop via Elastic MapReduce.

Design objectives. Dynamo’s primary concern is high availability where updates are not rejected even in the wake of network partitions or server failures. Dynamo has to deliver predictive performance, in addition to reliability and scalability. Services supported by Dynamo have stringent latency requirements, and this precludes supporting ACID properties. Indeed, data stores providing ACID guarantees tend to exhibit poor availability.

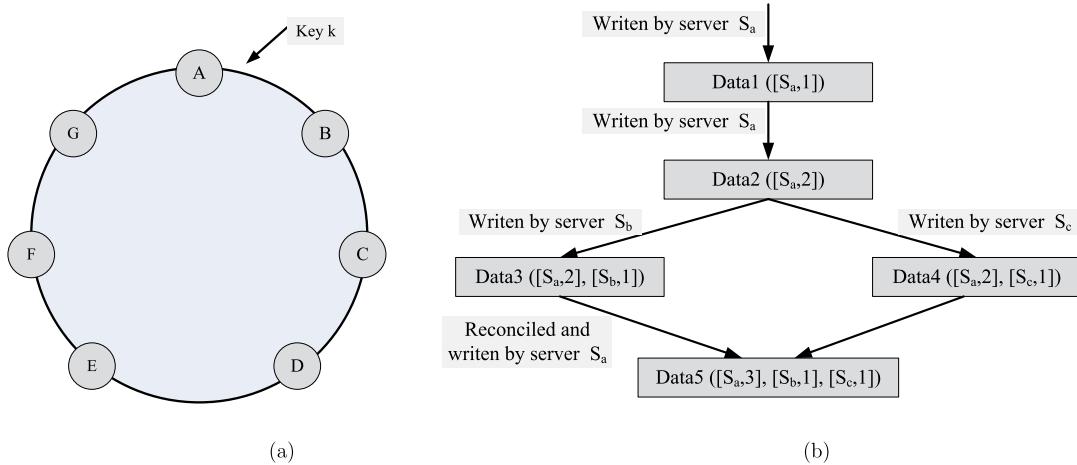
The most significant design considerations regard data replication and measures to increase availability in the wake of failures. Strong consistency and high data availability cannot be achieved simultaneously. Availability can be increased by optimistic replication, allowing changes to propagate to replicas in the background, while disconnected work is tolerated.

In traditional data stores, writes may be rejected if the data store cannot reach all, or a majority of, the replicas at a given time. This approach is not tolerated by many AWS applications. As a consequence, rather than implementing conflict resolution during writes and keeping the read complexity simple, Dynamo increases the complexity of conflict resolution of the read operations.

Dynamo supports simple read and write operations to data items uniquely identified by a key. The *get(key)* operation locates object replicas associated with the key in the storage system and returns a single object or a list of objects with conflicting versions along with a context. The *put(key, context, object)* operation determines where the replicas of the object should be placed based on the associated key and writes replicas to the secondary storage.

The *context* encodes system metadata about the object that is opaque to the caller and includes information such as the version of the object. Context information is stored along with the object so that the system can verify the validity of the context object supplied in the *put* request. The main techniques used to achieve Dynamo’s design objectives are:

1. *Incremental scalability* ensured by consistent hashing.
2. *High write availability* based on the use of vector clocks with reconciliation.
3. *Handling temporary failures* using sloppy quorum and hinted handoff. This provides high availability and durability guarantees when some of the replicas are not available.

**FIGURE 7.13**

- (a) Dynamo servers are organized as a ring. Ring nodes B, C, and D store keys in range (A, B), including the key, k .
 (b) Evolution of an object in time using vector clocks.

4. *Permanent failure recovery* based on anti-entropy and Merkle trees.⁶ This technique synchronizes divergent replicas in the background.
5. *Gossip-based membership protocol and failure detection*. This technique preserves symmetry and avoids having a centralized registry for storing membership and node liveness information.

Scaling, load balancing, and replication. Data partitioning scheme based on *consistent hashing* is designed to support incremental system scaling. The output of a hash function is treated as a ring, and each node in the system is assigned a random value representing its position on the ring.

Consistent hashing reduces the number of keys to be remapped when a hash table is resized. On average, only K/n keys need to be remapped, with K being the number of keys and n being the number of slots. In most traditional hash tables, a change in the number of slots causes nearly all keys to be remapped because the mapping between the keys and the slots is defined by a modular operation.

A data item identified by a key is assigned to a storage server by hashing the data item key to yield its position on the ring and then walking the ring clockwise to find the first node with a position larger than the position of the item. Each storage server is responsible for the ring region between itself and its predecessor in the ring; see Fig. 7.13(a). In this example, node B replicates key k at nodes C and D, in addition to storing it locally. Node D will store keys in the ranges $(A, B]$, $(B, C]$, and $(C, D]$.

Instead of mapping a storage server to a single point on the ring, the system uses the concept of “virtual nodes” and assigns it to multiple points on the ring. A physical server is mapped to multiple nodes of the ring. This form of virtualization supports:

⁶ A Merkle or hash tree is a tree in which every leaf node is labelled with the cryptographic hash of a data block, and every non-leaf node is labelled with the cryptographic hash of the labels of its child nodes. Hash trees allow efficient and secure verification of the contents of large data structures.

1. Load balancing. When a storage server is unavailable, its load is dispersed among available servers. When the server comes back again, it is added to the system and accepts a load roughly equivalent to the load of other servers.
2. System heterogeneity. The number of virtual nodes a physical server is mapped to depends on its capacity.

A data item is replicated at N servers. Each key, k , is assigned to a coordinator charged with the replication of the data items in its range. In addition to locally storing each key within its range, the coordinator replicates these keys at the $N-1$ clockwise successor nodes in the ring. Each node is responsible for the ring region between it and its N -th predecessor. The *preference list* is the list of all nodes responsible for storing a particular key.

Eventual consistency. This strategy allows *updates to be propagated to all replicas asynchronously*. A versioning system allows multiple versions of a data object to be present in the data store at the same time. The result of each modification is a new and immutable data version. New versions often subsume the older ones, and the system can use syntactic reconciliation to determine the authoritative version.

The system uses *vector clocks*, lists of (node, counter) pairs, to capture causality of each version of a data object. When a client updates an object, it must specify the version it is updating by passing the context it obtained from an earlier read operation, which contains the vector clock information. System failures combined with concurrent updates lead to conflicting versions of an object and version branching. Given two versions of the same object, the first is an ancestor of the second and can be forgotten if the counters on the first object's clock are less than or equal to all of the nodes in the second clock; otherwise, the two changes are in conflict and require reconciliation.

Fig. 7.13(b) illustrates the versioning mechanism for the following sequence of events. Data is written by server S_a , and object *Data1* with associated clock $[S_a, 1]$ is created. The same server S_a writes again, and object *Data2* with associated clock $[S_a, 2]$ is created. *Data2* is a descendent of *Data1* and overwrites it. There may be replicas of *Data1* at servers that have not yet seen *Data2*. Then, the same client updates the object, and server S_b handles the request; a new object data *Data3* and its associated clock $[(S_a, 2), (S_b, 1)]$ are created.

A different client reads *Data2* and tries to update it; this time, server S_c handles her request. A new object *Data4*, a descendent of *Data2*, with version clock $[(S_a, 2), (S_c, 1)]$ is created. Upon receiving *Data4* and its clock, a server aware of *Data1* or *Data2* could determine that both are overwritten by the new data and can be collected garbage.

A node aware of *Data3* and *Data4* will see that there is no causal relation between them because there are changes not reflected in each other. Both versions of the data must be kept and presented to a client for semantic reconciliation. If a client reads both *Data3* and *Data4*, its context will be $[(S_a, 2), (S_b, 1), (S_c, 1)]$, the summary of the virtual clocks of both data objects. If the client performs a reconciliation and the write request is handled by server S_a , then the vector clock of the new data, *Data5*, will be $[(S_a, 3), (S_b, 1), (S_c, 1)]$. The size of the vector clock grows, but in practice this growth is limited.

Sloppy quorum for handling failures. During server failures and network partitions a strict quorum membership is enforced by traditional systems. This conflicts with the durability requirement, and in Dynamo all read and write operations are performed on the first N *healthy nodes* from the preference list. In this *sloppy quorum*, the healthy nodes may not always be the first N nodes encountered while walking the consistent hashing ring.

For example, to maintain the desired availability and durability guarantees, when node A in Fig. 7.13(a) is unreachable during a write operation, a replica that would normally have been sent to A will be sent to D . The metadata of this replica will include a hint indicating the intended recipient of the replica.

Replica synchronization in case of permanent failures. Replica inconsistencies can be detected faster and with a minimum of data transfer using Merkle trees.⁷ This allows each branch of the tree to be checked independently without the need to download the entire tree. For example, if the hash values of the root of two trees are equal, then the values of the leaf nodes in the tree are equal, and the nodes require no synchronization.

Anti-entropy is a process of comparing the data of all replicas and updating each replica to the newest version; Merkle trees are used for anti-entropy. The key range is the set of keys covered by a virtual node. Each node maintains a separate Merkle tree for each key range, and two nodes exchange the root of the Merkle tree corresponding to the key ranges that they host in common to compare whether the keys within a key range are up-to-date.

Gossip-based node addition to or removal from the ring. The node receiving the request writes the change and its time of issue to persistent store. A gossip-based protocol propagates membership changes and maintains an eventually consistent view of the membership. Each node contacts a peer chosen at random every second, and the two reconcile their persisted membership change histories.

For example, consider the case when a new node Q is added between nodes A and B to the ring in Fig. 7.13(a). Now node Q will be storing keys in the ranges $(F, G]$, $(G, A]$ and $(A, Q]$ freeing nodes B , C , and D from storing the keys in these ranges. Upon confirmation from Q , nodes B , C , and D will transfer the appropriate set of keys to it. When a node is removed from the system, the reallocation of keys proceeds as a reverse process.

7.14 Disk locality versus data locality in computer clouds

Locality is critical for the performance of computing systems. Recall that a sequence of references is said to have spatial locality if the items referenced within a short time interval are close in space, e.g., they are at nearby memory addresses or nearby sectors on a disk. A sequence exhibits temporal locality if accesses to the same item are clustered in time.

Locality has major implications for memory hierarchies. The performance of a processor is highly dependent on the ability to access code and data in the cache rather than memory. Virtual memory can only be effective if the code and the data exhibit spatial and temporal locality so that page faults occur infrequently. Optimizing locality in cloud computing is a challenging problem, and a fair amount of effort has been devoted to algorithms and systems designed to increase the fraction of the tasks of a job enjoying locality. Locality improves the performance of cloud applications for two main reasons: (i) disk bandwidth is larger than the network bandwidth, and the off-rack communication bandwidth is oversubscribed and affects the off-rack disk access; and (ii) the better performance of I/O-intensive

⁷ A Merkle tree is a hash tree where leaves are hashes of the values of individual keys. Parent nodes higher in the tree are hashes of their respective children.

applications, when data is stored locally, is due to the lower latency and the higher bandwidth of a local disk versus the latency and the bandwidth of a remote disk.

An important question for cloud resource management is whether *disk locality* is important. Intuitively, we expect the answer to the question whether tasks should be dispatched to the cluster nodes, where their input data resides would be a resounding “Yes.” A slightly different view on the subject of disk locality is expressed in a paper with the blunt title “Disk-Locality in Datacenter Computing Considered Irrelevant” [23].

Maybe we tried to solve the wrong problem, and instead of focusing on disk locality, we should focus on data locality. In other words, we should look at the local memory as a “data cache” and make sure that data is stored in the local memory rather than the local disk of the processor where the task is scheduled to run. Data transfer through the network may still be necessary, but why store it on the disk and then load it in memory? Several arguments support this thesis:

- a. Networking technology improves at a faster pace than hard disk technology. Switches with aggregate link speeds of 40 Gbps and 100 Gbps are available today. Rates of 10 and 25 Gbps of server network interfaces will be available soon.
- b. The bandwidth available to applications will increase as data centers adopt bisection topologies for their interconnects [52]. Recall that the bisection bandwidth is the sum of the bandwidths of the minimal number of links that are cut when splitting the system into two parts.
- c. The latency of a read access to a local disk is only slightly lower than the latency of a read to a disk in the same rack; local disk access is about 8% faster [52]. Access to a disk in a different rack will not increase drastically due to faster networks.
- d. Though the cost of solid-state disk technology is dropping at a impressive rate, i.e., 50% per year, solid-state disks are unlikely to replace hard disk drives any time soon due to the sheer volume of data stored on computer clouds. To be competitive with HDDs, the cost-per-byte of SSD should be reduced by up to three orders of magnitude.
- e. Accessing data in local memory is two orders of magnitude faster than reading from a local disk. An increasing number of applications use the processor memory distributed across the nodes of large clusters to cache the data, rather than access the HDDs.
- f. There is at least a two orders of magnitude discrepancy between the capacity of the disks and the memory of today’s clusters, thus it seems reasonable to use processor memory as a cache for the much larger volume of data stored on the disks.

It makes sense to use processor memory as a cache for large data sets used as the input of an application if and only if: (1) the size of the application’s input is much larger than the memory size; and (2) the applications exhibit locality and have a modestly sized working set, in other words, if they access frequently only a relatively small fraction of data blocks.

Hadoop tasks running at a data center with some 3 000 machines interconnected by three-tiered networks are investigated in [23]. Most Facebook jobs are data-intensive and spend most of their execution time reading input data. The analysis of job traces led to the conclusion that there is an order of magnitude discrepancy between the input size and the memory size and that some 75% of the data blocks are accessed only once.

The analysis also shows that workloads have a heavy tail distribution of block access. Moreover, 96% of the data inputs can fit into a fraction of the total cluster memory for a majority of jobs. Some

64% of all jobs investigated perform well under a least frequently used (LFU) replacement policy applied to all their tasks.

The *task progress report*, \mathbb{T} , a metric for the effect of locality on the duration of a task, is defined as the ratio

$$\mathbb{T} = \frac{\text{data_read} + \text{data_written}}{\text{task_duration}}. \quad (7.1)$$

The results of measurements comparing node-local versus rack-local tasks show that 85% of the jobs have $0.9 \leq \mathbb{T} \leq 1.1$ and only 4% of jobs have $\mathbb{T} \leq 0.7$.

Data compression reduces the pressure on the data center disk space. Tasks read compressed data and uncompress it before processing. In spite of a network oversubscribed by a factor of ten at Facebook, data compression leads to very good results; running a task off-rack is only 1.2 times to 1.4 times slower compared to on-rack execution. The log analysis of Hadoop jobs suggests that prefetching data blocks could be beneficial because a large fraction of data blocks are accessed only once.

7.15 Database provenance

Data provenance or *lineage* describes the origins and the history of data and adds value to data by explaining how it was obtained. The lineage of a tuple \mathcal{T} in the result of a query is the set of items contributing to produce \mathcal{T} . The lineage is important for the extract-transform-load process and for incrementally adding and updating a database.

Before the Internet era the information in databases was trusted, and it was assumed that the organization maintaining the database was trustworthy and that every effort to ensure data veracity was made. This assumption is no longer true; there is no centralized control over data integrity because the Internet allows data to be indiscriminately created, copied, moved around, and combined. Establishing data provenance is necessary for all databases and is critical for cloud databases as the data owners relinquish control of their data to CSPs.

Data provenance has been practiced by the scientific and engineering community for some time, long before the disruptive effects of data democratization brought about by the Internet. Data collected by scientific experiments contains information about the experimental setup and the settings of measuring instruments for each batch of data. Ensuring that experiments can be replicated has always been an essential aspect of scientific integrity. The same requirements apply to engineering when test data, e.g., data collected during the testing of a new aircraft, include lineage information.

Data provenance could explain *why* an output record was produced, describing in detail *how* the record was produced, and/or explaining *where* output data comes from. The why-, how-, and where-provenance are analyzed in [105] and discussed briefly in this section.

The *witness* of a database record is the subset of database records ensuring that the record is the output of a query. The *why-provenance* includes information about the witnesses to a query. The number of witnesses can be exponentially large. To limit this number, the concept of *witness bases* of tuple \mathcal{T} in query \mathcal{Q} on database \mathcal{D} is defined as the particular set of witnesses that can be calculated efficiently from \mathcal{Q} and \mathcal{D} .

The *why-provenance* of an output tuple \mathcal{T} is the witness basis of \mathcal{T} according to \mathcal{Q} . The witness basis depends on the structure of the query, and it is sensitive to the query formulation. We now take a short

Table 7.5 Two queries Q and Q' of instance I are equivalent under R for different pairs of A and B elements of the three tuples t, t', t'' .

R	Output of $Q'(I)$		Output of $Q(I)$	
	A	B	A	B
$t:$	1	2	1	2
$t':$	1	3	1	3
$t'':$	4	2	4	2

detour for introducing *Datalog conjunctive query* notations before presenting an example showing that the why-provenance is sensitive to query rewriting.

Datalog is a declarative logic programming language. Query evaluation in Datalog is based on first order logic, thus it is sound and complete. A Datalog program includes *facts* and *rules*. A rule consists of two elements, the head and the body, separated by the “:-” symbol. A rule should be understood as: “head” if it is known that “body”. For example,

- the facts in the left box below mean: (1) Y is in relation R with X; (2) Z is in relation R with Y.
- the rules in the right box mean: (1) Y is in relation P with X if it is known that Y is in relation R with X; and (2) Y is in relation P with X if it is known that Z is in relation R with X AND rule P is satisfied, i.e., Y is in relation P with Z.

$R(X, Y).$
 $R(Y, Z).$

$P(X, Y) :- R(X, Y).$
 $P(X, Y) :- R(X, Z), P(Z, Y)$

Consider now an instance I defining relation R and three tuples t, t', t'' , two queries Q, Q' , and the output of the two queries $Q(I), Q'(I)$. The two queries are

$$\begin{aligned} Q : Ans(x, y) &:- R(x, y) \\ Q' : Ans(x, y) &:- R(x, y), R(x, z). \end{aligned} \tag{7.2}$$

Table 7.5 shows that queries Q and Q' are equivalent since they produce the same result. Table 7.6 shows that the why-provenance is sensitive to query rewriting. There exists a subset of the witness basis invariant under equivalent queries. This subset, called *minimal witness basis*, includes all *minimal witnesses* in the witness basis. A witness is minimal if none of its proper subinstances is also a witness in the witness basis. For example, $\{t\}$ is a minimal witness for the output tuple (1, 2) in Table 7.6, but $\{t, t'\}$ is not a minimal witness since $\{t\}$ is a sub-instance of it and it is a witness to (1, 2). Hence, the minimal witness basis is $\{t\}$ in this case.

The *why-provenance* describes the source tuples that witness the existence of an output tuple in the result of the query, but it does not show how an output tuple is derived according to the query. The *how-provenance* is more general than the *why-provenance* and it is also sensitive to query formulation as shown in Table 7.7. The *how-provenance* of the tuple (1, 2) is t in the output of Q and $t^2 + t \cdot t'$ in the output of Q' .

Table 7.6 The why-provenance of the two equivalent queries Q and Q' are different for the three tuples t , t' and t'' .

Instance I			Output of $Q(I)$			Output of $Q'(I)$		
R	A	B	A	B	<i>why</i>	A	B	<i>why</i>
$t:$	1	2	1	2	$\{\{t\}\}$	1	2	$\{\{t\}, \{t, t''\}\}$
$t':$	1	3	1	3	$\{\{t'\}\}$	1	3	$\{\{t'\}, \{t, t''\}\}$
$(t'')^2:$	4	2	4	2	$\{\{(t'')^2\}\}$	4	2	$\{\{(t'')^2\}\}$

Table 7.7 The how-provenance of the two equivalent queries Q and Q' are different for the three tuples t , t' and t'' .

Instance I			Output of $Q(I)$			Output of $Q'(I)$		
R	A	B	A	B	<i>how</i>	A	B	<i>how</i>
$t:$	1	2	1	2	t	1	2	$t^2 + t \cdot t'$
$t':$	1	3	1	3	t'	1	3	$(t')^2 + t \cdot t'$
$t'':$	4	2	4	2	t''	4	2	$(t'')^2$

The where-provenance describes the relationship between source and output locations, while why-provenance describes the relationship between source and output tuples. The where-provenance is also sensitive to query formulation [105].

7.16 History notes and further readings

A 1989 survey of distributed file systems can be found in [436]. NFS Versions 2, 3, and 4 are defined in RFCs 1094, 1813, and 3010, respectively. NFS Version 3 added a number of features including: support for 64-bit file sizes and offsets, support for asynchronous writes on the server, additional file attributes in many replies, and a *REaddirplus* operation. These extensions allowed the new version to handle files larger than 2 GB, to improve performance, and to get file handles and attributes along with file names when scanning a directory. NFS Version 4 borrowed a few features from the Andrew file system. WebNFS is an extension of NFS Version 2 and 3; it enables operations through firewalls and is easier integrated into Web browsers.

AFS was further developed as an open-source system by IBM under the name OpenAFS in 2000. Locus [500] was initially developed at UCLA in the early 1980s, and its development was continued by Locus Computing Corporation. Apollo [299] was developed at Apollo Computer Inc, established in 1980, and acquired in 1989 by HP. RFS (Remote File System) [41] was developed at Bell Labs in the mid-1980s. Documentation of a current version of GPFS and an analysis of caching strategy are discussed [250] and [440].

Several DBMS generations based on different models have been developed over the years. In 1968 IBM released the Information Management System (IMS) for the IBM 360 computers. IMS was based on the so-called *navigational model* supporting manual navigation in a linked data set where the data is organized hierarchically. In the RDBMS model, introduced by Codd, related records are linked together

and can be accessed using a unique *key*. Codd also introduced a *tuple calculus* as a basis for a query model for a RDBMS; this led to the development of the Structured Query Language.

In 1973, the Ingres research project at UC Berkeley developed a relational database management system. Several companies, including Sybase, Informix, NonStop SQL, and Ingres, were established to create SQL RDBMS commercial products based on the ideas generated by the Ingres project. IBM's DB2 and SQL/DS dominated the RDBMS market for mainframes during the last part of the 1980s. Oracle Corporation, founded in 1977, was also involved in the development of RDBMS.

ACID properties of database transactions were defined by Jim Gray in 1981 [206], and the term ACID was introduced in [227]. The object-oriented programming ideas of 1980s led to the development of Object-Oriented Data Base Management Systems (OODBMS) where the information is packaged as objects. Ideas developed by several research projects, including Encore-Ob/Server at Brown University, Exodus at the University of Wisconsin-Madison, Iris at HP, ODE at Bell Labs, and the Orion project at MCC-Austin, helped the development of several OODBMS commercial products [271].

NoSQL database management systems emerged in the 2000s. They do not follow the RDBMS model, do not use SQL as a query language, may not give ACID grantees, and have a distributed, fault-tolerant architecture.

Further readings. A 2011 article in the journal *Science* [236] discusses the volume of information stored, processed, and transferred through the networks. [359] is a comprehensive study of the storage technology up to 2003. Evolution of storage technology is presented in [251]. Network File System Versions 2, 3, and 4 are discussed in [433], [392], and [393], respectively. [358] and [245] provide a wealth of information about the Andrew File System, while [237] and [363] discuss in detail the Sprite File System. Other file systems such as Locus, Apollo, and the Remote File System (RFS) are discussed in [500], [299], and [41], respectively. The recovery in the Calypso file system is analyzed in [143]. The Lustre file system is analyzed in [378].

The General Parallel File System (GPFS) developed at IBM and its precursor, the TigerShark multimedia file system, are presented in [440] and [230]. A good source for information about the Google File System is [197]. Main memory OLTP recovery is covered in [324]. The development of Chubby is covered by [81]. NoSQL databases are analyzed in several papers including [458], [228], and [89]. BigTable and Megastore developed at Google are discussed in [94] and [43]. An evaluation of distributed data store is reported [78].

Attaching cloud storage to a campus grid is the subject of [150]. A cost analysis of storage in enterprise is discussed in [406], and [455] is an insightful discussion of main-memory OLTP databases. [496] presents VMware storage.

7.17 Exercises and problems

- Problem 1.** Analyze the reasons for the introduction of storage area networks (SANs) and their properties. *Hint:* read [359].
- Problem 2.** Block virtualization simplifies the storage management tasks in SANs. Provide solid arguments in support of this statement. *Hint:* read [359].
- Problem 3.** Analyze the advantages of memory-based checkpointing. *Hint:* read [261].
- Problem 4.** Discuss the security of distributed file systems including SUN NFS, Apollo Domain, Andrew, IBM AIX, RFS, and Sprite. *Hint:* read [436].

- Problem 5.** The designers of the Google file system (GFS) have reexamined the traditional choices for a file system. Discuss the four observations regarding these choices that have guided the design of GFS. *Hint:* read [197].
- Problem 6.** In his seminal paper on the virtues and limitations of the transaction concept [206], Jim Gray analyzes logging and locking. Discuss the main conclusions of his analysis.
- Problem 7.** Michael Stonebreaker argues that “blinding performance depends on removing overhead...” Discuss his arguments regarding the NoSQL concept. *Hint:* read [458].
- Problem 8.** Discuss the Megastore data model. *Hint:* read [43].
- Problem 9.** Discuss the use of locking in the BigTable. *Hint:* read [94] and [81].

Cloud security

8

Security has been a concern since the early days of computing, when a computer was isolated and threats could only be posed by someone with access to the computer room. Once computers were able to communicate with one another, the Pandora's box of threats was wide open. In an interconnected world, various embodiments of malware can migrate easily from one system to another, cross national borders, and infect systems all over the world.

Security of computer and communication systems takes on a new urgency as the society becomes increasingly more dependent on the information infrastructure. Even the critical infrastructure of a nation can be attacked by exploiting flaws in computer security and often human naivety. Malware, such as the Stuxnet virus, targets industrial systems controlled by software [102]. The concept of *cyberwarfare* meaning “actions by a nation-state to penetrate another nation’s computers or networks for the purposes of causing damage or disruption” [109] was recently added to the dictionary.

A computer cloud is a target-rich environment for malicious individuals and criminal organizations. It is, thus, no surprise that security is a major concern for existing users and for potential new users of cloud computing services. Some of the security risks faced by computer clouds are shared with other systems supporting network-centric computing and network-centric content, e.g., service-oriented architectures, grids, and web-based services.

Cloud computing is an entirely new computing paradigm based on a new technology. It is therefore reasonable to expect that new methods to deal with some of existing security threats will be developed, while other perceived threats will prove to be exaggerated [25]. Indeed, “early on in the life cycle of a technology, there are many concerns about how this technology will be used...they represent a barrier to the acceptance...over the time, however, the concerns fade, especially if the value proposition is strong enough” [248].

The breathtaking pace of developments in information science and technology has many side effects. One of them is that standards, regulations, and laws governing the activities of organizations supporting the new computing services, and in particular utility computing, have yet to be adopted. As a result, many issues related to privacy, security, and trust in cloud computing are far from being settled. For example, there are no international regulations related to data security and privacy. Data stored on a computer cloud can freely cross national borders throughout the data centers of the CSP.

The chapter starts with a discussion of cloud users concerns related to security in Section 8.1. In Section 8.2, we elaborate on the security threats perceived by cloud users, as mentioned in Section 2.10. Privacy and trust are covered in Sections 8.4 and 8.5. Encryption discussed in Section 8.6 protects data in cloud storage, but data must be decrypted for processing. Threats during processing caused by flaws in the hypervisors, rogue VMs, or a VMBR, discussed in Section 5.15, cannot be ignored.

Analysis of database service security in Section 8.7 is followed by a presentation of operating system security, VM security, and the security of virtualization in Sections 8.8, 8.9, and 8.10, respectively.

Sections 8.11 and 8.12 analyze the security risks posed by shared images and by a management OS. An overview of the Xoar hypervisor, a version of Xen that breaks the monolithic Design of the TCB, is discussed in Section 8.13, followed by mobile device security in Section 8.14. Section 8.15 and 8.16 cover the current cloud security threats and AWS security.

8.1 Security—the top concern for cloud users

The idea that moving to a cloud liberates an organization from concerns related to computer security and eliminates a wide range of threats to data integrity is accepted by some members of the IT community. Some argue that cloud security is in the hands of experts, hence cloud users are even better protected than before. As we shall see throughout this chapter, these seem to be rather naive points of view. Outsourcing computing to a cloud generates major new security and privacy concerns. Moreover, the Service Level Agreements do not provide adequate legal protection for cloud computer users who are often left to deal with events beyond their control.

Previously, cloud users were accustomed to operate inside a secure perimeter protected by a corporate firewall. Now, they have to extend their trust to the cloud service provider if they wish to benefit from the economical advantages of utility computing. The transition from a model where users have full control of all systems and where their sensitive information is stored and processed is a difficult one. The reality is that virtually all surveys report that security is the top concern of cloud users.

The major user concerns are: unauthorized access to confidential information and data theft. Data is more vulnerable in storage than while it is being processed. Data is kept in storage for extended periods of time, while during processing, it is exposed to threats for a relatively short time. Hence, close attention should be paid to the security of storage servers and to data in transit. There is also the risk of unauthorized access and data theft posed by rogue employees of a CSP. Cloud users are concerned about insider attacks because hiring and security screening policies of a CSP are totally opaque processes.

The next concern regards user control over the lifecycle of data. It is virtually impossible for a user to determine if data that should have been deleted actually was. Even if deleted, there is no guarantee that the media was totally wiped out and that the next user is not able to recover confidential data. This problem is exacerbated because the CSPs rely on seamless backups to prevent accidental data loss. Such backups are done without user knowledge or consent. During this exercise data records can be lost, accidentally deleted, or accessible to an attacker.

Lack of standardization is next on the list of concerns. Today, there are no interoperability standards, as we have discussed in Section 2.6. Many questions do not have satisfactory answers at this time; for example: What can be done when the service provided by the CSP is interrupted? How should one access critically needed data in case of a blackout? What if the CSP drastically raises its prices? What is the cost of moving to a different CSP?

It is undeniable that auditing and compliance pose an entirely different set of challenges in cloud computing. These challenges are not yet resolved. A full audit trail on a cloud is an unfeasible proposition at this time.

Another, less analyzed, user concern is that cloud computing is based on a new technology expected to evolve into the future. Case in point, autonomic computing is likely to enter the scene. When this happens, self-organization, self-optimization, self-repair, and self-healing could generate additional se-

curity threats. In an autonomic system, it will be even more difficult than at the present time to determine when an action occurred, what was the reason for that action, and if and how it created the opportunity for an attack or for data loss. It is still unclear how autonomic computing can be compliant with privacy and legal issues.

There is no doubt that multitenancy is the root cause of many user concerns. Nevertheless, multitenancy enables a higher server utilization, thus lower costs. The users have to learn to live with multitenancy, one of the pillars of utility computing. The threats caused by *multitenancy* differ from one cloud delivery model to another. For example, in the case of SaaS, private information such as names, addresses, phone numbers, and possibly credit card numbers of many users are stored on one server; when the security of that server is compromised, a large number of users are affected.

Users are greatly concerned about the legal framework for enforcing cloud computing security. The cloud technology has moved much faster than cloud security, and privacy legislators and users have legitimate concerns regarding the ability to defend their rights. As the data centers of a CSP may be located in several countries, it is difficult to understand which laws apply: the laws of the country where information is stored and process; the laws of the countries the information crossed when sent by the user; or the laws of the user's country.

To further complicate legal aspects of cloud computing security, a CSP may outsource handling of personal and/or sensitive information. Existing laws stating that the CSP must exercise reasonable security may be difficult to implement when there is a chain of outsourcing to companies in different countries. Lastly, a CSP may be required to share with law enforcement agencies private data. For example, Microsoft was served a subpoena to provide emails exchanged by users of the Hotmail service.

The question is: What can and should cloud users do to minimize the security risks regarding the data handling by the CSP? First, a user should evaluate the security policies and the mechanisms the CSP has in place to enforce these policies. Then, the user should analyze the information that would be stored and processed on the cloud. Finally, contractual obligations should be clearly spelled out. The contract between the user and the CSP should [396]:

1. State explicitly CSP's obligations for handling securely sensitive information and to comply with privacy laws.
2. Spell out CSP liabilities for mishandling sensitive information, e.g., data loss.
3. Spell out the rules governing ownership of the data.
4. Specify the geographical regions where information and backups can be stored.

To minimize security risks, a user may try to avoid processing sensitive data on a cloud. Google's Secure Data Connector carries out an analysis of the data structures involved and allows access to data protected by a firewall. This solution is not feasible for some applications, e.g., processing of medical or personnel records; it may not be feasible when the cloud processing workflow requires cloud access to the entire volume of user data. When the volume of sensitive data or the processing workflow requires sensitive data to be stored on the cloud then, whenever feasible, data should be encrypted [192,526].

8.2 Cloud security risks

Some believe that it is very easy, possibly too easy, to start using cloud services without a proper understanding of the security risks and without the commitment to follow user responsibilities spelled

out by Service Level Agreements (SLA).¹ A first question is: What are the security risks faced by cloud users? A cloud could be used to launch large-scale attacks against other components of the cyber infrastructure, so the next question is: How can nefarious use of cloud resources be prevented?

There are multiple ways to look at the security risks for cloud computing. Three broad classes of security risks for cloud computing are: traditional security threats, threats related to system availability, and threats related to third-party data control [107].

Traditional threats are those that have been experienced for some time by any system connected to the Internet, but with some cloud-specific twists. The impact of traditional threats is amplified due to the vast amount of cloud resources and the large user population that can be affected. The long list of cloud user concerns includes also the fuzzy bounds of responsibility between the providers of cloud services and users and the difficulties to accurately identify the cause of a problem. The traditional threats begin at the user site. The user must protect the infrastructure used to connect to the cloud and to interact with the application running on the cloud. This task is more difficult because some components of this infrastructure are outside the firewall protecting the user.

The next threat is related to the authentication and authorization process. The procedures in place for one individual do not extend to an enterprise. In this case, the cloud access of the members of an organization must be nuanced; various individuals should be assigned distinct levels of privilege based on their role in the organization. It is also nontrivial to merge or adapt the internal policies and security metrics of an organization with the ones of the cloud.

Moving from the user to the cloud, we see that the traditional attacks have already affected cloud service providers. The favorite means of attack are: distributed denial of service (DDoS) attacks that prevent legitimate users to access cloud services, phishing, SQL injection, or cross-site scripting. Phishing is an attack intended to gain information from a database by masquerading as a trustworthy entity. Such information could be names and credit card numbers, social security numbers, and other personal information stored by online merchants or other service providers.

SQL injection is a form of attack typically used against a web site. In this case, an SQL command entered in a web form causes the contents of a database used by the web site to be either dumped to the attacker or altered. SQL injection can be used against other transaction processing systems, and it is successful when the user input is not strongly typed or rigorously filtered. Cross-site scripting is the most popular form of attack against web sites; a browser permits the attacker to insert client-scripts into the web pages, and thus bypass the access controls at the web site.

Identifying the path followed by an attacker is much more difficult in a cloud environment. Cloud servers host multiple VMs, and multiple applications may run under one VM. Multitenancy in conjunction with hypervisor vulnerabilities could open new attack channels for malicious users. Traditional investigation methods based on digital forensics cannot be extended to a cloud where the resources are shared among a large user population and the trace of events related to a security incident is wiped out due to the high rate of write operations on any storage media.

Availability of cloud services is another concern. System failures, power outages, and other catastrophic events could shutdown cloud services for extended periods of time. Data lock-in discussed in Section 2.6 could prevent a large organization whose business model depends on these data to function properly, when such a rare event occurs.

¹ AWS SLAs can be found at <https://aws.amazon.com/legal/service-level-agreements/>.

Clouds can also be affected by phase transition phenomena and other effects specific to complex systems. Another critical aspect of availability is that the users cannot be assured that an application hosted on the cloud returns correct results.

Third-party control generates a spectrum of concerns caused by a lack of transparency and limited user control. For example, a cloud provider may subcontract some resources from a third party whose level of trust is questionable. There are examples when subcontractors failed to maintain the customer data. There are also examples when the third party was not a subcontractor but a hardware supplier, and the loss of data was caused by poor-quality storage devices [107].

Storing proprietary data on the cloud is risky because cloud provider espionage poses real dangers. The terms of contractual obligations usually place all responsibilities for data security with the user. The Amazon Web Services customer agreement does not help user's confidence because it states "We...will not be liable to you for any direct, indirect, incidental...damages...nor...be responsible for any compensation, reimbursement, arising in connection with: (A) your inability to use the services...(B) the cost of procurement of substitute goods or services...or (D) any unauthorized access to, alteration of, or deletion, destruction, damage, loss or failure to store any of your content or other data."

It is very difficult for a cloud user to prove that data has been deleted by the service provider. The lack of transparency makes auditability a very difficult proposition for cloud computing. Auditing guidelines elaborated by the National Institute of Standards (NIST), such as the Federal Information Processing Standard (FIPS) and the Federal Information Security Management Act (FISMA), are mandatory for US Government agencies.

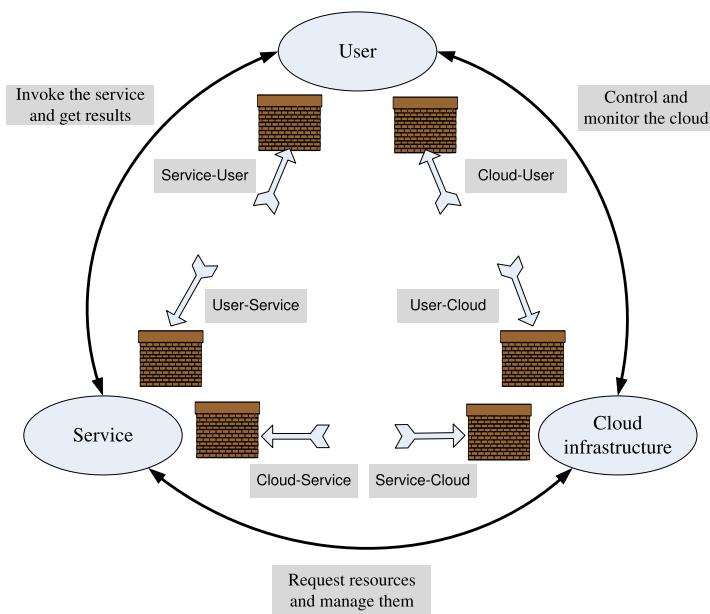
The 2010 Cloud Security Alliance (CSA) report. The report identifies seven top threats to cloud computing. These threats are: the abusive use of the cloud; APIs that are not fully secure, malicious insiders; shared technology; account hijacking; data loss or leakage; and unknown risk profile [121]. According to this report, the IaaS delivery model can be affected by all threats and PaaS can be affected by all, except the shared technology, while SaaS is affected by all, except abuse and shared technology.

Abusing the cloud refers to conducting nefarious activities from the cloud, for example, using multiple AWS instances or applications supported by IaaS to launch distributed denial of service attacks or to distribute spam and malware. *Shared technology* considers threats due to multitenant access supported by virtualization. Hypervisors can have flaws allowing a guest OS to affect the security of the platform shared with other VMs.

Insecure APIs may not protect the users during a range of activities, starting with authentication, and access control to monitoring and control of the application during runtime. The cloud service providers do not disclose their hiring standards and policies, thus the risks of *malicious insiders* cannot be ignored. The potential harm due to this particular form of attack is high.

Data loss and data leakage are two risks with devastating consequences for an individual or an organization using cloud services. Maintaining copies of the data outside the cloud is often unfeasible due to the sheer volume of data. If the only copy of the data is stored on the cloud, then sensitive data is permanently lost when cloud data replication fails, followed by a storage media failure. As some of the data often includes proprietary or sensitive data, access to such information by third parties could have severe consequences.

Account or service hijacking is a significant threat, and cloud users must be aware of and guard against all methods to steal credentials. Lastly, *unknown risk profile* refers to exposure to the ignorance or underestimation of the risks of cloud computing.

**FIGURE 8.1**

Surfaces of attacks in a cloud computing environment.

The 2011 CSA report. The report “Security Guidance for Critical Area of Focus in Cloud Computing V3.0” provides a comprehensive analysis of the risks and makes recommendations to minimize the risk in cloud computing [122].

An attempt to identify and classify the attacks in a cloud computing environment is discussed in [212]. The three actors involved in the model are: the user, the service, and the cloud infrastructure, and there are six types of attacks possible; see Fig. 8.1. The user can be attacked from two directions, the service and the cloud. Secure Sockets Layer (SSL) certificate spoofing, attacks on browser caches, or phishing attacks are example of attacks that originate at the service. The user can also be a victim of attacks that either truly originate or that spoof origination from the cloud infrastructure.

The service can be attacked by the user. Buffer overflow, SQL injection, and privilege escalation are the common types of attacks from the service. The service can also be the subject of attacks by the cloud infrastructure, and this is probably the most serious line of attack. Limiting access to resources, privilege-related attacks, data distortion, and injecting additional operations are only a few of the many possible lines of attacks originated at the cloud.

The cloud infrastructure can be attacked by a user that targets the cloud control system. These types of attacks are the same a user would direct toward any other cloud service. The cloud infrastructure may also be targeted by a service requesting an excessive amount of resources and causing the exhaustion of the resources.

Top cloud security threats. The 2016 CSA report lists top security threats [407]:

1. Data breaches. The most damaging breaches concern sensitive data, including financial and health information, trade secrets, and intellectual property. The ultimate responsibility rests with the organizations maintaining data on the cloud, and CSA recommends that organizations use multifactor authentication and encryption to protect against data breaches. Multifactor authentication, such as one-time passwords, phone-based authentication, and smart card protection, make it harder for attackers to use stolen credentials.
2. Compromised credentials and broken authentication. Such attacks are due to lax authentication, weak passwords, and poor key and/or certificate management.
3. Hacked interfaces and APIs. Cloud security and service availability can be compromised by a weak API. When third parties rely on APIs, more services and credentials are exposed.
4. Exploited system vulnerabilities. Resource sharing and multitenancy create new attack surfaces, but the cost to discover and repair vulnerabilities is small compared to the potential damage.
5. Account hijacking. All accounts should be monitored so that every transaction can be traced to the individual requesting it.
6. Malicious insiders. This threat can be difficult to detect, and system administrator errors could sometimes be falsely diagnosed as threats. A good policy is to segregate duties and enforce activities, such as logging, monitoring, and auditing administrator activities.

The other six threats are: advanced persistent threats (APTs), permanent data loss, inadequate diligence, cloud service abuse, DoS attacks, and shared technology.

An update of the “Cloud Controls Matrix” spells out the control specification impact on the cloud architecture, cloud delivery models, and other aspects of the cloud ecosystem; see <https://cloudsecurityalliance.org/download/cloud-controls-matrix-v3-0-1/>.

Cloud vulnerability incidents reported over a period of four years and data breach incidents in 2014 identified several other threats including [471]: hardware failures, natural disasters, cloud-related malware, inadequate infrastructure design and planning, point-of-sale (POS) intrusions and payment card skimmers, crimeware and cyber espionage, insider and privilege misuse, web app attacks, and physical theft/loss.

2021-global CSIO report. The report, commissioned by Dynatrace in 2021, stresses that the need to innovate faster and shift to cloud-native application architectures built on microservices, containers, and Kubernetes is driving complexity and creates vulnerability blind spots. It is becoming nearly impossible to distinguish between potential vulnerabilities and critical exposures. It is also virtually impossible for developers to conduct manual security scans or have the means to automatically identify potential vulnerabilities. Microservices are constantly changing and sometimes migrate from one cloud to another.

The report <https://assets.dynatrace.com/en/docs/report/2021-global-ciso-report.pdf> also acknowledges that “alternative approaches, such as scanning source code or container images in pre-production environments, cannot offer real-time visibility into live exposures. This means vulnerabilities often remain undetected in production, where attackers can exploit them.” The only hope for vulnerability management is future fully automated runtime security.

8.3 Security as a service (SecaaS)

SecaaS includes security services offered by a CSP. Several potential benefits of SecaaS [123] are: (i) intelligence-sharing since it protects multiple clients simultaneously and shares data intelligence and data across them; (ii) insulation of clients as it can intercept attacks before they hit an organization; and last but not least (iii) scaling and cost.

Potential concerns of SecaaS are: (i) regulation differences because the service may be unable to assure compliance in all jurisdictions an organization operates in; and (ii) lack of visibility as the service provider may not reveal details of how it implements its own security and manages its own environment. SecaaS offerings include:

1. *Identity, Entitlement, and Access Management Services* based on Federated Identity Brokers, strong authentication, multifactor authentication, etc.
2. *Cloud Access and Security Brokers*. Cloud Security Gateways connect to on-premises tools to help an organization detect, assess, and potentially block cloud usage and unapproved services.
3. *Web Security Gateways* providing an added layer of protection in addition to anti-malware software to prevent malware from entering the enterprise via activities such as web browsing.
4. *Email Security* including protection against phishing and malicious attachments, as well as enforcing corporate policies like acceptable use and spam prevention.
5. *Security Assessment* including traditional security/vulnerability assessments of assets, application security assessments, and cloud platform assessment tools.
6. *Web Application Firewalls*.
7. *Intrusion Detection/Prevention*.
8. *Security Information and Event Management*. Aggregates log and event data from virtual and real networks, applications, and systems and provides real-time reporting and alerts.
9. *Encryption and Key Management*.
10. *Business Continuity and Disaster Recovery*.
11. *Distributed Denial of Service Protection*.

CSA recommends that, before engaging a SecaaS, an organization should ensure that the service is compatible with the organization's own current and future plans, such as supported cloud platforms, the workstation, and the mobile operating systems it accommodates. Also, it is important to understand any security-specific requirements for data-handling and investigative and compliance support.

8.4 Privacy and privacy impact assessment

The term *privacy* refers to the right of an individual, a group of individuals, or an organization to keep information of a personal nature or proprietary information from being disclosed. Many nations view privacy as a basic human right. The Universal Declaration of Human Rights, article 12, states: "No 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."

U.S. Constitution contains no express right to privacy, however the Bill of Rights reflects framers concern for protecting specific aspects of privacy.² In the UK, privacy is guaranteed by the Data Protection Act. The European Court of Human Rights has developed several documents defining the right to privacy. Privacy laws differ from country to country; laws in one country may require public disclosure of information considered private in other countries and cultures.

At the same time, the right to privacy is limited by laws. For example, taxation laws require individuals to share information about personal income or earnings. Individual privacy may conflict with other basic human rights, e.g., with freedom of speech. The digital age has confronted legislators with significant challenges related to privacy as new threats have emerged. For example, personal information voluntarily shared, but stolen from sites granted access to it, or misused can lead to *identity theft*.

Some countries have been more aggressive in addressing the new privacy concerns than others. For example, European Union (EU) countries have strict laws governing handling of personal data in the digital age. A sweeping new privacy right, the “right to be forgotten,” is codified as part of a broad new proposed data-protection regulation in EU. This right addresses the following problem: It is very hard to escape your past now when every photo, status update, and tweet lives forever on some web site.

Our discussion targets primarily public clouds where privacy has an entirely new dimension because data, often in an unencrypted form, resides on servers owned by a CSP. Services based on individual preferences, location of individuals, membership in social networks, or other personal information present a special risk. The data owner cannot rely exclusively on the CSP to guarantee data privacy.

Privacy concerns are different for the three cloud delivery models and also depend on the actual context. For example, consider Gmail, a widely used SaaS delivery model; Gmail privacy policy reads (see <http://www.google.com/policies/privacy/> accessed on October 6, 2012): “We collect information in two ways: information you give us... like your name, email address, telephone number or credit card; information we get from your use of our services such as: ...device information, ...log information, ...location information, ...unique application numbers, ...local storage, ...cookies and anonymous identifiers... We will share personal information with companies, organizations or individuals outside of Google if we have a good-faith belief that access, use, preservation or disclosure of the information is reasonably necessary to: meet any applicable law, regulation, legal process or enforceable governmental request; ...protect against harm to the rights, property or safety of Google, our users or the public as required or permitted by law. We may share aggregated, non-personally identifiable information publicly and with our partners—like publishers, advertisers or connected sites. For example, we may share information publicly to show trends about the general use of our services.”

The main aspects of cloud privacy are: the lack of user control, potential unauthorized secondary use, data proliferation, and dynamic provisioning [396]. The lack of user control refers to the fact that user-centric data control is incompatible with cloud usage. Once data is stored on the servers of the CSP, the user loses control of the exact location and, in some instances, it could lose access to the data. For example, in the case of the Gmail service, the account owner has no control on where the data is stored or how long old emails are stored on some backups of the servers.

² The 1st Amendment covers the protection of beliefs, the 3rd Amendment privacy of homes, the 4th Amendment the privacy of persons and possessions against unreasonable searches, and the 5th Amendment the protection from self-incrimination. Thus, the privacy of personal information, and, according to some justices, the 9th Amendment that reads “The enumeration in the Constitution, of certain rights, shall not be construed to deny or disparage others retained by the people” can be viewed as protection of privacy in ways not explicitly specified by the first eight amendments in the Bill of Rights.

A CSP may obtain revenues from unauthorized secondary usage of the information, e.g., for targeted advertising. There are no technological means to prevent this use. Dynamic provisioning refers to threats due to outsourcing. A range of issues are very fuzzy, e.g.: how to identify the sub-contractors of a CSP, what rights to the data they have, and what rights to data are transferable in case of a bankruptcy or merger.

It is imperative for legislative bodies to address the multiple aspects of privacy in the digital age. A document elaborated by the Federal Trading Commission for the US Congress states [172]: “Consumer-oriented commercial web sites that collect personal identifying information from or about consumers online would be required to comply with the four widely-accepted fair information practices:

1. Notice—web sites should be required to provide consumers clear and conspicuous notice of their information practices, including what information they collect, how they collect it (e.g., directly or through nonobvious means such as cookies), how they use it, how they provide Choice, Access, and Security to consumers, whether they disclose the information collected to other entities, and whether other entities are collecting information through the site.
2. Choice—web sites should be required to offer consumers choices as to how their personal identifying information is used beyond the use for which the information was provided, e.g., to consummate a transaction. Such choices would encompass both internal secondary uses (such as marketing back to consumers) and external secondary uses, such as disclosing data to other entities.
3. Access—web sites would be required to offer consumers reasonable access to the information a web site has collected about them, including a reasonable opportunity to review information and to correct inaccuracies or delete information.
4. Security—web sites would be required to take reasonable steps to protect the security of the information they collect from consumers. The Commission recognizes that the implementation of these practices may vary with the nature of the information collected and the uses to which it is put, as well as with technological developments. For this reason, the Commission recommends that any legislation be phrased in general terms and be technologically neutral. Thus, the definitions of fair information practices set forth in the statute should be broad enough to provide flexibility to the implementing agency in promulgating its rules or regulations.”

There is the need for tools capable of identifying privacy issues in information systems, the so-called *Privacy Impact Assessment (PIA)*. As of mid-2017, there are no international standards for such a process, though different countries and organization require PIA reports. An example of an analysis is to assess the legal implications of the UK–US Safe Harbor process to allow US companies to comply with the European Directive 95/46/EC³ on the protection of personal data.

Such an assessment forces a proactive attitude towards privacy. An ab initio approach for embedding privacy rules in new systems is preferable to painful changes that could affect the functionality of existing systems. A PIA tool that could be deployed as web-based service proposed in [469] includes: project information, an outline of project documents, privacy risks, and stakeholders. The tool will produce a PIA report consisting of a summary of findings, a risk summary, security, transparency, and cross-borders data flows.

³ See <http://eur-lex.europa.eu/LexUriServ/LexUriServ.do?uri=CELEX:31995L0046:en:HTML>.

The centerpiece of the PIA tool is a knowledge base (KB) created and maintained by domain experts. The users of the SaaS service providing access to the PIA tool must complete a questionnaire. The system uses templates to generate additional questions necessary to complete the PIA report. An expert system infers which rules are satisfied by the facts in the database as provided by the users and executes the rule with the highest priority.

8.5 Trust

In the context of cloud computing trust is intimately related to the general problem of trust in online activities. We first discuss the traditional concept of trust and then the trust in online activities.

Trust. According to the Merriam–Webster dictionary, trust means “assured reliance on the character, ability, strength, or truth of someone or something.” Trust is a complex phenomenon; it enables cooperative behavior, promotes adaptive organizational forms, reduces harmful conflict, decreases transaction costs, facilitates formulation of ad hoc work groups, and promotes effective responses to crisis [424].

Two conditions must exist for trust to develop. The first is *risk*, the perceived probability of loss. Indeed, trust would not be necessary if there is no risk involved, if there is a certainty that an action will succeed. The second is *interdependence*, the interests of one entity cannot be archived without reliance on other entities. A trust relationship goes through three phases: (i) building phase, when trust is formed; (ii) stability phase, when trust exists; and (iii) dissolution phase, when trust declines.

There are various reasons for and forms of trust. Utilitarian reasons could be based on the belief that the costly penalties for breach of trust exceed any potential benefits from opportunistic behavior. This is the essence of *deterrence-based* trust. Another reason is the belief that the action involving the other party is in the self-interest of that party. This is the so-called *calculus-based* trust. After a long sequence of interactions *relational trust* between entities can be developed based on the accumulated experience of dependability and reliance on each other.

The common wisdom is that an entity must work very hard to build trust, but may lose the trust very easily. A single violation of trust can lead to irreparable damage. *Persistent trust* is based on the long-term behavior of an entity, while *dynamic trust* is based on a specific context, e.g., the state of the system or the effect of technological developments.

Trust in the Internet The trust in the Internet “obscures or lacks entirely the dimensions of character and personality, nature of relationship, and institutional character” of the traditional trust [365]. The missing identity, personal characteristics, and role definitions are elements we have to deal in the context of online trust.

The Internet offers individuals the ability to obscure or conceal their identities. The resulting anonymity reduces the clues normally used in judgments of trust. Identity is critical for developing trust relationships; it allows us to base our trust on the past history of interactions with an entity. Anonymity causes mistrust because identity is associated with accountability, and in the absence of identity, accountability cannot be enforced.

The opacity extends immediately from identity to personal characteristics. It is impossible to infer if the entity or individual we transact with is who it pretends to be because the transactions occur between entities separated in time and distance. Lastly, there are no guarantees that the entities we transact with fully understand the role they have assumed.

To remedy the loss of clues, we need security mechanisms for access control, transparency of identity, and surveillance. The mechanisms for access control are designed to keep out intruders and mischievous agents. Identity transparency requires that the relationship between a virtual agent and a physical person should be carefully checked security through methods such as biometric identification. Digital signatures and digital certificates are used for identification. Surveillance could be based on *intrusion detection* or can be based on logging and auditing. The first is based on real-time monitoring, while the second relies on offline sifting through audit records.

Credentials are used when the entity is not known; credentials are issued by a trusted authority and describe the qualities of the entity using the credential. A DDS (Doctor of Dental Surgery) diploma hanging on the wall of a dentist's office is a credential that the individual has been trained by an accredited university and hence capable to perform a set of procedures. A digital signature is a credential used in many distributed applications.

Policies and *reputation* are two ways of determining trust. Policies reveal the conditions to obtain trust and the actions when some of the conditions are met. Policies require the verification of credentials. Reputation is a quality attributed to an entity based on a relatively long history of interactions or possibly observations of the entity. Recommendations are based on trust decisions made by others and filtered through the perspective of the entity assessing the trust.

In a computer science context, “trust of a party A to a party B for a service X is the measurable belief of A in that B behaves dependably for a specified period within a specified context (in relation to service X),” [374]. An assurance about the operation of particular hardware or software component leads to persistent social-based trust in that component.

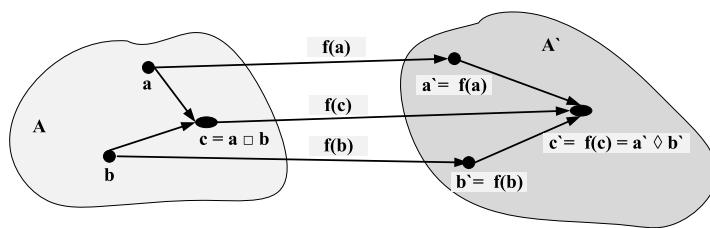
A comprehensive discussion of trust in computer services in the semantic web can be found in [33]. In Section A.1, we discuss the concept of trust in the context of cognitive radio networks where multiple transmitters compete for communication channels. Then, in Section A.2, we present a cloud-based trust management service.

8.6 Cloud data encryption

Individual users, large corporations, and the government ponder if it is safe to store sensitive information on a public cloud. Encryption is the obvious solution to protect outsourced data. The seminal RSA paper [418] and the survey of existing public-key cryptosystems in [429] are some of the notable publications in the vast literature dedicated to cryptosystems.

Cloud service providers offer encryption services. Amazon offers AWS Key Management Service (KMS) to create and control the encryption keys used by clients to encrypt their data. KMS is integrated with other AWS services including EBS, S3, RDS, Redshift, Elastic Transcoder, and WorkMail. AWS also offers Encryption SDK for developers.

Several recent research results in cryptography are important to data security in cloud computing. In 1999, Pascal Paillier proposed a trapdoor mechanism for public-key cryptosystems based on composite degree residuosity for hard-to-factor numbers $n = pq$ where p and q are large prime numbers [386]. This solution exploits the homomorphic properties of composite residuosity classes to design distributed cryptographic protocols. A major breakthrough are the algorithms for Fully Homomorphic Encryption (FHE) proposed by Craig Gentry in his seminal 2009 dissertation at Stanford University

**FIGURE 8.2**

A homomorphism $f : A \rightarrow A'$ is a structure-preserving map between sets A and A' with the composition operations \square and \diamond , respectively. Let $a, b, c \in A$ with $c = a \square b$ and $a', b', c' \in A'$ with $c' = a' \diamond b'$. Let $a' = f(a)$, $b' = f(b)$, $c' = f(c)$ be the results of the mapping $f(\cdot)$. The operation \diamond in the target domain produces the same result as the mapping the result of the operation \square applied to the two elements in the original domain, $f(a) \diamond f(b) = f(a \square b)$.

[192,193]. In recent years, searchable symmetric encryption protocols have been reported in [85] and [168].

Sensitive data is safe while in storage, provided that it is encrypted with strong encryption. But encrypted data must be decrypted for processing, and this opens a window of vulnerability. So, here is the first question examined in this section: Is it feasible to operate on encrypted data? The *homomorphic encryption*, a long-time dream of security experts, reflects the concept of homomorphism, a structure-preserving map $f(\cdot)$ between two algebraic structures of the same type; see Fig. 8.2.

When $f(\cdot)$ in Fig. 8.2 is a one-to-one mapping, the inverse of $f(\cdot)$ is $f^{-1} : A' \rightarrow A$ and $a = f^{-1}(a')$, $b = f^{-1}(b')$, $c = f^{-1}(c')$. In this case, we can carry out the operation \diamond in the target domain and apply the inverse mapping to get the same result produced by the \square operation in the original domain, $f^{-1}(a) \diamond f^{-1}(b) = f(a \square b)$. In the case of homomorphic encryption, the map is a one-to-one transformation, $f(\cdot)$ is the encryption procedure, its inverse, $f^{-1}(\cdot)$, is the decryption procedure, and the composition operation can be any arithmetic and logic operation. The homomorphism shows that we can conduct arithmetic and/or logic operations with encrypted data and the decryption of the result of these operation replicates the result of carrying out the same operations with the plaintext data. This implies that the window of vulnerability created when data is decrypted for processing disappears.

General computations with encrypted data are theoretically feasible using FHE algorithms. Unfortunately, the homomorphic encryption is not a practical solution at this time. Existing algorithms for homomorphic encryption increase the processing time with encrypted data by many orders of magnitude compared with the processing of plaintext data. A recent implementation of FHE [223] requires about six minutes per batch; the processing time for a simple operation on encrypted data dropped to almost one second after improvements in other experiments [153].

Users send a variety of queries to many large databases stored on clouds. Such queries often involve logic and arithmetic functions, so an important question is if it is feasible and practical to search encrypted databases. Application of widely used encryption techniques to database systems could lead to significant performance degradation. For example, if an entire column of a NoSQL database table contains sensitive information and is encrypted, then a query predicate with a comparison operator requires a scan of the entire table to evaluate the query. This is due to the fact that existing encryption algorithms do not preserve order, and database indices such as B-tree can no longer be used.

Order Preserving Encryption (OPE) [66] can be used for encryption of numeric data. OPE maps a range of numerical values into a much larger and sparse range of values. Let a order-preserving function $f : \{1, \dots, M\} \rightarrow \{1, \dots, N\}$ with $N \gg M$ be uniquely represented by a combination of M out of N ordered items. Given N balls in a bin, M black and $N - M$ white, we draw a ball at random without replacement at each step. The random variable X describing the total number of balls in our sample after we collect the k -th black ball follows the negative hypergeometric distribution. One can show that an order preserving $f(x)$ for a given point $x \in \{1, \dots, M\}$ has an NHG distribution over a random choice of f .

To encrypt plaintext x , the OPE encryption algorithm performs a binary search down to x . Given the secret key K , the algorithm first assigns $\text{Encrypt}(K, M/2)$, then $\text{Encrypt}(K, M/4)$ if the index $m < M/2$ and $\text{Encrypt}(K, 3M/4)$ otherwise, and so on, until $\text{Encrypt}(K, x)$ is assigned. Each ciphertext assignment is made according to the output of the negative hypergeometric sampling algorithm. One can prove by strong induction on the size of the plaintext space that the resulting scheme induces a random order-preserving function from the plaintext to ciphertext space.

It is sufficient to have an order-preserving hash function H (not necessarily invertible) to allow efficient range queries on encrypted data. The algorithm would use a secret key $(K_{\text{Encrypt}}, K_H)$, where K_{Encrypt} is a key for a normal (randomized) encryption scheme and K_H is a key for H . Then $\text{Encrypt}(K_{\text{Encrypt}}, x) \parallel H(K_H, x)$ will be the encryption of x [66].

Searching encrypted databases is of particular interest [13]. Several types of searches are frequently conducted including: single-keyword, multikeyword, fuzzy-keyword, ranked, authorized, and verifiable search. Searchable symmetric encryption (SSE) is used when an encrypted databases \mathcal{E} is outsourced to a cloud or to a different organization. SSE hides information about the database and the queries.

A client only stores the cryptographic key. To search a database, the client encrypts the query, sends it to the database server, receives the encrypted result of the query, and decrypts it using the cryptographic key. Information leakage from these searches is confined to query patterns, while disclosure of explicit data and query plaintext values is prevented.

An SSE protocol supporting conjunctive search and general Boolean queries on symmetrically encrypted data was proposed in [85]. This SSE protocol scales to very large databases. It can be used for arbitrarily structured data including free text search with the moderate and well-defined leakage to the outsourced server. The performance results of a prototype applied to encrypted search over the entire English Wikipedia are reported. The protocol was extended with support for range, substring, wildcard, and phrase queries [168].

The next question is if sensitive data stored on the servers of a private cloud is vulnerable. The threat posed by an outsider attacker is diminished if the private cloud is protected by an effective firewall. Nevertheless, there are dangers posed by an insider. If such an attacker has access to log files, it can infer the location of database hot spots, copy data selectively, and use the data for a nefarious activity. To minimize the risks posed by an insider, a set of protections rings should be enforced to restrict the access of each member of the staff to a limited area of the data base.

8.7 Security of database services

Cloud users who delegate the control of their data to the database services supported by virtually all CSP are concerned with Data Base as a Service (DBaaS) security aspects. The model used to evaluate

DBaaS security includes several group of entities: the data owners, the users of data, the cloud service providers, and third-party agents or Third Party Auditors (TPA).

Data owners and DBaaS users fear compromised integrity and confidentiality, as well as data unavailability. Insufficient authorization, authentication and accounting mechanisms, inconsistent use of encryption keys and techniques, alteration or deletion of records without maintaining backup, and operational failures are the major causes of data loss in DBaaS.

Some data integrity and privacy issues are due to the absence of authentication, authorization and accounting controls, or poor key management for encryption and decryption. Confidentiality means that only authorized users should have access to the data. Unencrypted data is vulnerable to bugs, errors, and attacks from external entities affecting data confidentiality.

Insider attacks are a concern for DBaaS users and data owners. Superusers have unlimited privileges, and misuse of superuser privileges poses a considerable threat to confidential data, such as medical records, sensitive business data, proprietary product data, etc. Malicious external attackers use spoofing, sniffing, man-in-the-middle attacks, side channeling, and illegal transactions to launch DoS attacks.

Another concern is illegal data recovery from storage devices, a side effect of multi-tenancy. CSP often carry out sanitation operations after deleting data from physical devices, but, unless a thorough scrubbing operation is carried out, sophisticated attackers can still recover information from storage devices. Data is also vulnerable during transfer from the data owner to the DBaaS through public networks. Encryption before data transmission can reduce the risks posed to the data in transit to the cloud.

Data provenance, the process of establishing the origin of data and its movement between databases, uses metadata to determine the data accuracy, but the security assessments are time-sensitive. Moreover, analyzing large provenance metadata graphs is computationally expensive.

Cloud users are not aware of the physical location of their data. This lack of transparency allows cloud service providers to optimize the use of resources, but in the case of security breaches, it is next to impossible for users to identify compromised resources. DBaaS users do not have fine-grained control of the remote execution environment and cannot inspect the execution traces to detect the occurrence of illegal operations.

To increase availability, performance, and to enhance reliability, DBaaS replicate data. Ensuring consistency among the replicas is challenging. Another critical function of DBaaS is to carry out timely backups of all sensitive and confidential data to facilitate quick recovery in case of disasters. Auditing and monitoring are important functions of a DBaaS but generate their own security risks when delegated to TPAs. Conventional methods for auditing and monitoring demand detailed knowledge of the network infrastructure and physical devices. Data privacy laws can be violated because consumers are unaware where the data is actually stored. Privacy laws in Europe and South America prohibit storing data outside the country of origin.

In summary, DBaaS data *availability* is affected by several threats including: (i) resource exhaustion caused by imprecise specification of user needs or incorrect evaluation of user specifications; and (ii) failures of the consistency management; multiple hardware and/or software failures lead to inconsistent views of user data; and failure of the monitoring and auditing system. DBaaS data *confidentiality* is affected by insider and outsider attacks, access control issues, illegal data recovery from storage, network breaches, third-party access, and inability to establish the provenance of the data.

8.8 Operating system security

An operating system allows multiple applications to share the hardware resources of a physical system subject to a set of policies. A critical function of an OS is to protect applications against a wide range of malicious attacks such as unauthorized access to privileged information, tampering with executable code, and spoofing. Such attacks can now target even single-user systems such as personal computers, tablets, or smartphones. Data brought into the system may contain malicious code; this could be the case of a Java applet, or of data imported by a browser from a malicious web site.

The *mandatory security* of an OS is considered to be [296]: “any security policy where the definition of the policy logic and the assignment of security attributes is tightly controlled by a system security policy administrator.” Access control, authentication usage, and cryptographic usage policies are all elements of the mandatory OS security.

An access control policy specifies how the OS controls the access to various system objects; authentication usage defines the authentication mechanisms used by the OS to authenticate a principal; and cryptographic usage policies specify the cryptographic mechanisms used to protect the data. A necessary but not sufficient condition for security is that the subsystems tasked to perform security-related functions are tamper-proof and cannot be bypassed. The OS should confine an application to a unique security domain.

Applications with special privileges that perform security-related functions are called *trusted applications*. Such applications should only be allowed the lowest level of privileges required to perform their functions. For example, type enforcement is a mandatory security mechanism that can be used to restrict a trusted application to the lowest level of privileges.

Enforcing mandatory security through mechanisms left to the user’s discretion can lead to a breach of security, sometimes due to malicious intent or in other cases due to carelessness or lack of understanding. Discretionary mechanisms place the burden of security on individual users. Moreover, an application may change a carefully defined discretionary policy without the consent of the user, while a mandatory policy can only be changed by a system administrator.

Unfortunately, commercial operating systems do not support multilayered security. They only distinguish between a completely privileged security domain and a completely unprivileged one. Some operating systems, e.g., Windows NT, allow a program to inherit all the privileges of the program invoking it, regardless of the level of trust in that program.

The existence of *trusted paths*, mechanisms supporting user interactions with trusted software, is critical for system security. Malicious software can impersonate trusted software when such mechanisms do not exist. Some systems allow servers to authenticate their clients and provide trusted paths for login authentication and password changing.

The solution discussed in [296] is to decompose a complex mechanism in several components with well-defined roles. For example, the access control mechanism for the application space could consist of *enforcer* and *decider* components. To access a protected object, the enforcer will gather the required information about the agent attempting the access and will pass this information to the decider together with the information about the object and the elements of the policy decision; finally, it will carry out the actions requested by the decider.

A trusted-path mechanism is required to prevent malicious software invoked by an authorized application to tamper with the attributes of the object and/or with the policy rules. A trusted path is also required to prevent an impostor from impersonating the decider agent. A similar solution is proposed

for the cryptography usage that should be decomposed into an analysis of the invocation mechanisms and an analysis of the cryptographic mechanism.

Another question is how an OS can protect itself and the applications running under it from malicious mobile code attempting to gain access to the data and the other resources and compromise system confidentiality and/or integrity. Java Security Manager uses the type-safety attributes of Java to prevent unauthorized actions of an application running in a “sandbox.” Yet, the Java Virtual Machine (JVM) accepts byte code in violation of language semantics; moreover, it cannot protect itself from tampering from other applications.

Even if all these security problems could be eliminated, the security relies on the ability of the file system to preserve the integrity of Java class code. The approach to requiring digitally signed applets and accepting them only from trusted sources could fail due to the all-or-nothing security model. A solution to securing mobile communications could be confining a browser to a distinct security domain.

Specialized *closed-box platforms*, such as the ones on some cellular phones, game consoles and Automatic Teller Machines (ATMs), could have embedded cryptographic keys that allow themselves to reveal their true identity to remote systems and authenticate the software running on them. Such facilities are not available to *open-box platforms*, the traditional hardware designed for commodity operating systems.

A highly secure operating system is necessary but not sufficient. Application-specific security is also necessary. Sometimes, security implemented above the operating system is better, e.g., electronic commerce requires a digital signature on each transaction.

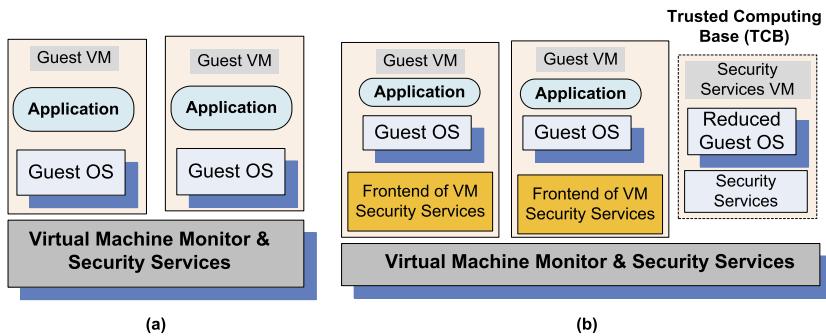
We conclude that commodity operating systems offer low assurance. Indeed, an OS is a complex software system consisting of millions of lines of code, and it is vulnerable to a wide range of malicious attacks. An OS poorly isolates one application from another; once an application has been compromised, the entire physical platform and all applications running on it can be affected. Thus, the platform security level is reduced to the security level of the most vulnerable application running on the platform.

Operating systems provide only weak mechanisms for applications to authenticate one another and do not have a trusted path between users and applications. These shortcomings add to the challenges of providing security in a distributed computing environment. For example, a financial application cannot determine if a request comes from an authorized user or from a malicious program; in turn, a human user cannot distinguish a response from a malicious program impersonating the service from the response provided by the service.

8.9 Virtual machine security

Our discussion of VM security is restricted to the traditional system VM model in Fig. 5.1(b) when the hypervisor controls the access to the hardware. The hybrid and the hosted VM models shown in Figs. 5.1(c) and (d), respectively, expose the entire system to the vulnerability of the host operating system, thus we will not analyze these models.

Virtual security services are typically provided by the hypervisor as shown in Fig. 8.3(a); another alternative is to have a dedicated VM providing security service as in Fig. 8.3(b). A secure TCB (Trusted Computing Base) is a necessary condition for security in a VM environment. When TCB is compromised, the security of the entire system is affected.

**FIGURE 8.3**

(a) Virtual security services provided by the hypervisor/Virtual Machine Monitor; (b) A dedicated security VM.

The analysis of Xen and *vBlades* in Sections 5.8 and 5.12 shows that the VM technology provides a stricter isolation of VMs from one another than the isolation of processes in a traditional operating system. Indeed, a hypervisor controls the execution of privileged operations and can thus enforce memory isolation, as well as disk and network access.

Hypervisors are less complex and better structured than traditional operating systems and in a better position to respond to security attacks. A major challenge is that a hypervisor sees only raw data regarding the state of a guest OS, while security services typically operate at a higher logical level, e.g., at the level of a file rather than a disk block.

A guest OS runs on simulated hardware, and the hypervisor has access to the state of all VMs operating on the same hardware. The state of a guest VM can be saved, restored, cloned, and encrypted by the hypervisor. Replication can ensure not only reliability but also support security, while cloning could be used to recognize a malicious application by testing it on a cloned system and observing if it behaves normally.

We can also clone a running system and examine the effect of potentially dangerous applications. Another interesting possibility is to have the guest VM's files moved to a dedicated VM and, thus, protect it from attacks [541]. This solution is possible because inter-VM communication is faster than communication between two physical machines.

Sophisticated attackers are able to fingerprint VMs and avoid VM honey pots designed to study the methods of attack. They can also attempt to access VM-logging files and, thus, recover sensitive data; such files have to be very carefully protected to prevent unauthorized access to cryptographic keys and other sensitive data.

We expect to pay some price for the better security provided by virtualization. This price includes: higher hardware costs because a virtual system requires more resources such as CPU cycles, memory, disk, and network bandwidth; the cost of developing hypervisors and modifying the host operating systems in the case of paravirtualization; and the overhead of virtualization since the hypervisor is involved in privileged operations.

VM-based intrusion detection systems such as Livewire and Siren that exploit the three capabilities of a VM for intrusion detection, isolation, inspections, and interposition are surveyed in [541]. Resource isolation was examined in Section 4.11. Inspection means that the hypervisor has the ability to review

the state of the guest VMs, and interposition means that hypervisor can trap and emulate the privileged instruction issued by the guest VMs. [541] also discusses VM-based intrusion prevention systems, such as SVFS, NetTop, and IntroVirt.

NIST security group distinguishes two groups of threats, hypervisor-based and VM-based. There are several types of hypervisor-based threats:

- a.** Resources starvation and denial of service for some VMs. Probable causes: (1) badly configured resource limits for some VMs; (2) a rogue VM with the capability of bypassing resource limits set in hypervisor.
- b.** VM side-channel attacks: malicious attack on one or more VMs by a rogue VM under the same hypervisor. Probable causes: (1) lack of proper isolation of inter-VM traffic due to misconfiguration of the virtual network residing in the hypervisor; (2) limitation of packet-inspection devices to handle high-speed traffic, e.g., video traffic; (3) presence of VM instances built from insecure VM images, e.g., a VM image having a guest OS without the latest patches.
- c.** Buffer overflow attacks.

There are two types of VM-based threats: (1) deployment of rogue or insecure VM; unauthorized users may create insecure instances from images or may perform unauthorized administrative actions on existing VMs. Probable cause: improper configuration of access controls on VM administrative tasks such as instance creation, launching, suspension, re-activation, etc.; (2) Presence of insecure and tampered VM images in the VM image repository. Probable causes: (i) lack of access control to the VM image repository; (ii) lack of mechanisms to verify the integrity of the images, e.g., digitally signed image.

8.10 Security of virtualization

The complex relationship between virtualization and security has two distinct aspects: virtualization of security and security of virtualization [303]. In Chapter 5, we praised the virtues of virtualization. We also discussed two problems associated with virtual environments: (a) the negative effect on performance due to the additional overhead; and (b) the need for more powerful systems to run multiple VMs. In this section, we take a closer look at the security of virtualization.

The complete state of an operating system running under a VM is captured by the VM. The VM state can be saved in a file, and then the file can be copied and shared. There are several useful implications of this important virtue of virtualization:

1. Supports the IaaS delivery model. An IaaS user selects an image matching the local application environment and then uploads and runs the application on the cloud using this image.
2. Increased reliability. An operating system with all the applications running under it can be replicated and switched to a hot standby in case of a system failure. Recall that a hot standby is a method to achieve redundancy. The primary and the backup systems run simultaneously and have identical state information.
3. Straightforward mechanisms for implementing resource management policies. An OS and the applications running under it can be moved to another server to balance the load of a system. For

- example, the load of lightly loaded servers can be moved to other servers and then lightly loaded servers can be switched off or placed in standby mode to reduce power consumption.
- 4. Improved intrusion detection. In a virtual environment a clone can look for known patterns in system activity and detect intrusion. The operator can switch a server to hot standby when suspicious events are detected.
 - 5. Secure logging and intrusion protection. When implemented at the OS level, intrusion detection can be disabled and logging can be modified by an intruder. When implemented at the hypervisor layer, the services cannot be disabled or modified. In addition, the hypervisor may be able to log only events of interest for a post-attack analysis.
 - 6. More efficient and flexible software maintenance and testing. Virtualization allows the multitude of OS instances to share a small number of physical systems, instead of a large number of dedicated systems running under different operating systems, different versions of each OS, and different patches for each version.

Is there a price to pay for the benefits of virtualization? There is always the other side of a coin, so we should not be surprised that the answer to this question is a resounding Yes. In a 2005 paper [189], Garfinkel and Rosenblum argue that the serious implications of virtualization on system security cannot be ignored. This theme is revisited in 2008 by Price [404], who reached similar conclusions.

A first type of undesirable effects of virtualization lead to a diminished ability of an organization to manage its systems and track their status:

- i. The number of physical systems in the inventory of an organization is limited by cost, space, energy consumption, and human support. The explosion of the number of VMs is a fact of life; to create a VM, one simply copies a file. The only limitation for the number of VMs is the amount of storage space available.
- ii. In addition to a quantity, there is also a qualitative side to the explosion of the number of VMs. Traditionally, organizations install and maintain the same version of system software. In a virtual environment such a uniformity cannot be enforced, and the number of different operating systems, their versions, and the patch status of each version will be very diverse. This diversity will tax the support team.
- iii. One of the most critical problems posed by virtualization is related to the software lifecycle. The traditional assumption is that the software lifecycle is a straight line, hence the patch management is based on a monotonic forward progress. The virtual execution model *maps to a tree structure* rather than a line. Indeed, at any point in time, multiple instances of the VM can be created and then each one of them can be updated, different patches installed, and so on. This problem has serious implications for security, as we shall soon see.

What are the direct implications of virtualization on security? A first question is: How can the support team deal with the consequences of an attack in a virtual environment? Also, do we expect the infection with a computer virus or a worm to be less manageable in a virtual environment? The surprising answer is that an infection may last indefinitely.

Infected VMs may be dormant at the time when the measures to clean up the systems are taken. Then, at a later time, the infected VMs wake up and infect other systems. The scenario can repeat itself and guarantee that infection will last indefinitely. This is in stark contrast with the manner an infection is treated in non-virtual environments. Once an infection is detected, the infected systems are quarantined and then cleaned up; the systems will then behave normally until the next infection occurs.

A more general observation is that in a traditional computing environment a steady state can be reached. In this steady state all systems are brought up to a “desirable” state, whereas “undesirable” states, states when some of the systems are either infected by a virus or display an undesirable pattern of behavior, are only transient. The desirable state is reached by installing the latest system software version and then applying the latest patches to all systems.

A virtual environment may never reach a steady state due to the lack of control. In a nonvirtual environment, the security can be compromised when an infected laptop is connected to the network protected by a firewall, or when a virus is brought in on a removable media. But, unlike a virtual environment, the system can still reach a steady state.

A side effect of the ability to record the complete state of a VM in a file is the possibility to roll back a VM. This opens wide the door for a new type of vulnerability caused by events recorded in the memory of an attacker. Two such situations are discussed in [189].

The first is that one-time passwords are transmitted in the clear, and the protection is guaranteed only if the attacker does not have the possibility to access passwords used in previous sessions. An attacker can replay rolled-back versions and access past sniffed passwords if a system runs the S/KEY password system. S/KEY is a password system based on Leslie Lamport’s scheme. It is used by several operating systems including, Linux, OpenBSD, and NetBSD. The real password of the user is combined with a short set of characters and a counter that is decremented at each use to form a single-use password.

The second is related to the requirement of some cryptographic protocols and even non-cryptographic protocols regarding the “freshness” of the random number source used for session keys and nonces. This occurs when a VM is rolled back to a state when a random number has been generated but not yet used. A nonce is a random or pseudorandom number issued in an authentication protocol to ensure that old communications cannot be reused in replay attacks. For example, nonces are used to calculate an *MD5* of the password for *HTTP* digest access authentication. Each time the authentication challenge response code is presented, the nonces are different, thus replay attacks are virtually impossible. This guarantees that an online order to Amazon or other online store cannot be replayed.

Even noncryptographic use of random numbers may be affected by the roll-back scenario. For example, the initial sequence number for a new TCP connection must be “fresh.” The door to TCP hijacking is left open when the initial sequence number is not fresh.

Another undesirable effect of virtual environment affects trust. Recall from Section 8.5 that *trust is conditioned by the ability to guarantee the identity of entities involved*. Each computer system in a network has a unique physical, or MAC address. The uniqueness of the MAC address guarantees that an infected or a malicious system can be identified and then shut down, or denied network access. This process breaks down for virtual systems when VMs are created dynamically. Often, a random MAC address is assigned to a newly created VM to avoid name collision. The other effect discussed at length in Section 8.11 is that popular VM images are shared by many users.

The ability to guarantee confidentiality of sensitive data is yet another pillar of security affected by virtualization. Virtualization undermines the basic principle that time-sensitive data stored on any system should be reduced to a minimum. First, the owner has very limited control over where sensitive data is stored: It could be spread across many servers and may be left on some of them indefinitely. A hypervisor records the state of a VM to be able to roll it back; this process allows an attacker to access sensitive data the owner attempted to destroy.

8.11 Security risks posed by shared images

Even if we assume that a cloud service provider is trustworthy, there are other sources of concern many users either ignore, or they underestimate the danger they pose. One of them, especially critical for the IaaS cloud delivery model, is image sharing. For example, a user of AWS has the option to choose between Amazon Machine Images (AMIs) accessible through the Quick Start or the Community AMI menus of the EC2 service. The option of using one of these AMIs is especially tempting for a first-time user, or for a less sophisticated one.

First, we review the process to create an AMI. We can start from a running system, from another AMI, or from the image of a VM and copy the contents of the file system to the S3, so-called *bundling*. The first of the three steps of bundling is to create an image, the second step is to compress and encrypt the image, and the last step is to split the image into several segments and then upload the segments to the S3 storage.

Two procedures *ec2-bundle-image* and *ec2-bundle-volume* are used for the creation of an AMI. The first is used for images prepared as loopback files⁴ when the data is transferred to the image in blocks. To bundle a running system, the creator of the image can use the second procedure when bundling works at the level of the file system and files are copied recursively to the image.

To use an image, a user has to specify the resources, provide the credentials for login, a firewall configuration, and specify the region, as discussed in Section 2.2. Once the image is instantiated, the user is informed about the public DNS, and the VM is available. A Linux system can be accessed using ssh at port 22, while the Remote Desktop at port 3389 is used for Windows.

The results of an analysis carried over a period of several months, from November 2010 to May 2011, of over 5 000 AMIs available through the public catalog at Amazon are reported in [44]. Many images analyzed allowed a user to *undelete* files, recover credentials, private keys, or other types of sensitive information with little effort and using standard tools. The results of this study were shared with the Amazon's Security Team, which acted promptly to reduce the threats posed to AWS users.

The details of the testing methodology can be found in [44], but here, we only discuss the results. The study was able to audit some 5 303 images out of the 8 448 Linux AMIs and 1 202 Windows AMIs at Amazon sites in the US, Europe, and Asia. The audit covered software vulnerabilities and security and privacy risks.

The average duration of an audit was 77 minutes for a Windows image and 21 minutes for a Linux image; the average disk space used was about 1 GB and 2.7 GB, respectively. The entire file system of a Windows AMI was audited because most of the malware targets Windows systems. Only directories containing executables for Linux AMIs were scanned; this strategy and the considerably longer start-up time of Windows explain the time discrepancy of the audits for the types of AMIs.

A *software vulnerability* audit revealed that 98% of the Windows AMIs (249 out of 253) and 58% (2 005 out of 3 432) Linux AMIs audited had critical vulnerabilities. The average number of vulnerabilities per AMI were 46 for Windows and 11 for Linux AMIs. Some of the images were rather old; 145,

⁴ A *loopback file system* (LOFS) is a virtual file system providing an alternate path to an existing file system. When other file systems are mounted onto an LOFS file system, the original file system does not change. One useful purpose of LOFS is to take a CDROM image file, a file of type “.iso”, mount it on the file system and then access it without the need to record a CD-R. It is somewhat equivalent to the Linux *mount -o loop* option but adds a level of abstraction; most commands that apply to a device can be used to handle the mapped file.

38, and 2 Windows AMIs and 1197, 364, and 106 Linux AMIs were older than two, three, and four years, respectively. The tool used to detect vulnerabilities, the Nessus system, available from <http://www.tenable.com/productus/nessus>, classifies the vulnerabilities based on the severity in four groups, at levels zero to three. The audit reported only vulnerabilities of the highest severity level, e.g., remote code execution.

Three types of *security risks* are analyzed: (1) backdoors and leftover credentials, (2) unsolicited connections, and (3) malware. An astounding finding, about 22% of scanned Linux AMIs contained credentials allowing an intruder to remotely login to the system. Some 100 passwords, 995 ssh keys, and 90 cases when both could be retrieved were identified.

To rent a Linux AMI, a user must provide the public part of her ssh key, and this key is stored in the *authorized_keys* in the home directory. This opens a backdoor for a malicious creator of an AMI who does not remove her own public key from the image and can remotely login to any instance of this AMI. Another backdoor is opened when the ssh server allows password-based authentication and the malicious creator of an AMI does not remove her own password. This backdoor is open even wider since one can extract the password hashes and then crack the passwords using a tool such as John the Ripper, see <http://www.openwall.com/john>.

Another threat is posed by the omission of the *cloud-init* script that should be invoked when the image is booted. This script provided by Amazon regenerates the host key an ssh server uses to identify itself; the public part of this key is used to authenticate the server. When this key is shared among several systems, these systems become vulnerable to man-in-the middle attacks.

An attacker impersonates the agents at both ends of a communication channel in the *man-in-the-middle* attack and makes them believe that they are communicating through a secure channel. For example, if B sends her public key to A, but C is able to intercept it, such an attack proceeds as follows: C sends a forged message to A claiming to be from B, but instead includes C's public key. Then A encrypts her message with C's key, believing that she is using B's key, and sends the encrypted message to B. The intruder, C, intercepts, deciphers the message using her private key, possibly alters the message, and re-encrypts with the public key B originally sent to A. When B receives the newly encrypted message, she believes it came from A.

When this script does not run, an attacker can use the *NMap* tool⁵ to match the *ssh* keys discovered in the AMI images with the keys obtained with *NMap*. The study reports that the authors were able to identify more than 2 100 instances following this procedure.

Unsolicited connections pose a serious threat to a system. Outgoing connections allow an outside entity to receive privileged information, e.g., the IP address of an instance and events recorded by a *syslog* daemon to files in the *var/log* directory of a Linux system. Such information is available only to users with administrative privileges.

The audit detected two Linux instances with modified *syslog* daemon that forwarded to an outside agent information about events such as login and incoming requests to a web server. Some of the unsolicited connections are legitimate, for example, connections to a software update site. It is next to impossible to distinguish legitimate from malicious connections.

⁵ *NMap* is a security tool running on most operating systems including *Linux*, *Microsoft Windows*, *Solaris*, *HP-UX*, *SGI-IRIX*, and *BSD* variants, such as *Mac OS X* to map the network. Mapping the network means to discover hosts and services in a network.

Malware including viruses, worms, spyware, and Trojans were identified using *ClamAV*, a software tool with a database of some 850 000 malware signatures, available from <http://www.clamav.net>. Two infected Windows AMIs were discovered, one with a *Trojan-Spy* (variant 50112) and a second one with a *Trojan-Agent* (variant 173287). The first Trojan carries out key logging and allows stealing data from the files system and monitoring processes; the AMI also included a tool to decrypt and recover passwords stored by the *Firefox* browser, called *Trojan.Firepass*.

The creator of a shared AMI assumes some *privacy risks*. Her private keys, IP addresses, browser history, shell history, and deleted files can be recovered from the published images. A malicious agent can recover the AWS API keys that are not password protected. Then the malicious agent can start AMIs and run cloud applications at no cost to herself because the computing charges are passed on to the owner of the API key. The search can target files with names such as *pk-[0-9A-Z]*.pem* or *cert-[0-9A-Z]*.pem* used to store API keys.

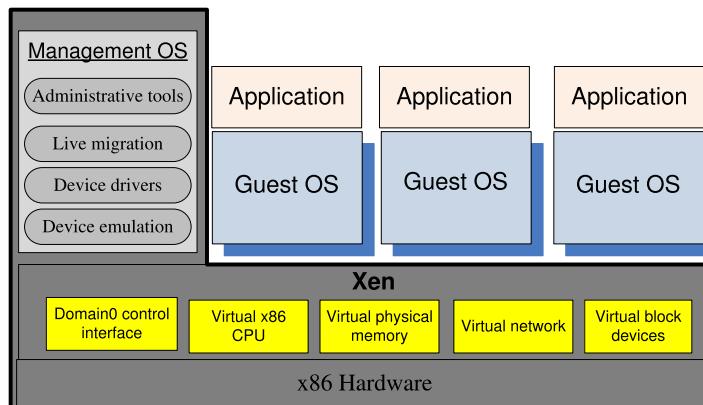
Another avenue for a malicious agent is to recover ssh keys stored in files named *id_dsa* and *id_rsa*. Though ssh keys can be protected by a passphrase, the audit determined that the majority of ssh keys (54 out of 56) were not password protected. A *passphrase* is a sequence of words used to control access to a computer system. A passphrase is the analog of a password but provides added security. For high-security nonmilitary applications, NIST recommends an 80-bit strength passphrase. Therefore, a secure passphrase should consist of at least 58 characters including uppercase and alphanumeric characters. The entropy of written English is less than 1.1 bits per character.

Recovery of IP addresses of other systems owned by the same user requires access to the *lastlog* or the *lastb* databases. The audit found 187 AMIs with a total of more than 66 000 entries in their *lastb* databases. Nine AMIs contained Firefox browser history and allowed the auditor identify the domains contacted by the user.

Six hundred and twelve AMIs contained at least one shell history file. The audit analyzed 869 history files named *~/.history*, *~/.bash_history*, and *~/.sh_history* containing some 160 000 lines of command history and identified 74 identification credentials. The users should be aware that, when the HTTP protocol is used to transfer information from a user to a web site, the GET requests are stored in the logs of the web server. Passwords and credit card numbers communicated via a GET request can be exploited by a malicious agent with access to such logs. When remote credentials, such as the DNS management password, are available, then a malicious agent can redirect traffic from its original destination to her own system.

Recovery of deleted files containing sensitive information poses another risk. When the sectors on the disk containing sensitive information are actually overwritten by another file, recovery of sensitive information is much harder. To be safe, the creator of the image effort should use utilities, such as *shred*, *scrub*, *zerofree* or *wipe*, to make recovery of sensitive information next to impossible. If the image is created with the block-level tool discussed at the beginning of this section, the image will contain blocks of the file system marked as free; such block may contain information from deleted files. The audit process was able to recover files from 98% of the AMIs using the *exundelete* utility. The number of files recovered from an AMI were as low as 6 and as high as 40 000.

We conclude that the users of published AMIs, as well as the providers of images, may be vulnerable to a wide range of security risks and must be fully aware of the dangers posed by image sharing.

**FIGURE 8.4**

The trusted computing base of a Xen-based environment includes the hardware, Xen, and the management operating system running in *Dom0*. The management OS supports administrative tools, live migration, device drivers, and device emulators. A guest OS and applications running under it reside in a *DomU*.

8.12 Security risks posed by a management OS

We often hear that virtualization enhances security because a VM hypervisor is considerably smaller than an operating system. For example, the Xen hypervisor discussed in Section 5.8 has approximately 60 000 lines of code, one to two orders of magnitude fewer than a traditional operating system.⁶

A hypervisor supports a stronger isolation between the VMs running under it than the isolation between processes supported by a traditional operating system. Yet the hypervisor must rely on a management OS to create VMs and to transfer data in and out from a guest VM to storage devices and network interfaces.

A small hypervisor can be carefully analyzed, thus one could conclude that the security risks in a virtual environment are diminished. We have to be cautious with such sweeping statements. Indeed, the Trusted Computer Base (TCB)⁷ of a cloud computing environment includes not only the hypervisor but also the management OS. The management OS supports administrative tools, live migration, device drivers, and device emulators.

For example, the TCB of an environment based on Xen includes not only the hardware and the hypervisor, but also the management operating system running in *Dom0*; see Fig. 8.4. System vulnerabilities can be introduced by both Xen and the management operating system. An analysis of Xen vulnerabilities reports that 21 of 23 attacks were against service components of the control VM [114];

⁶ The number of lines of code of the *Linux* operating system evolved over time from 176 250 for *Linux 1.0.0*, released in March 1995 to 1 800 847 for *Linux 2.2.0*, released in January 1999, 3 377 902 for *Linux 2.4.0*, released in January 2001, and to 5 929 913 for *Linux 2.6.0*, released in December 2003.

⁷ The TCB is defined as the totality of protection mechanisms within a computer system, including hardware, firmware, and software. The combination of all these elements is responsible for enforcing a security policy.

11 attacks were attributed to problems in the guest OS caused by buffer overflow and 8 were denial of service attacks. Buffer overflow allows execution of arbitrary code in the privileged mode.

Dom0 manages building of user domains (*DomU*), a process consisting of several steps:

- a. Allocate memory in *Dom0* address space and load the kernel of guest OS from secondary storage.
- b. Allocate memory and use foreign mapping to load the kernel to the new VM. The foreign mapping mechanism of Xen is used by *Dom0* to map arbitrary memory frames of a VM into its page tables.
- c. Set up the initial page tables for the new VM.
- d. Release the foreign mapping on the new VM memory, set up the virtual CPU registers, and launch the new VM.

A malicious *Dom0* can play several nasty tricks at the time that it creates a *DomU*:

- i. Refuse to carry out the steps necessary to start the new VM, an action that can be considered a *denial-of-service* attack.
- ii. Modify the kernel of the guest OS in ways that will allow a third party to monitor and control the execution of applications running under the new VM.
- iii. Undermine the integrity of the new VM by setting the wrong page tables and/or setup wrong virtual CPU registers.
- iv. Refuse to release the foreign mapping, and access the memory while the new VM is running.

We now discuss the run-time interaction between *Dom0* and a *DomU*. Recall that *Dom0* exposes a set of abstract devices to the guest operating systems using *split drivers*; the front end of such a driver is in *DomU*, its back end in *Dom0*, and the two communicate via a ring in shared memory; see Section 5.8.

In the original implementation of Xen, a service running in a *DomU* sends data to, or receives data from, a client located outside the cloud using a network interface in *Dom0*; it transfers the data to I/O devices using a device driver in *Dom0*. Later implementations of Xen offer the pass-through option.

We have to ensure that run-time communication through *Dom0* is encrypted. Yet, Transport Layer Security (TLS) does not guarantee that *Dom0* cannot extract cryptographic keys from the memory of the OS and applications running in *DomU*. A significant security weakness of *Dom0* is that the entire state of the system is maintained by XenStore; see Section 5.8. A malicious VM can deny access to this critical element of the system to other VMs; it can also gain access to the memory of a *DomU*. This brings us to additional requirements for confidentiality and integrity imposed on *Dom0*.

Dom0 should be prohibited to use foreign mapping for sharing memory with a *DomU*, unless *DomU* initiates the procedure in response to a hypercall from *Dom0*. When this happens, *Dom0* should be provided with an encrypted copy of the memory pages and of the virtual CPU registers. The entire process should be closely monitored by the hypervisor which, after the access, should check the integrity of the affected *DomU*.

A virtualization architecture that guarantees confidentiality, integrity, and availability for the TCB of a Xen-based system is presented in [303]. A secure environment, when *Dom0* cannot be trusted, can only be ensured if the guest application is able to store, communicate and process data safely. The guest software should have access to a secure secondary storage on a remote storage server and to the network interfaces when communicating with the user. A secure run-time system is also needed.

To implement a secure run-time system, we have to intercept and control the hypercalls used for communication between a *Dom0* that cannot be trusted and a *DomU* we want to protect. Hypercalls issued by *Dom0* that do not read from or write to the memory of a *DomU* or to its virtual registers should

be allowed. Other hypercalls should be restricted either completely or during a specific time window. For example, hypercalls used by *Dom0* for debugging or for the control of the IOMMU should be prohibited. The Input/Output Memory Management Unit (IOMMU) connects the main memory with a DMA-capable I/O bus; it maps device-visible virtual addresses to physical memory addresses and provides memory protection from misbehaving devices.

We cannot restrict some of the hypercalls issued by *Dom0*, even though they can be harmful to the security of *DomU*. For example, foreign mapping and access to the virtual registers are needed to save and restore the state of *DomU*. We should check the integrity of *DomU* after the execution of such security-critical hypercalls.

New hypercalls are necessary to protect:

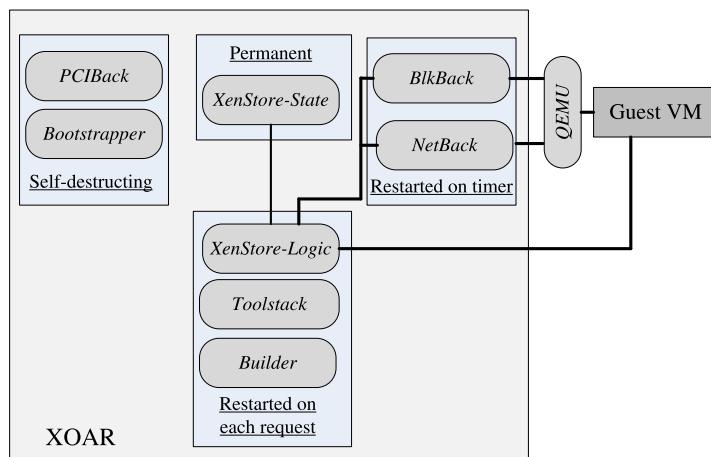
1. The privacy and integrity of the virtual CPU of a VM. When *Dom0* wants to save the state of the VM, the hypercall should be intercepted and the contents of the virtual CPU registers should be encrypted. The virtual CPU context should be decrypted, and then an integrity check should be carried out when *DomU* is restored.
2. The privacy and integrity of the VM virtual memory. The *page table update* hypercall should be intercepted and the page should be encrypted so that *Dom0* handles only encrypted pages of the VM. The hypervisor should calculate a hash of all the memory pages before they are saved by *Dom0* to guarantee the integrity of the system. Address translation is necessary because a restored *DomU* may be allocated a different memory region [303].
3. The freshness of the virtual CPU and the memory of the VM. The solution is to add to the hash a version number.

As expected, the increased level of security and privacy leads to an increased overhead. Measurements reported in [303] show increases by a factor of: 1.7 to 2.3 for the domain build time, 1.3 to 1.5 for the domain save time, and 1.7 to 1.9 for the domain restore time.

8.13 Xoar—breaking the monolithic design of the TCB

Xoar is a modified version of Xen designed to boost system security [114]. The security model of Xoar assumes that the system is professionally managed and that a privileged access to the system is granted only to system administrators. The model also assumes that the administrators have neither financial incentives, nor the desire to violate the trust of the user. The security threats come from a guest VM that could attempt to violate the data integrity or the confidentiality of another guest VM on the same platform, or to exploit the code of the guest. Other sources of threats are bugs in initialization code of the management VM. Xoar is based on micro-kernel⁸ design principles. Xoar modularity makes exposure to risk explicit and enables the guests to configure the access to services based on their needs. Modularity allows the designers of Xoar to reduce the size of the permanent footprint of the system and

⁸ A microkernel (μ -kernel) supports only the basic functionality of an operating system kernel including low-level address space management, thread management, and interprocess communication. Traditional operating system components, such as device drivers, protocol stacks, and file systems, are removed from the microkernel and run in the user space.

**FIGURE 8.5**

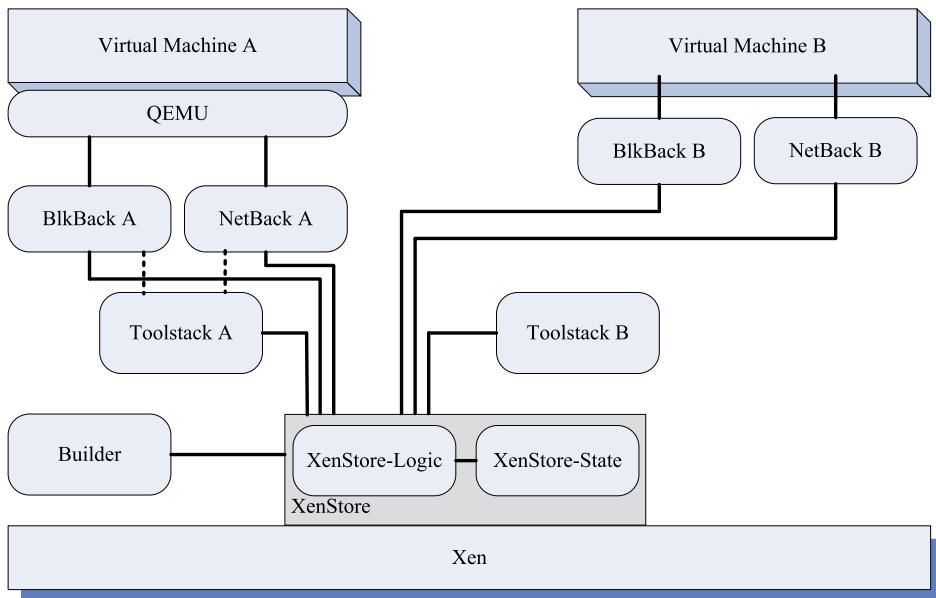
Xoar has nine classes of components of four types: permanent, self-destructing, restarted upon request, and restarted on timer. A guest VM is started using the *Toolstack* by the *Builder* and is controlled by the *XenStore-Logic*. The devices used by the guest VM are emulated by the *QEMU* component.

increase the level of security of critical components. The ability to record a secure audit log is another critical function of a hypervisor facilitated by a modular design. The design goals of Xoar are:

- Maintain the functionality provided by Xen.
- Ensure transparency with existing management and VM interfaces.
- Tight control of privileges. Each component should only have the privileges required by its function.
- Minimize the interfaces of all components to reduce the possibility that a component can be used by an attacker.
- Eliminate sharing. Make sharing explicit whenever it cannot be eliminated to allow meaningful logging and auditing.
- Reduce the opportunity of an attack targeting a system component by limiting the time window during which the component runs.

These design principles aim to break the monolithic TCB design of a *Xen*-based system. Inevitably, this strategy has an impact on performance, but the implementation should attempt to keep the modularization overhead to a minimum.

A close analysis shows that booting the system is a complex activity, but the fairly large modules used during booting are no longer needed once the system is up and running. In Section 5.8, we have seen that *XenStore* is a critical system component because it maintains the state of the system, and, thus, it is a prime candidate for hardening. The *ToolStack* is only used for management functions and can only be loaded upon request. Xoar has four types of components: permanent, self-destructing, restarted upon request, and restarted on timer; see Fig. 8.5:

**FIGURE 8.6**

Component sharing between guest VM in Xoar. Two VMs share only the *XenStore* components. Each one has a private version of the *BlkBack*, *NetBack*, and *Toolstack*.

1. Permanent components. *XenStore-State* maintains all information regarding the state of the system.
2. Components used to boot the system; they self-destruct before any user VM is started. The two components discover the hardware configuration of the server including the PCI drivers and then boot the system: *PCIBack*—virtualizes access to PCI bus configuration and *Bootstrapper*—coordinates booting of the system.
3. Components restarted on each request: *XenStore-Logic*; *Toolstack*—handles VM management requests, e.g., it requests the *Builder* to create a new guest VM in response to a user request; and *Builder*—initiates user VMs.
4. Components restarted on a timer: the two components export physical storage device drivers and the physical network driver to a guest VM. *Blk-Back*—exports physical storage device drivers using *udev*⁹ rules. *NetBack*—exports the physical network driver.

Another component, QEMU, is responsible for device emulation. *Bootstrapper*, *PCIBack*, and *Builder* are the most privileged components, but the first two are destroyed once Xoar is initialized. The *Builder* is very small, it consists of only 13 000 lines of code. XenStore is broken into two components, *XenStore-Logic* and *XenStore-State*. Access control checks are done by a small monitor module in *XenStore-State*. Guest VMs share only the *Builder*, *XenStore-Logic*, and *XenStore-State*; see Fig. 8.6.

⁹ *udev* is the device manager for the Linux kernel.

Users of Xoar are able to share only service VMs with guest VMs that they control; to do so, they specify a tag on all of the devices of their hosted VMs. Auditing is more secure: whenever a VM is created, deleted, stopped, or restarted by Xoar, the action is recorded in an append-only database on a different server accessible via a secure channel.

Rebooting provides the means to ensure that a VM is in a known good state. To reduce the overhead and the increased startup time demanded by a reboot, Xoar uses *snapshots* instead of rebooting. The service VM snapshots itself when it is ready to service a request; similarly, snapshots of all components are taken immediately after their initialization and before they start interacting with other services or guest VMs. Snapshots are implemented using a copy-on-write (COW) mechanism¹⁰ to preserve any page about to be modified.

8.14 Mobile devices and cloud security

Mobile devices are an integral part of the cloud ecosystem; mobile applications use cloud services to access and store data or to carry out a multitude of computational tasks. The security challenges for mobile devices are common to all computer and communication systems and include: (i) Confidentiality—ensure that transmitted and stored data cannot be read by unauthorized parties; (ii) Integrity—detect intentional or unintentional changes to transmitted and stored data; (iii) Availability—ensure that users can access cloud resources whenever needed; and (iv) Nonrepudiation—the ability to ensure that a party to a contract cannot deny the sending of a message that they originated.

The technology stack of a mobile device consists of the hardware, the firmware, the operating system, and the applications. The separation between the firmware and the hardware of a mobile device is blurred. A baseband processor is used solely for telephony services involving data transfers over cellular networks operating outside the control of the mobile OS that runs on the application processor. Security-specific hardware and firmware stores encryption keys, certificates, credentials, and other sensitive information on some mobile devices.

The nature of mobile devices places them at higher exposure to threats than stationary ones. Mobile devices are designed to easily install applications, to use third-party applications from application stores, and to communicate with computer clouds via often untrusted cellular and WiFi networks. Mobile devices interact frequently with other systems to exchange data and often use untrusted content.

Mobile devices often require a short authentication passcode and may not support strong storage encryption. Location services increase the risk of targeted attacks. Potential attackers are able to determine user's location, correlate the location with information from other sources on the individuals the user associates with, and infer other sensitive information. Special precautions must then be taken due to exposure to the unique security threats affecting mobile devices, including: (i) mobile malware; (ii) stolen data due to loss, theft, or disposal; (iii) unauthorized access; (iv) electronic eavesdropping and electronic tracking; and (v) access to data by third-party applications. Some of these threats can

¹⁰ Copy-on-write (COW) is used by virtual memory operating systems to minimize the overhead of copying the virtual memory of a process when a process creates a copy of itself. Then, the pages in memory that might be modified by the process or by its copy are marked as COW. When one process modifies the memory, the operating system's kernel intercepts the operation and copies the memory so that changes in one process's memory are not visible to the other.

propagate to the cloud infrastructure a mobile device is connected to. For example, files stored on the mobile devices subject to ransomware and encrypted by a malicious intruder can migrate to the backup stored on the cloud. The risks posed to the cloud infrastructure by mobile devices are centered around data leakage and compromise. Such security risks are due to a set of reasons including:

1. Loss of the mobile device, lock-screen protection, enabling smudge attacks, and other causes leading to mobile access control. A smudge attack is a method to discern the password pattern of a touchscreen device such as a cell phone or tablet computer.
2. Lack of confidentiality protection for data in transit in unsafe or untrusted WiFi or cellular networks.
3. Unmatched firmware or software including operating system and application software bypassing the security architecture, e.g., rooted/jailbroken devices.
4. Malicious mobile applications bypassing access control mechanisms.
5. Misuse or misconfiguration of location services, such as GPS.
6. Acceptance of fake mobility management profiles.

An in-depth discussion of the Enterprise Mobile Management (EMM) lists various EMM services, including the Mobile Device Management (MDM) and the Mobile Application Management (MAM), and suggests a number of functional and security capabilities of the system [184]. Some of these policies and mechanisms should also be applied to mobile devices connected to computer cloud:

- i. Use device encryption, application-level encryption, and remote wipe capabilities to protect storage.
- ii. Use Transport Layer Security (TLS) for all communication channels.
- iii. Isolate user-level applications from each other to prevent data leakage between applications using sandboxing.
- iv. Use device integrity checks for boot validation, verified application, and OS updates.
- v. Use auditing and logging.
- vi. Enforce authentication of the device owner.
- vii. Automatic, regular device integrity and compliance checks for threats and compliance.
- viii. Automated alerts for policy violations.

A system for Microsoft Outlook mobile application requires individuals who wish to participate in a managed scenario to download the Microsoft Community Portal application and input the required information including local authentication to the mobile OS via a lock screen and the encryption capabilities provided by the mobile OS to protect data on the device. The cloud MDM portal discussed in [184] is available to administrators through a web interface.

8.15 Mitigating cloud vulnerabilities in the age of ransomware

In Section 2.13, we have seen that computer clouds are changing the enterprise computing landscape at stunning speed. The massive migration to cloud computing is motivated by the desire to reduce costs and operate in a more secure environment. The latter rationale is being challenged by this explosive migration; the number of potentially vulnerable and attractive targets increases and hackers are becoming more sophisticated and able to exploit vulnerabilities in cloud systems.

Data breaches due to cloud misconfigurations cost businesses nearly \$3.18 trillion in 2019. Some vulnerabilities are exposed by the increasing larger size of the cloud user population. Gartner predicts that by 2025, 99% of cloud security failures will be due to customers fault.

Nearly two-thirds of businesses and organizations view cloud security as the biggest impediment to adoption. According to CSA, the top security issues affecting organizations using the cloud are: malware proliferation (63%), advanced persistent threats, (53%) compromised accounts (43%), and insider threats (42%). We start the discussion of cloud vulnerabilities in 2021 with a few relevant statistics regarding enterprises on the cloud¹¹:

1. 29% of enterprises on the cloud have experienced potential account compromises.
2. 27% of enterprises on the cloud permit root user activities.
3. 49% of enterprises on the cloud do not encrypt their cloud databases.
4. 21% of cloud-hosted files include sensitive data, a 17% increase in the past two years.
5. 83% of businesses worldwide store sensitive data on the cloud.
6. 41% of access keys remain unchanged on the cloud in the past 90 days.
7. 22% of cloud users share files, and as much as 48% of cloud files eventually get shared.
8. 20% of the sensitive data passing through clouds is due to email services; this volume has risen by 59% in the past two years.

The most significant cloud security threats are: (a) access management; (b) data breaches and data leaks; (c) data loss; (d) insecure APIs; and (e) misconfigured cloud storage.

Access management enforces a rigorous access policy and uses authentication and identity verification tools. To mitigate the threats, many CSPs now offer multifactor authentication systems as part of their standard packages. At the same time, IT administrators in charge of access management need to periodically audit the level of access of all employees and promptly remove the privileges of departing employees; they should deploy the most secure authentication and identity verification tools available.

Data breaches are a significant threat in cloud systems. Attackers take advantage of expired digital certificates, as in the case of *Equifax* when personal data of more than 148 million individuals was stolen by hackers in 2017. To mitigate this threat, organizations should use encryption for email servers and messages. Organizations should also use digital certificates, such as SSL/TLS website certificates and S/MIME (secure/multipurpose Internet mail extension) certificates.

Data loss is a major concern for cloud users. An average of 51% of organizations have publicly exposed at least one cloud storage service. 84% of organizations report that traditional security solutions do not work in cloud environments. Regular and thorough backups on secure storage devices are critical to protect valuable data from ransomware attacks. A hacker with access to the systems used to access the cloud could encrypt user data files that are then automatically backed up to the cloud; a hacker may also encrypt cloud storage and demand payment for returning data.

Many APIs still have security vulnerabilities, e.g., some APIs give CSPs access to user data. It was recently revealed that Facebook and Google stored user passwords in plaintext, allowing staff to access them. Before moving to a CSP any organization should make sure that the CSP adheres to OWASP API security guidelines see <https://owasp.org/www-project-api-security/>.

¹¹ The statistics and the data on security threats are from <https://cloud-standards.org/cloud-computing-statistics/> and <https://securityboulevard.com/2020/05/cloud-security-5-serious-emerging-cloud-computing-threats-to-avoid/>, respectively.

Misconfiguration of cloud storage exposes confidential data. This can be caused by either data stored in large and confusing structures leaving important files unprotected, or by failing to change the default security settings on the cloud storage. An example of misconfigured cloud storage is the National Security Agency mistake that made a number of top-secret documents available via an external browser.

The shared security responsibility (SSR) model discussed in Section 2.1 and illustrated by Table 2.1 is a powerful concept for cloud security with different implications for different cloud delivery models. It may seem counterintuitive that SaaS is a most vulnerable target though the cloud infrastructure and the software stack are all controlled by the CSP.

A closer look reveals that SaaS access to data is solely the responsibility of the large number of users, some lacking basic elements of computer literacy. Even the most sophisticated SaaS uses have minimal control or visibility of the underlying infrastructure of the services they use. Malicious attacks such as *GoldenEye* and *XcodeGhost* ransomware show that hackers are focused on high-value targets—such as health care institutions and critical services provided by local communities.

The most relevant SaaS security issues are: (i) the lack of total control over who can access sensitive data in the cloud; (ii) the inability to control and keep track of data in transit to and from cloud applications; (iii) the inability to prevent the misuse of data or theft by insiders; (iv) the lack of skilled staff to manage the security of cloud applications used by an organization; and (v) the inability to assess the security operations of the cloud service providers.

Any attack on the IaaS infrastructure affects a large number of users. IaaS user responsibilities reflect the extended degrees of freedom allowed to IaaS users. A hostile takeover of resources allocated to a user and their use as bots to launch attacks against third parties by malicious individuals is possible. Advanced persistent attacks mounted on the IaaS infrastructure are possible.

8.16 AWS security

To conclude this chapter, we take a closer look at the security support available at one IaaS cloud service provider. We choose AWS due to the large number of organizations that demand high security standards hosted by Amazon. Security experts monitor the vast AWS infrastructure and build and maintain security services to prevent, detect, respond, and remediate the effects of security threats. We start with a review of AWS security tools.¹²

CloudTrail is a service that enables governance, compliance, operational auditing, and risk auditing of an AWS account. The event history, including Management Console activity, AWS SDK (Software Development Kit), and all AWS services invoked, simplifies security analysis, resource change tracking, and troubleshooting.

The Inspector (IS) is an automated security assessment service; it assesses applications for exposure, vulnerabilities, and deviations from best practices. After performing an assessment, IS produces a detailed list of security findings prioritized by level of severity and checks for unintended network accessibility of EC2 instances and for vulnerabilities on those EC2 instances.

¹² An in-depth of discussion of AWS security can be found at <https://aws.amazon.com/security/>.

Web Application Firewall (WAF) protects web applications or APIs against common web exploits and bots that may compromise security, consume excessive resources, or affect availability. WAF provides control over the traffic generated by applications and blocks attack patterns, such as SQL injection¹³ or cross-site scripting (XSS).¹⁴

Cognito is used for identity management. It can detect brute-force authentication and fraudulent login attempts. Cognito User Pools supports identity and access management standards, such as Oauth 2.0, SAML 2.0, and OpenID Connect.

CloudHSM helps generate encryption keys using managed hardware security modules. An AWS user can configure the Key Management Service (KMS) to use CloudHSM cluster as a custom key store, rather than the default KMS key store. It automatically load balances requests and securely duplicates keys stored in any HSM to all of the other HSMs in a cluster.

CloudFront is a content delivery network. It protects applications from DDoS attacks and transfers data securely at high speeds. It offers security capabilities, such as field level encryption and HTTPS support, seamlessly integrated with Shield, WAF, and Route 53 to protect against multiple types of attacks including network and application layer DDoS (Distributed Denial of Service) attacks. With edge compute features and Lambda@Edge, one can run code across AWS locations globally.

We believe that some of the most consequential aspects of AWS security are:

- i. Fine-grain identity, access controls, and continuous monitoring for near real-time security information.
- ii. Automation to reduce human configuration errors.
- iii. Extended use of encryption. Data flowing across the AWS global network is automatically encrypted at the physical layer before leaving a secured facilities.
- iv. The added level of competence provided by a global program of Technology and Consulting Partners specialized in security—the AWS Partner Network.

8.17 Further readings

Cloud Security Alliance (CSA) is an organization with more than 100 corporate members. It aims to address all aspects of cloud security and serve as a cloud security standards incubator. The reports, available from the web site of the organization, are periodically updated; the original report was published in 2009 [120], and subsequent reports followed [121] and [122]. An open security architecture is presented in [380].

A seminal paper [189] on the negative implications of virtualization on system security “When virtual is harder than real: security challenges in VM based computing environments” by Garfinkel and Rosenblum was published in 2005, followed by another one that reached similar conclusions, [404]. Risk and trust are analyzed in [152], [254], and [319]. Cloud security is also discussed in [344], [459], and [465]. Managing cloud information leakage is the topic of [512].

¹³ Malicious SQL statements are inserted into an entry field for execution.

¹⁴ Malicious scripts are injected into trusted websites enabling web application to send malicious code as a browser script to a different end user.

A 2010 paper [212] presents a taxonomy of attacks on computer clouds and [137] covers the management of security services lifecycle. Security issues vary depending on the cloud model as discussed in [375]. The privacy impact in cloud computing is the topic of [469]. A 2011 book [515] gives a comprehensive look at cloud security. Privacy and protection of personal data in the EU is discussed in a document available at <http://ec.europa.eu/justice/policies/privacy>.

The paper [35] analyzes the inadequacies of current risk controls for the cloud. Intercloud security is the theme of [59]. Secure collaborations are discussed in [62]. The paper [304] presents an approach for secure VM execution under untrusted management OS. The social impact of privacy in cloud computing is analyzed in [166]. An anonymous access control scheme is presented in [260].

An empirical study into the security exposure to hosts of hostile virtualized environments can be found at <http://taviso.decsystem.org/virtsec.pdf>. A model-based security testing approach for cloud computing is presented in [535]. Cold boot attacks on encryption keys are discussed in [222]. Cloud security concerns and mobile device security are covered in several NIST documents. [401] introduces an encryption system for query processing and [435] discusses a pragmatic security discipline.

8.18 Exercises and problems

- Problem 1.** Identify the main security threats for the SaaS cloud delivery model on a public cloud. Discuss the various aspects of these threats on a public cloud as compared to the threats posed to similar services provided by a traditional service-oriented architecture running on a private infrastructure.
- Problem 2.** Analyze how the six attack surfaces discussed in Section 8.2 and illustrated in Fig. 8.1 apply to the SaaS, PaaS, and IaaS cloud delivery models.
- Problem 3.** Analyze Amazon privacy policies, and design a service level agreement you would sign if you were to process confidential data using AWS.
- Problem 4.** Analyze the implications of the lack of trusted paths in commodity operating systems and give one or more examples showing the effects of this deficiency. Analyze the implications of the two-level security model of commodity operating systems.
- Problem 5.** Compare the benefits and the potential problems due to virtualization on public, private, and hybrid clouds.
- Problem 6.** Read [44] and discuss the measures taken by Amazon to address the problems posed by shared images available from AWS. Would it be useful to have a cloud service to analyze images and sign them before being listed and made available to the general public?
- Problem 7.** Analyze the risks posed by foreign mapping and the solution adopted by Xoar. What is the security risk posed by XenStore?
- Problem 8.** Read [114] and discuss the performance of the system. What obstacles to its adoption by the providers of IaaS services can you foresee?
- Problem 9.** Discuss the impact of international agreements regarding privacy laws on cloud computing.
- Problem 10.** Propagation of the malware in the Internet has similarities with the propagation of an infectious disease. Discuss the three models for the propagation of an infectious disease in a finite population, *SI*, *SIR*, and *SIS*. Justify the formulas describing the dynamics of the system for each model. *Hint:* read [73, 160] and [270].

Cloud resource management and scheduling

9

Resource management is a core function of any manmade system; it affects the three dimensions of the system evaluation space: performance, functionality, and cost. Efficient resource management has a direct effect on performance and cost and an indirect effect on the functionality of the system because some of the functions may be avoided due to the poor performance and/or high cost.

A cloud is a complex system with a large number of shared resources subject to unpredictable requests and affected by external events it cannot control. Cloud resource management requires complex policies and decisions for multiobjective optimization. Effective resource management is extremely challenging due to the scale of the cloud infrastructure and to the unpredictable interactions of the system with a large population of users. The scale makes it impossible to have accurate global state information, and the large user population makes it nearly impossible to predict the type and intensity of system workload.

Resource management becomes even more complex when resources are oversubscribed and users are uncooperative, driven by self-interest. In addition to external factors, resource management is affected by internal factors, such as the heterogeneity of the hardware and software systems, the ability to approximate the global state of the system and to redistribute the load, the failure rates of various components, and many other factors.

Strategies for resource management associated with IaaS, PaaS, and SaaS delivery models differ. In all cases, cloud service providers are faced with large fluctuating workloads challenging the claim of cloud elasticity. When a spike can be predicted, resources can be provisioned in advance, e.g., for web services subject to seasonal spikes. The situation is slightly more complicated for unplanned spikes.

Auto-scaling can be used for unplanned spike loads provided that: (a) there is a pool of resources that can be released or allocated on demand; and (b) there is a monitoring system that allows a control loop to decide in real time to reallocate resources. Auto-scaling is supported by PaaS services, such as Google App Engine. Auto-scaling for IaaS is complicated due to the lack of standards.

Centralized control cannot provide adequate solutions for management policies when changes in the environment are frequent and unpredictable. Distributed control poses its own challenges since it requires some form of coordination between the entities in control. Autonomic policies are of great interest due to the scale of the system and the unpredictability of the load when the ratio of peak to mean resource demands can be very large.

Throughout this text we use the term *bandwidth* in a broad sense to mean the number of operations or the amount of data transferred per time unit. For example, Mips (Million Instructions Per Second) or Mflops (Million Floating Point Instructions Per Second) measure the CPU speed, and Mbps (Megabits per second) measure the speed of a communication channel. The *latency* is defined as the time elapsed from the instance an operation is initiated until the instance its effect is sensed. Latency is context dependent. For example, the latency of a communication channel is the time it takes a bit to traverse

the communication channel from its source to its destination; memory latency is the time elapsed from the instance a memory *read* instruction is issued until the time the data becomes available in a memory register. The demand for computing resources, such as CPU cycles, primary and secondary storage, and network bandwidth, depend heavily on the volume of data processed by an application.

This chapter covers a broad range of topics related to cloud resource management and scheduling. The sections marked (R) introduce the reader to cloud resource management research topics. An overview of policies and mechanisms for cloud resource management in Section 9.1 is followed by a discussion of scheduling algorithms for computer clouds in Section 9.2 and by an analysis of delay scheduling and of data-aware scheduling in Sections 9.3 and 9.4, respectively.

The Apache capacity scheduler is presented in Section 9.5, and the start-time fair queuing and the borrowed virtual time scheduling algorithms are analyzed in Sections 9.6 and 9.7, respectively. Cloud scheduling subject to deadlines and MapReduce scheduling with deadlines are discussed in Sections 9.8 and 9.9, respectively. Resource bundling and combinatorial auctions are covered in Section 9.10.

Cloud energy efficiency and cloud resource utilization and the impact of application scaling on resource management are analyzed in Sections 9.11 and 9.12, respectively. A control theoretic approach to resource allocations is discussed in Sections 9.13, 9.14, and 9.15, followed by a machine learning algorithm for coordination of specialized autonomic performance managers in Section 9.16.

A utility model for resource allocation for a web service is then presented in Section 9.17. Cloud infrastructure will most likely continue to scale up to accommodate the demands of an increasingly larger cloud user community. A fair question is whether current cloud management systems can sustain this expansion. Self-organization and self-management alternatives come to mind, but the very slow progress made by the autonomic computing initiative is likely to dampen the enthusiasm of those believing in self-management. We discuss emergence and self-organization in Section 9.18 and conclude this chapter with an in-depth analysis of cloud interoperability in Section 9.19.

9.1 Policies and mechanisms for resource management

A *policy* refers to the *principles guiding the decisions*, while *mechanisms* represent the *means to implement policies*. Separation of policies from mechanisms is a guiding principle in computer science. Butler Lampson [295] and Per Brinch Hansen [225] offer solid arguments for this separation in the context of operating system design.

Cloud resource management policies can be loosely grouped into five classes: admission control, capacity allocation, load balancing, energy optimization, and QoS guarantees. The explicit goal of an *admission control policy* is to prevent the system from accepting workload in violation of high-level system policies; for example, a system may not accept additional workload that would prevent it from completing work already in progress or contracted.

Limiting the workload requires some knowledge of the global state of the system; in a dynamic system such knowledge, when available, is at best obsolete. *Capacity allocation* is concerned with resource allocation to individual instances; an instance is an activation of a service. Locating resources for an instance is subject to multiple global optimization constraints and requires the search of a very large search space when the state of individual systems are changing rapidly.

Load balancing and energy optimization can be done locally, but global load balancing and energy optimization policies encounter the same difficulties as the ones we have already discussed. Load balancing and energy optimization are correlated and affect the cost for providing services [145].

The common meaning of the term “load balancing” is an even distribution of the workload to a set of servers. For example, consider the case of four identical servers, A , B , C , and D , whose relative loads are 80%, 60%, 40%, and 20%, respectively, of their capacity. As a result of a perfect load balancing all servers would end up with the same load, 50% of each server’s capacity.

An important goal of cloud resource management is minimization of the cost for providing cloud service and, in particular, minimization of cloud energy consumption. This leads to a different meaning of “load balancing;” instead of evenly distributing the load amongst all servers, we wish to *concentrate the workload and use the smallest number of servers* while switching the others to a standby mode, a state where a server uses very little energy. In our example the load from D will migrate to A , and the load from C will migrate to B ; thus, A and B will be loaded to full capacity, while C and D will be switched to standby mode. *Quality of Service* is the resource management facet probably the most difficult to address and, at the same time possibly the most critical for the future of cloud computing.

Often, as we shall see in this section, resource management strategies jointly target performance and power consumption. The Dynamic Voltage and Frequency Scaling (DVFS)¹ techniques, such as Intel’s SpeedStep and AMD’s PowerNow, lower the voltage and the frequency to decrease the power consumption.² Motivated initially by the need to save power for mobile devices, these techniques have migrated virtually to all processors, including the ones used for high-performance servers.

Processor performance decreases, but at a substantially lower rate than energy consumption, as a result of lower voltages and clock frequencies [301]. Table 9.1 shows the dependence of normalized performance and normalized energy consumption of a typical modern processor on the clock rate. As we can see, at 1.8 GHz, we save 18% of the energy required for maximum performance, while the performance is only 5% lower than the peak performance achieved at 2.2 GHz. This seems a reasonable energy–performance tradeoff!

Virtually all optimal, or near-optimal, mechanisms to address the five classes of policies do not scale up and typically target a single aspect of resource management, e.g., admission control, but ignore energy conservation. Many require complex computations that cannot be done effectively in the time available to respond. Performance models are very complex, analytical solutions are intractable, and the monitoring systems used to gather state information for these models can be too intrusive and unable to provide accurate data.

Many techniques are concentrated on system performance in terms of throughput and time in system, but rarely include energy trade-offs or QoS guarantees. Some techniques are based on unrealistic assumptions. For example, capacity allocation is viewed as an optimization problem, but under the assumption that servers are protected from overload.

Allocation techniques in computer clouds must be based on a systematic approach, rather than ad hoc methods. The four basic mechanisms for the implementation of resource management policies are:

¹ Dynamic voltage and frequency scaling is a power management technique to increase or decrease the operating voltage or frequency of a processor to increase the instruction execution rate and, respectively, to reduce the amount of heat generated and to conserve power.

² The power consumption P of a CMOS-based circuit is: $P = \alpha \cdot C_{eff} \cdot V^2 \cdot f$ with: α —the switching factor, C_{eff} —the effective capacitance, V —the operating voltage, and f —the operating frequency.

Table 9.1 Normalized performance and energy consumption, function of processor speed; performance decreases at a lower rate than energy consumption when the clock rate decreases.

CPU speed (GHz)	Normalized energy (%)	Normalized performance (%)
0.6	0.44	0.61
0.8	0.48	0.70
1.0	0.52	0.79
1.2	0.58	0.81
1.4	0.62	0.88
1.6	0.70	0.90
1.8	0.82	0.95
2.0	0.90	0.99
2.2	1.00	1.00

- *Control theory*. Control theory uses the feedback to guarantee system stability and predict transient behavior [264,287], but can be used only to predict local, rather than global, behavior; Kalman filters have been used for unrealistically simplified models.
- *Machine learning*. A major advantage of machine learning techniques is that they do not need a performance model of the system [480]; this technique could be applied for coordination of several autonomic system managers, as discussed in [268].
- *Utility-based*. Utility-based approaches require a performance model and a mechanism to correlate user-level performance with cost as discussed in [11].
- *Market-oriented/economic mechanisms*. Such mechanisms do not require a model of the system, e.g., combinatorial auctions for bundles of resources discussed in [456].

A distinction should be made between interactive and noninteractive workloads. Management techniques for interactive workloads, e.g., web services, involve flow control and dynamic application placement, while those for noninteractive workloads are focused on scheduling. A fair amount of work reported in the literature is devoted to resource management of interactive workloads, some to noninteractive ones, and only a few to heterogeneous workloads, a combination of the two, e.g. [467].

9.2 Scheduling algorithms for computer clouds

Scheduling is a critical component of the cloud resource management responsible for resource sharing/multiplexing at several levels. A server can be shared among several VMs, each VM can support several applications, and each application may consist of multiple threads. CPU scheduling supports the virtualization of a processor, the individual threads acting as virtual processors; a communication link can be multiplexed among a number of virtual channels, one for each flow.

In addition to the need to meet its design objectives, a scheduling algorithm should be efficient, fair, and starvation-free. The objectives of a scheduler for a batch system are to maximize the throughput (the number of jobs completed in one unit of time, e.g., in one hour) and to minimize the turnaround time (the time between job submission and its completion). The objectives of a real-time system scheduler are to meet the deadlines and to be predictable.

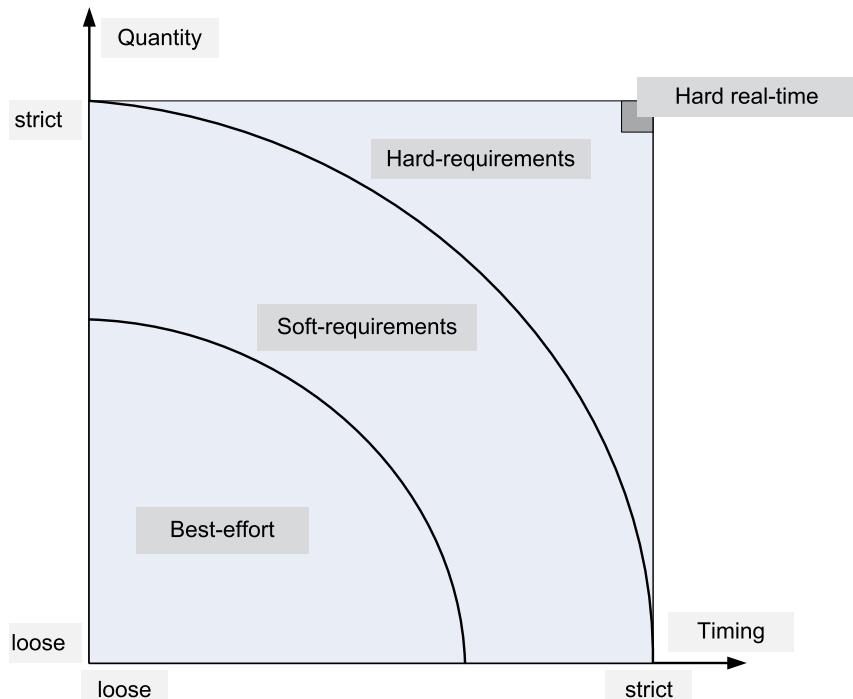


FIGURE 9.1

Resource requirement policies. *Best-effort* policies do not impose requirements regarding either the amount of resources allocated to an application, or the timing when an application is scheduled. *Soft-requirements* policies require statistically guaranteed amounts of resources and timing constraints. *Hard-requirements* policies demand strict timing and precise amounts of resources.

Schedulers for systems supporting a mixture of tasks, some with hard real-time constraints, others with soft, or no timing constraints, are often subject to contradictory requirements. Some schedulers are *preemptive*, allowing a high-priority task to interrupt the execution of a lower priority one; others are *nonpreemptive*.

Two distinct dimensions of resource management must be addressed by a scheduling policy: (a) the amount/quantity of resources allocated; and (b) the timing when access to resources is granted. Fig. 9.1 identifies several broad classes of resource allocation requirements in the space defined by these two dimensions: best-effort, soft requirements, and hard requirements. Hard real-time requirements are the most challenging, as they require strict timing and precise amounts of resources.

There are multiple definitions of a fair scheduling algorithm. First, we discuss the *max-min fairness criterion* [185]. Consider a resource with bandwidth B shared among n users who have equal rights; each user requests an amount b_i and receives B_i . Then, according to the max-min criterion, the following conditions must be satisfied by a fair allocation:

- C_1 —the amount received by any user is not larger than the amount requested, $B_i \leq b_i$.

- C_2 —if the minimum allocation of any user is B_{min} , no allocation satisfying condition C_1 has a higher B_{min} than the current allocation.
- C_3 —when we remove the user receiving the minimum allocation B_{min} and then reduce the total amount of the resource available from B to $(B - B_{min})$, the condition C_2 remains recursively true.

A fairness criterion for CPU scheduling [205] requires that the amount of work $\Omega_a(t_1, t_2)$ and $\Omega_b(t_1, t_2)$ in the time interval from t_1 to t_2 of two runnable threads a and b minimize the expression

$$\left| \frac{\Omega_a(t_1, t_2)}{w_a} - \frac{\Omega_b(t_1, t_2)}{w_b} \right|, \quad (9.1)$$

where w_a and w_b are the weights of the threads a and b , respectively.

The QoS requirements differ for different classes of cloud applications and demand different scheduling policies. Best-effort applications, such as batch applications and analytics,³ do not require QoS guarantees. Multimedia applications such as audio and video streaming have soft real-time constraints and require statistically guaranteed maximum delay and throughput. Applications with hard real-time constraints do not use a public cloud at this time, but may do so in the future.

Round-robin, first-come-first-serve (FCFS), shortest-job-first (SJF), and priority algorithms are among the most common scheduling algorithms for best-effort applications. Each thread is given control of the CPU for a definite period of time, called a *time-slice*, in a circular fashion in case of round-robin scheduling; the algorithm is fair and starvation-free. The threads are allowed to use the CPU in the order they arrive in the case of the FCFS algorithms and in the order of their running time in the case of SJF algorithms.

Earliest Deadline First (EDF) and *Rate Monotonic Algorithms* (RMA) are used for real-time applications. Integration of scheduling for the three classes of applications is discussed in [71] and two new algorithms for integrated scheduling, the *Resource Allocation/Dispatching* (RAD) and the *Rate-Based Earliest Deadline* (RBED) are proposed.

Several algorithms of special interest for computer clouds are discussed below. These algorithms illustrate the evolution in thinking regarding the fairness of scheduling and the need to accommodate multi-objective scheduling, in particular scheduling for Big Data and for multimedia applications.

9.3 Delay scheduling (R)

How does one simultaneously ensure fairness and maximize resource utilization without compromising locality and throughput for Big Data applications running on large computer clusters? This is a question faced early on by large IT service providers and we discuss it next.

Hadoop scheduler. A Hadoop job consists of multiple Map and Reduce tasks and the question is how to allocate resources to the tasks of newly submitted jobs. Recall from Section 11.7 that the *job tracker* of the *Hadoop master* manages a number of slave servers running under the control of *task trackers*

³ The term *analytics* is overloaded; sometimes it means discovery of patterns in the data; it could also mean statistical processing of the results of a commercial activity.

with slots for Map and Reduce tasks. A FIFO scheduler with five priority levels assigns slots to tasks based on their priority. The lower the number of tasks of a job already running on slots of all servers, the higher is the priority of the remaining tasks.

An obvious problem with this policy is that priority-based allocation does not consider data locality, namely the need to place tasks close to their input data. The network bandwidth in a large cluster is considerably lower than the disk bandwidth; also the latency for local data access is much lower than the latency of a remote disk access. Locality affects the throughput: *server locality*, i.e., getting data from the local server is significantly better in terms of time and overhead than *rack locality*, i.e. getting input data from a different server in the same rack.

In steady-state, priority scheduling leads to the tendency to assign the same slot repeatedly to the next task(s) of the same job. As one of the job's tasks completes execution its priority decreases and the available slot is allocated to the next task of the same job. Input data for every job are striped over the entire cluster so spreading the tasks over the cluster can potentially improve data locality, but the priority scheduling favors the occurrence of *sticky slots*.

According to [532] sticky slots did not occur in Hadoop at the time of the report “due to a bug in how Hadoop counts running tasks. Hadoop tasks enter a *commit pending* state after finishing their work, where they request permission to rename their output to its final filename. The job object in the master counts a task in this state as running, whereas the slave object doesn't. Therefore, another job can be given the task's slot.” Data gathered at Facebook shows that only 5% of the jobs with a low number, 1–25, of Map tasks achieve server locality and only 59% show rack locality.

Task locality and average job locality. A task assignment satisfies the locality requirement if the input task data are stored on the server hosting the slot allocated to the task. We wish to compute the expected locality of job \mathcal{J} with the fractional cluster share $f_{\mathcal{J}}$, assuming that a server has L slots and that each block of the files system has R replicas. By definition, $f_{\mathcal{J}} = n/N$ with n the number of slots allocated to job \mathcal{J} and N the number of servers in the cluster.

The probability that a slot does not belong to \mathcal{J} is $(1 - f_{\mathcal{J}})$, there are R replicas of block $\mathcal{B}_{\mathcal{J}}$ of job \mathcal{J} and each replica is on a node with L slots. Thus, the probability that none of the slots of job \mathcal{J} has a copy of block $\mathcal{B}_{\mathcal{J}}$ is $(1 - f_{\mathcal{J}})^{RL}$. It follows that $\mathcal{L}_{\mathcal{J}}$, the locality of job \mathcal{J} , can be at most

$$\mathcal{L}_{\mathcal{J}} = 1 - (1 - f_{\mathcal{J}})^{RL}. \quad (9.2)$$

How should a fair scheduler operate on a shared cluster? What is the number n of slots of a shared cluster the scheduler should allocate to jobs assuming that tasks of all jobs take an average of T seconds to complete? A sensible answer is that the scheduler should provide enough slots such that the response time on the shared cluster should be the same as $R_{n,\mathcal{J}}$, the completion time of job \mathcal{J} would experience on a fictitious private cluster with n available slots for the n tasks of \mathcal{J} as soon as job \mathcal{J} arrives.

The job completion time can be approximated by the sum of the job processing time for all but the last task plus the waiting time until a slot becomes available for the last task of the job. There is no waiting time for running on the private cluster with n slots, while on a shared cluster the last task will have to wait before a slot is available.

A slot allocated to a task on the shared cluster will be free on average every T/N seconds; thus, the time the job will have to wait until all its n tasks have found a slot on the shared cluster will be $n \times T/N$. This implies that a fair scheduler should guarantee that the waiting time of job \mathcal{J} on the

shared cluster is much smaller than the completion time, $R_{n,\mathcal{J}}$ on the fictitious private cluster

$$R_{n,\mathcal{J}} \gg f_{\mathcal{J}} \times T. \quad (9.3)$$

Eq. (9.3) is satisfied if one of the following three conditions is satisfied:

1. $f_{\mathcal{J}}$ is small—there are many jobs sharing the cluster and the fraction of slots allocated to each job is small;
2. T is small—the individual tasks are short;
3. $R_{n,\mathcal{J}} \gg T$ —the completion time of a job is much longer than the average task completion time; the large jobs dominate the workload.

The cumulative distribution function of running time of MapReduce jobs at Facebook resembles a sigmoid function⁴ with the middle stage of a job duration starting at about 10 seconds and ending at about 1 000 seconds. The median completion time of a Map task is much shorter than the median completion time of a job, 19 versus 84 seconds. There are fewer Reduce tasks, but of a longer average duration, 231 seconds. Eighty-three percent of the jobs are launched within 10 seconds. Results reported in [532] show that delay scheduling performs well when most tasks are short relative to job duration, and when a running task can read a given data block from multiple locations.

Delay scheduling. A somewhat counterintuitive scheduling policy, *delay scheduling*, is proposed in [532]. As the name implies, the new policy delays scheduling the tasks of a new job for a relatively short time to address the conflict between fairness and locality.

This policy skips the task at the head of the priority queue if its input data are not available on the server where the slot is located. It repeats this process up to D times as specified by the delay scheduling algorithm pseudocode in the box on the next page. The almost doubling of the throughput under the new policy, while ensuring fairness for workloads at Yahoo and Facebook, is a good indication of delay scheduling policy merits.

An analysis of the new policy assumes a cluster with N servers and L slots per server, thus with a total number of slots $S = NL$. A job \mathcal{J} prefers slots on servers where its data are stored, call this set of slots $\mathcal{P}_{\mathcal{J}}$. Call $p_{\mathcal{J}}$ the probability that a task of job \mathcal{J} has data on the server with the slot allocated to it

$$p_{\mathcal{J}} = \frac{|\mathcal{P}_{\mathcal{J}}|}{N}. \quad (9.4)$$

The probability that a task skipped D times does not have the input data on the server the slot allocated to run the task resides decreases exponentially with D ; after being skipped D times, this probability is $(1 - p_{\mathcal{J}})^D$. For example, if $p_{\mathcal{J}} = 0.1$ and $D = 40$ then the probability of having data on the slot allocated to the task is $1 - (1 - p_{\mathcal{J}})^D = 0.99$, a 99% chance.

⁴ A sigmoid function $S(t)$ is “S-shaped.” It is defined as $S(t) = \frac{1}{1+e^{-t}}$, and its derivative can be expressed as function of itself, $S'(t) = S(t)(1 - S(t))$. $S(t)$ can describe biological evolution with an initial segment describing early/childhood stage, a median stage representing maturity, and a third stage describing the late life stage.

 Delay scheduling algorithm

```

1 Initialize  $j.skipcount$  to 0 for all jobs  $j$ 
2     when a heartbeat is received from node  $n$ 
3         if  $n$  has a free slot then
4             sort jobs in increasing order of number of running tasks
5             for  $j$  in  $jobs$  do
6                 if  $j$  has unlaunched task  $t$  with data on  $n$  then
7                     launch  $t$  on  $n$ 
8                     set  $j.skipcount = 0$ 
9                 else if  $j$  has unlaunched task  $t$  then
10                    if  $j.skipcount > D + 1$  then
11                        launch  $t$  on  $n$ 
12                    else
13                        set  $j.skipcount = j.skipcount + 1$ 
14                    end if
15                end if
16            end for
17        end if

```

How does one achieve the desired level of locality for job \mathcal{J} with n tasks? An approximate analysis reported in [532] assumes that all tasks are of the same length and that the preferred location sets $\mathcal{P}_{\mathcal{J}}$ are uncorrelated. If \mathcal{J} has k tasks left to launch and the replication factor is as before equal to R , then

$$p_{\mathcal{J}} = 1 - \left(1 - \frac{k}{N}\right)^R \quad (9.5)$$

and the probability of launching a task of \mathcal{J} 's after D skips is

$$p_{\mathcal{J},D} = 1 - (1 - p_{\mathcal{J}})^D = 1 - \left(1 - \frac{k}{N}\right)^{RD} \geq 1 - e^{-RDk/N}. \quad (9.6)$$

The expected value of $p_{\mathcal{J},D}$ is then

$$\mathcal{L}_{\mathcal{J},D} = \frac{1}{N} \sum_{k=1}^N \left(1 - e^{-RDk/N}\right) = 1 - \frac{1}{N} \sum_{k=1}^N e^{-RDk/N}. \quad (9.7)$$

Then

$$\mathcal{L}_{\mathcal{J},D} \geq 1 - \frac{1}{N} \sum_{k=1}^{\infty} e^{-RDk/N} = 1 - \frac{e^{-RD/N}}{N(1 - e^{-RD/N})} \quad (9.8)$$

It follows that a locality $\mathcal{L}_{\mathcal{J},D} \geq \lambda$ requires job \mathcal{J} to forgo D times its turn for a new slot while the head of the priority queue, with D satisfying the following condition

$$D \geq -\frac{N}{R} \ln \frac{n(1-\lambda)}{1+n(1-\lambda)} \quad \text{or} \quad D \leq \frac{N}{R} \ln \left[1 + \frac{1}{n(1-\lambda)} \right]. \quad (9.9)$$

Hadoop fair scheduler (HFS). The next objective of [532] is the development of a more complex Hadoop scheduler with several new capabilities:

1. Fair sharing at the level of users rather than jobs. This requires a two-level scheduling: the first level allocates task slots to pools of jobs using a fair sharing policy; at the second level each pool allocates its slots to jobs in the pool.
2. User controlled scheduling; the second level policy can be either FIFO or fair sharing of the slots in the pool.
3. Predictable turnaround time. Each pool has a guaranteed minimum share of slots. To accomplish this goal HFS defines a *minimum share timeout* and a *fair share timeout* and when the corresponding timeout occurs it kills buggy jobs or tasks taking a very long time. Instead of using a minimum skip count D use a *wait time* to determine how long a job waits to allocate a slot to its next ready-to-run task.

HFS creates a sorted list of jobs ordered according to its scheduling policy. Then, it scans down this list to identify the job allowed to schedule a task next, and within each pool, it applies the pool's internal scheduling policy. Pools missing their minimum share are placed at the head of the sorted list and the other pools are sorted to achieve a weighted fair sharing. The pseudocode of the HFS scheduling algorithm maintains three variables for each job j initialized as $j.level = 0$, $j.wait = 0$, and $j.skipped = false$ when a heartbeat is received from node n :

HFS scheduling algorithm

```

1 for each job  $j$  with  $j.skipped = true$ 
2   increase  $j.wait$  by the time since the last heartbeat and set  $j.skipped = false$ 
3   if  $n$  has a free slot then
4     sort jobs using hierarchical scheduling policy
5     for  $j$  in jobs do
6       if  $j$  has a node-local task  $t$  on  $n$  then
7         set  $j.wait = 0$  and  $j.level = 0$ 
8         return  $t$  to  $n$ 
9       else
10      if  $j$  has a rack-local task  $t$  on  $n$  and  $j.level > 2$  or  $j.wait > W1 + 1$  then
11        set  $j.wait = 0$  and  $j.level = 1$  return  $t$  to  $n$ 
12      else
13        if  $j.level = 2$  and  $j.level = 1$  and  $j.wait > W2 + 1$ 
14          or  $j.level = 0$  and  $j.wait > W1 + W2 + 1$  then
15            set  $j.wait = 0$  and  $j.level = 2$  return any unlaunched task  $t$  in  $j$  to  $n$ 
16        else
17          set  $j.skipped = true$ 
18      end if
19    end for
20  end if
```

A job starts at locality level 0 and can only launch node-local tasks. After at least $W1$ seconds, the job advances at level 1 and may launch rack-local tasks, then after a further $W2$ seconds, it goes to level

2 and may launch off-rack tasks. If a job launches a local task with a locality higher than the level it is on, it goes back down to a previous level.

In summary, delay scheduling can be generalized to preferences other than locality and the only requirement is to have a sorted list of jobs according to some criteria. It can also be applied to resource types other than slots. Measurements reported in [532] show that near 100% locality can be achieved by relaxing fairness.

9.4 Data-aware scheduling (R)

The analysis of delay scheduling emphasizes the important role of data locality for I/O-intensive application performance. This is the topic discussed in this section and addressed by a paper covering task scheduling of I/O-intensive applications [492].

There are many I/O-intensive applications translated into jobs described as a DAG (Direct Acyclic Graph) of tasks. Examples of such applications include MapReduce, approximate query processing applied to exploratory data analysis and interactive debugging, and machine learning applied to spam classification and machine translation.

Jobs running such applications consist of multiple sets of tasks running in every stage; later-stage tasks consume data generated by early-stage tasks. For some of these applications different subsets of tasks can be scheduled independently to optimize data locality, without affecting the correctness of the results. This is the case of an application using the gradient descent algorithm.

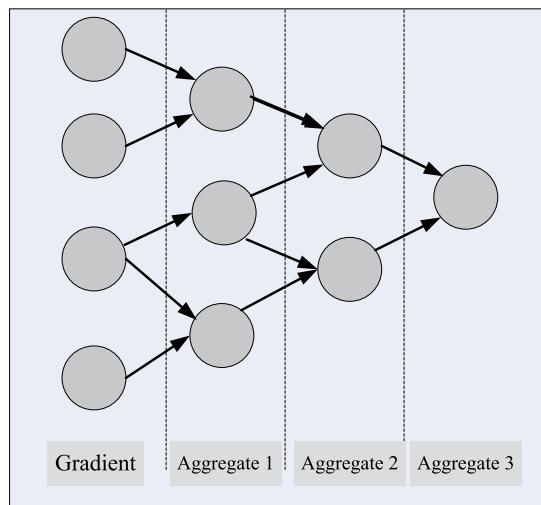
The gradient descent is a first-order iterative optimization algorithm. To find local minima the algorithm takes steps proportional with the negative of the gradient of the function at each iteration. Fig. 9.2 shows that such an application has multiple stages and that groups of tasks at each stage can be scheduled independently. Data-aware scheduling improves locality hence, response time and performance.

In this context, *late binding* means correlating tasks with data dynamically, depending on the state of the cluster. *Data-aware scheduling* improves early-stage tasks locality and, whenever possible, the locality of the later-stage tasks of the job. Locality of all tasks is important because the job completion time is determined by the slowest task; thus, a special attention should be paid to *straggler*⁵ tasks.

Recall from Chapter 4 that communication latency increases, while the communication bandwidth decreases with the “distance” between the server running a task and the one where the data is stored. An I/O-intensive task operates efficiently when its input data is already loaded in local memory and less efficiently when the data is stored on a local disk. The task efficiency decreases even further when the data resides on a different server of the same rack, and it is significantly lower when the data is on a server in a different rack.

A number of racks are interconnected by a cell switch; thus, one of the scheduler’s objective is to balance cross-rack traffic. Intermediate stages of a job often involve group communication among tasks running on servers in different racks, e.g., as one-to-many, many-to-one, and many-to-many data exchanges. A second important goal of data-aware scheduling is to reduce cross-rack communication through the placement of the producer and consumer tasks.

⁵ According to the Merriam Webster dictionary, “straggling” means to “walk or move in a slow and disorderly way,” or “spread out from others of the same kind.”

**FIGURE 9.2**

Stages of an application using the gradient descent algorithm. Tasks in the later stages of the computation use data produced by the early stage tasks. Of the four tasks in the *Gradient* stage, the top group of two tasks can be scheduled independently of the group of the lower two tasks.

KMN, the scheduler discussed in [492], launches a number of additional tasks in the early stages thus, allows choices for the later stage tasks. The name of the system comes from its basic ideas, it chooses K out of N blocks from the input data and schedules $M > K$ first-stage tasks on servers where these blocks reside. The “best” K out of $\binom{M}{K}$ choices increases the likelihood that upstream tasks outputs are distributed across racks and the next-stage tasks using these data as input are scheduled such that the cross-rack traffic is balanced. This heuristic is justified because selecting the best location of next stage tasks based on the output produced by earlier tasks is an NP-complete problem.

KMN is implemented in Scala,⁶ and it is built on top of Spark, a system for in-memory cluster computing discussed in Section 4.12. One of the novel ideas of the KMN scheduler is to choose K out of N input data blocks to improve locality of a cluster with S slots per server. If u denotes the utilization of a slot, then server utilization when all its slots are busy is u^S . The probability that one of the S tasks running on the server has a block of input data on the local disk is $p_t = 1 - u^S$.

When we choose K out of N the scheduler can choose $\binom{N}{K}$ input block combinations. The probability that K out of N tasks enjoy locality, $p_{K|N}$, is given by the binomial distribution assuming that

⁶ Scala is a general-purpose programming language with a strong static type system and support for functional programming. Scala code is compiled as Java byte code and runs on JVM (Java Virtual Machine). KMN consists of some 1 400 lines of Scala code.

the probability of success is p_t .

$$p_{K|N} = 1 - \sum_{i=0}^{K-1} \binom{N}{i} p_t^i (1-p_t)^{N-i}, \quad (9.10)$$

or

$$p_{K|N} = 1 - \sum_{i=0}^{K-1} \binom{N}{i} (1-u^S)^i u^{S(N-i)}. \quad (9.11)$$

It is easy to see that the probability of achieving locality is high even for very large utilization of the server slots, e.g., when $u = 90\%$.

The probability that all K blocks in one of the $f = \binom{N}{K}$ samples achieve locality is p_t^K , and, as the samples are independent, the probability that at least one of the samples achieves locality is

$$p_{K|N}^{(1)} = (1 - p_t^K)^f. \quad (9.12)$$

This probability increases as f increases.

Selection of the best K outputs from the M upstream tasks using a round-robin strategy is described by the following pseudocode from [492]:

```
//Given: upstreamTasks - list with rack, index within rack for each task
//Given: K - number of tasks to pick
// Number of upstream tasks in each rack
upstreamRacksCount = map()
// Initialize
for task in upstreamTasks do
    upstreamRacksCount[task:rack] += 1
end for
// Sort the tasks in round-robin fashion
roundRobin = upstreamTasks.sort(CompareTasks)
chosenK = roundRobin[0 : K]
return chosenK
procedure COMPARETASKS(task1; task2)
    if task1:idx != task2:idx then
        // Sort first by index
        return task1:idx < task2:idx
    else
        // Then by number of outputs
        numRack1 = upstreamRacksCount[task1:rack]
        numRack2 = upstreamRacksCount[task2:rack]
        return numRack1 > numRack2
    end if
end procedure
```

A hash map includes the list of the upstream tasks and how many tasks should run on each rack. Then the tasks are sorted first by their index in the rack and then by the number of tasks in the rack.

Experiments conducted on an EC2 cluster with 100 servers show that when the KMN scheduler is used instead of the native *Spark* scheduler the average job completion time is reduced by 81%. This reduction is due to the 98% locality of input tasks and a 48% improvement in data transfer. The overhead of the KMN scheduler is small, it uses 5% additional resources.

9.5 Apache capacity scheduler

Apache capacity scheduler [28] is a pluggable MapReduce scheduler for Hadoop. It supports multiple queues, job priorities, and guarantees to each queue a fraction of the capacity of the cluster. Other features of the scheduler are:

1. Free resources can be allocated to any queue beyond its guaranteed capacity. Excess allocated resources can be reclaimed and made available to another queue to meet its guaranteed capacity.
2. Excess resources taken from a queue will be restored to the queue within N minutes from the instance they are need.
3. Higher priority jobs in a queue have access to resources allocated to the queue before jobs with lower priority have access.
4. Does not support preemption; once running, a job will not be preempted for a higher priority job.
5. Each queue enforces a limit on the percentage of resources allocated to a user at any given time, if there is competition for them.
6. Supports memory-intensive jobs. A job can specify higher memory requirements than the default, and the tasks of the job will only be run on *TaskTrackers*⁷ that have enough memory to spare.

When a *TaskTracker* is free, the scheduler chooses the queue that needs to reclaim any resources the earliest and if no such queue exists it then picks a queue whose ratio of the number of running slots to the guaranteed capacity is the lowest. Once a queue is selected, the scheduler chooses a job in the queue. Jobs are sorted based on submission time and priority (if supported). Once a job is selected, the scheduler chooses a task to run. Periodically, the scheduler takes actions allowing queues to reclaim capacity as follows:

- (a) A queue reclaims capacity when it has at least one task pending and a fraction of its guaranteed capacity is used by another queue; the scheduler determines the amount of resources to reclaim within the reclaim time for the queue.
- (b) The scheduler kills the tasks that started the last when a queue which has already received Fair Queue resources, it is allowed to reclaim, but its reclaim time is about to expire.

The scheduler can be configured with several properties for each queue using the file *conf/capacity-scheduler.xml*. Queue properties can be defined by concatenating the string *mapred.capacity-scheduler.queue.<queue-name>* with the property name:

⁷ In *Hadoop*, the *JobTracker* and *TaskTracker* daemons handle the processing of *MapReduce* jobs.

.guaranteed capacity—percentage of the number of slots guaranteed to be available for jobs in the queue i .

.reclaim-time-limit—the amount of time, in seconds, before resources distributed to other queues will be reclaimed.

.supports-priority; if true, priorities of jobs will be taken into account in scheduling decisions.

.user-limit-percent; if there is competition for resources each queue enforces a limit on the percentage of resources allocated to a user at any given time. If two users have submitted jobs to a queue, no single user can use more than 50% of the queue resources. If a third user submits a job, no single user can use more than 33% of the queue resources. With 4 or more users, no user can use more than 25% of the queue's resources. A value of 100% implies no user limits are imposed.

9.6 Start-time fair queuing (R)

A hierarchical CPU scheduler for multimedia operating systems was proposed in [205]. The basic idea of the *start-time fair queuing* (SFQ) algorithm is to organize the consumers of the CPU bandwidth in a tree structure. The root node is the processor and the leaves of this tree are the threads of each application. A scheduler acts at each level of the hierarchy; the fraction of the processor bandwidth, B , allocated to the intermediate node i is

$$\frac{B_i}{B} = \frac{w_i}{\sum_{j=1}^n w_j} \quad (9.13)$$

with w_j , $1 \leq j \leq n$, the weight of the n children of node i ; see Fig. 9.3.

When a VM is not active, its bandwidth is reallocated to the other VMs active at the time; when one of the applications of a VM is not active, its allocation is transferred to the other applications running on the same VM. Similarly, if one of the threads of an application is not runnable, then its allocation is transferred to the other threads of the application.

Call $v_a(t)$, $v_b(t)$ the virtual time of threads a , b and $v(t)$ the virtual time of the scheduler at time t . Call q the time quantum of the scheduler, in milliseconds. The threads a and b have their time quantum, q_a and q_b , weighted by w_a and w_b , respectively; thus, in our example, the time quantum of the two threads are q/w_a and q/w_b , respectively. The i -th activation of thread a will start at virtual time S_a^i and will finish at virtual time F_a^i . We call τ^j the real time of the j -th invocation of the scheduler.

An SFQ scheduler follows several rules:

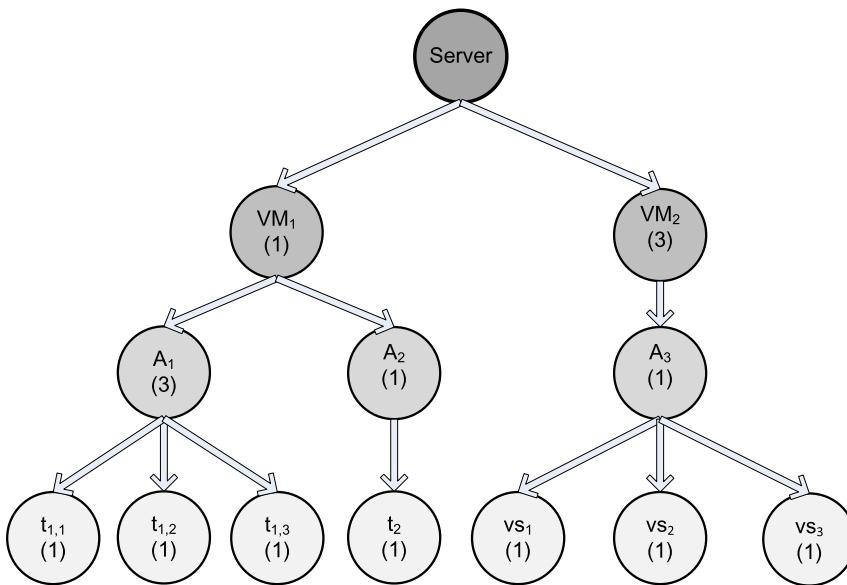
1. Threads are serviced in order of their virtual start up time; ties are broken arbitrarily.
2. The virtual startup time of the i -th activation of thread x is

$$S_x^i(t) = \max \left\{ v(\tau^j), F_x^{(i-1)}(t) \right\} \text{ and } S_x^0 = 0. \quad (9.14)$$

The condition for thread i to be started is that thread $(i - 1)$ has finished and that the scheduler is active.

3. The virtual finish time of the i -th activation of thread x is

$$F_x^i(t) = S_x^i(t) + \frac{q}{w_x}. \quad (9.15)$$

**FIGURE 9.3**

SFQ tree for scheduling when VM_1 and VM_2 run on a powerful server. VM_1 runs two best-effort applications A_1 , with three threads $t_{1,1}$, $t_{1,2}$, and $t_{1,3}$, and A_2 with a single thread t_2 ; VM_2 runs a video-streaming application A_3 with three threads vs_1 , vs_2 , and vs_3 . The weights of VMs, applications, and individual threads are shown in parenthesis.

A thread is stopped when its time quantum has expired; its time quantum is the time quantum of the scheduler divided by the weight of the thread.

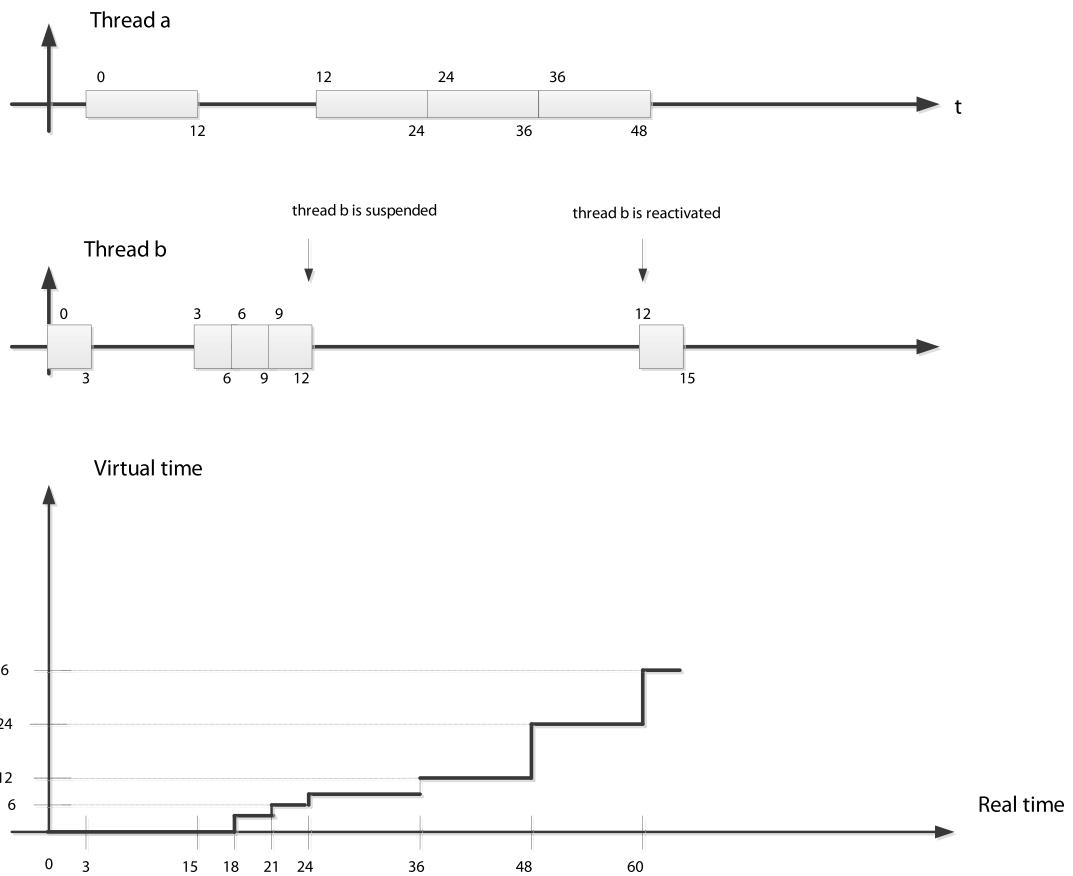
4. The virtual time of all threads is initially zero, $v_x^0 = 0$. The virtual time $v(t)$ at real time t is computed as follows:

$$v(t) = \begin{cases} \text{Virtual start time of the thread in service at time } t, & \text{if CPU busy} \\ \text{Maximum finish virtual time of any thread,} & \text{if CPU idle} \end{cases} \quad (9.16)$$

In this description of the algorithm, we have included the real time t to stress the dependence of all events in virtual time on the real time. To simplify the notation we'll use in our examples the real time as the index of the event, in other words at S_a^6 means the start up time of thread a at real time $t = 6$.

Example. The following example illustrates the application of the SFQ algorithm when there are two threads with the weights $w_a = 1$ and $w_b = 4$ and the time quantum is $q = 12$; see Fig. 9.4.

Initially, $S_a^0 = 0$, $S_b^0 = 0$, $v_a(0) = 0$, and $v_b(0) = 0$. Thread b blocks at time $t = 24$ and wakes up at time $t = 60$. The scheduling decisions at successive intervals of time are:

**FIGURE 9.4**

(Top) Virtual startup time $S_a(t)$ and $S_b(t)$ and virtual finish time $F_a(t)$ and $F_b(t)$ function of real time t for each activation of threads a and b , respectively, are marked at the top and, respectively, at the bottom of the box representing a running thread. (Bottom) Scheduler virtual time $v(t)$ function of the real time.

$t=0$: we have a tie, $S_a^0 = S_b^0$ and arbitrarily thread b is chosen to run first; the virtual finish time of thread b is

$$F_b^0 = S_b^0 + q/w_b = 0 + 12/4 = 3. \quad (9.17)$$

$t=3$: both threads are runnable and thread b was in service, thus, $v(3) = S_b^0 = 0$; then

$$S_b^1 = \max \{v(3), F_b^0\} = \max \{0, 3\} = 3. \quad (9.18)$$

But $S_a^0 < S_b^1$ thus thread a is selected to run. Its virtual finish time is

$$F_a^0 = S_a^0 + q/w_a = 0 + 12/1 = 12. \quad (9.19)$$

$t = 15$: both threads are runnable and thread a was in service at this time thus,

$$v(15) = S_a^0 = 0 \quad \text{and} \quad S_a^1 = \max \{v(15), F_a^0\} = \max \{0, 12\} = 12. \quad (9.20)$$

As $S_b^1 = 3 < 12$, thread b is selected to run; the virtual finish time of thread b is now

$$F_b^1 = S_b^1 + q/w_b = 3 + 12/4 = 6. \quad (9.21)$$

$t = 18$: both threads are runnable, and thread b was in service at this time, thus,

$$v(18) = S_b^1 = 3 \quad \text{and} \quad S_b^2 = \max \{v(18), F_b^1\} = \max \{3, 6\} = 6. \quad (9.22)$$

As $S_b^2 < S_a^1 = 12$, thread b is selected to run again; its virtual finish time is

$$F_b^2 = S_b^2 + q/w_b = 6 + 12/4 = 9. \quad (9.23)$$

$t = 21$: both threads are runnable, and thread b was in service at this time, thus,

$$v(21) = S_b^2 = 6 \quad \text{and} \quad S_b^3 = \max \{v(21), F_b^2\} = \max \{6, 9\} = 9. \quad (9.24)$$

As $S_b^2 < S_a^1 = 12$, thread b is selected to run again; its virtual finish time is

$$F_b^3 = S_b^3 + q/w_b = 9 + 12/4 = 12. \quad (9.25)$$

$t = 24$: Thread b was in service at this time, thus,

$$v(24) = S_b^3 = 9 \quad \text{and} \quad S_b^4 = \max \{v(24), F_b^3\} = \max \{9, 12\} = 12. \quad (9.26)$$

Thread b is suspended till $t = 60$, thus, the thread a is activated; its virtual finish time is

$$F_a^1 = S_a^1 + q/w_a = 12 + 12/1 = 24. \quad (9.27)$$

$t = 36$: thread a was in service, and it is the only runnable thread at this time, thus,

$$v(36) = S_a^1 = 12 \quad \text{and} \quad S_a^2 = \max \{v(36), F_a^1\} = \max \{12, 24\} = 24 \quad \text{then} \quad (9.28)$$

$$F_a^2 = S_a^2 + q/w_a = 24 + 12/1 = 36. \quad (9.29)$$

$t = 48$: thread a was in service, and it is the only runnable thread at this time thus,

$$v(48) = S_a^2 = 24 \quad \text{and} \quad S_a^3 = \max \{v(48), F_a^2\} = \max \{24, 36\} = 36, \quad \text{then} \quad (9.30)$$

$$F_a^3 = S_a^3 + q/w_a = 36 + 12/1 = 48. \quad (9.31)$$

$t = 60$: thread a was in service at this time, thus,

$$v(60) = S_a^3 = 36 \quad \text{and} \quad S_a^4 = \max \{v(60), F_a^3\} = \max \{36, 48\} = 48. \quad (9.32)$$

But now thread b is runnable and $S_b^4 = 12$. Thus, thread b is activated and

$$F_b^4 = S_b^4 + q/w_b = 12 + 12/4 = 15. \quad (9.33)$$

The algorithm allocates CPU fairly when the available bandwidth varies in time and provides throughput, as well as delay guarantees. The algorithm schedules the threads in the order of their virtual startup time, the shortest one first; the length of the time quantum is not required when a thread is scheduled, but only after the thread has finished its current allocation. The overhead of SFQ algorithm is comparable to that of the Solaris scheduling algorithm [205].

9.7 Borrowed virtual time (R)

The objective of the *borrowed virtual time* (BVT) scheduling algorithm is to support low-latency dispatching of real-time applications, as well as a weighted sharing of the CPU among several classes of applications [154]. Like SFQ, the BVT algorithm supports scheduling a mixture of applications, some with hard, some with soft real-time constraints, and applications demanding only a best-effort.

Thread i has an *effective virtual time*, E_i , an *actual virtual time*, A_i , as well as a *virtual time warp*, W_i . The scheduler thread maintains its own *scheduler virtual time* (SVT) defined as the minimum actual virtual time A_j of any thread. The threads are dispatched in the order of their effective virtual time, E_i , a policy called the Earliest Virtual Time (EVT).

Virtual time warp allows a thread to acquire an earlier effective virtual time, in other words, to borrow virtual time from its future CPU allocation. The virtual warp time is enabled when the variable *warpBack* is set; in this case a latency-sensitive thread gains dispatching preference as

$$E_i \leftarrow \begin{cases} A_i & \text{if } \textit{warpBack} = \text{OFF} \\ A_i - W_i & \text{if } \textit{warpBack} = \text{ON} \end{cases} \quad (9.34)$$

The algorithm measures the time in *minimum charging units*, mcu , and uses a time quantum called *context switch allowance* (C), which measures in multiples of mcu , the real time a thread is allowed to run when competing with other threads; typical values for the two quantities are $mcu = 100$ μ sec and $C = 100$ msec. A thread is charged an integer number of mcu . Context switches are triggered by traditional events, the running thread is blocked waiting for an event to occur, the time quantum expires, an interrupt occurs; context switching also occurs when a thread becomes runnable after sleeping.

This policy prevents a thread that has been sleeping for a long time to claim control of the CPU for a longer period of time than it deserves. When thread τ_i becomes runnable after sleeping, its actual virtual time is updated $A_i \leftarrow \max \{A_i, SVT\}$.

If there are no interrupts, threads are allowed to run for the same amount of virtual time. Individual threads have weights; a thread with a larger weight consumes its virtual time more slowly. In practice, each thread τ_i maintains a constant k_i and uses its weight w_i to compute the amount Δ used to advance

Table 9.2 The real time and the effective virtual time $E_a(t)$ and $E_b(t)$ at time of a context switch. There is no time warp, thus, the effective virtual time is the same as the actual virtual time. At time $t = 0$, $E_a(0) = E_b(0) = 0$ and we choose thread a to run. CS—context switch; rt—real time; RT—running thread.

CS	rt	RT	Effective virtual time of running thread	Thread running next
1	$t = 2$	a	$E_a(2) = A_a(2) = A_a(0) + \Delta/3 = 30$	$b, E_b(2) = 0 < E_a(2) = 30$
2	$t = 5$	b	$E_b(5) = A_b(5) = A_b(0) + \Delta = 90$	$a, E_a(5) = 30 < E_b(5) = 90$
3	$t = 11$	a	$E_a(11) = A_a(11) = A_a(2) + \Delta = 120$	$b, E_b(11) = 90 < E_a(11) = 120$
4	$t = 14$	b	$E_b(14) = A_b(14) = A_b(5) + \Delta = 180$	$a, E_a(14) = 120 < E_b(14) = 180$
5	$t = 20$	a	$E_a(20) = A_a(20) = A_a(11) + \Delta = 210$	$b, E_b(20) = 180 < E_a(20) = 210$
6	$t = 23$	b	$E_b(23) = A_b(23) = A_b(14) + \Delta = 270$	$a, E_a(23) = 210 < E_b(23) = 270$
7	$t = 29$	a	$E_a(29) = A_a(29) = A_a(20) + \Delta = 300$	$b, E_b(29) = 270 < E_a(29) = 300$
8	$t = 32$	b	$E_b(32) = A_b(32) = A_b(23) + \Delta = 360$	$a, E_a(32) = 300 < E_b(32) = 360$
9	$t = 38$	a	$E_a(38) = A_a(38) = A_a(29) + \Delta = 390$	$b, E_b(11) = 360 < E_a(11) = 390$
10	$t = 41$	b	$E_b(41) = A_b(41) = A_b(32) + \Delta = 450$	$a, E_a(41) = 390 < E_b(41) = 450$

its actual virtual time upon completion of a run $A_i \leftarrow A_i + \Delta$. Given two threads a and b , then $\Delta = k_a/w_a = k_b/w_b$.

The BVT policy requires that every time the actual virtual time is updated, a context switch from the current running thread τ_i to a thread τ_j occurs if

$$A_j \leq A_i - \frac{C}{w_i}. \quad (9.35)$$

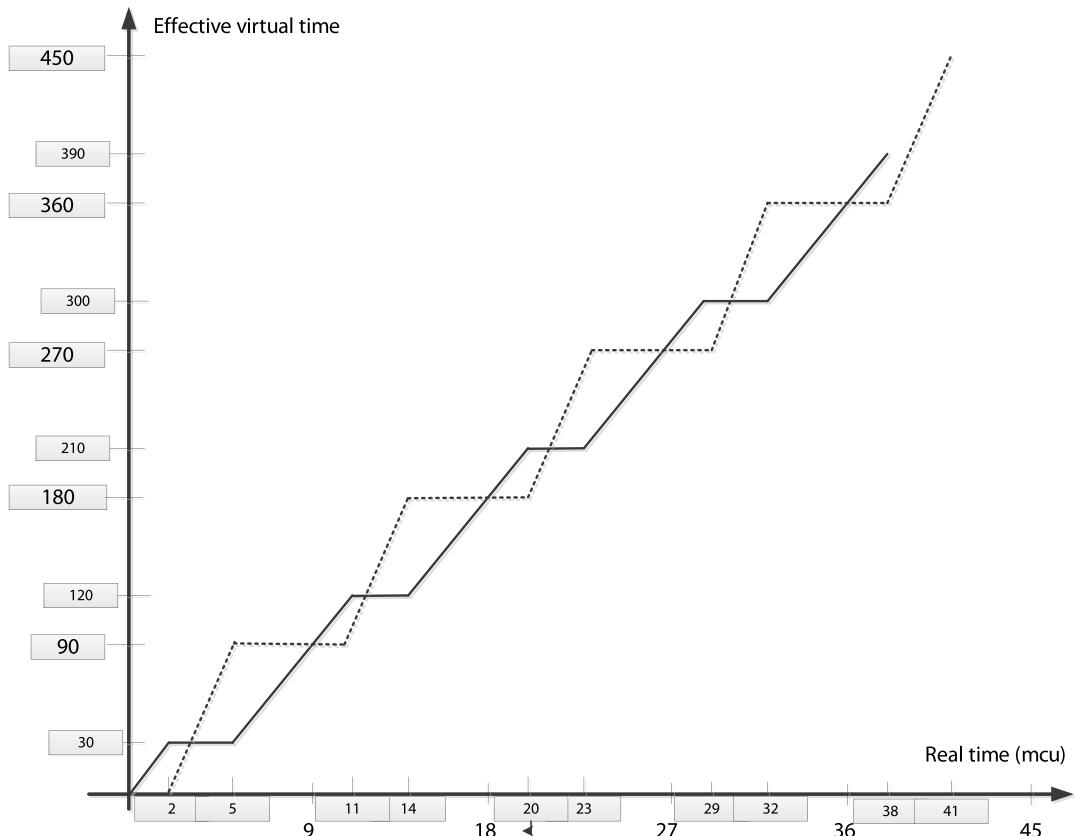
Example 1. Application of BVT algorithm for scheduling two threads a and b of best-effort applications. The first thread has a weight twice of the weight of the second, $w_a = 2w_b$; when $k_a = 180$ and $k_b = 90$, then $\Delta = 90$. We consider periods of real time allocation of $C = 9 \text{ mcu}$; the two threads a and b are allowed to run for $2C/3 = 6 \text{ mcu}$ and $C/3 = 3 \text{ mcu}$, respectively. Threads a and b are activated at times

$$\begin{aligned} a : 0, 5, 5 + 9 = 14, 14 + 9 = 23, 23 + 9 = 32, 32 + 9 = 41, \dots \text{ and} \\ b : 2, 2 + 9 = 11, 11 + 9 = 20, 20 + 9 = 29, 29 + 9 = 38, \dots \end{aligned}$$

Context switches occur at real times 2, 5, 11, 14, 20, 23, 29, 32, 38, 41, The time is expressed in units of mcu . The initial run is a shorter one, consisting of only 3 mcu ; a context switch occurs when a , which runs first, exceeds b by 2 mcu .

Table 9.2 shows the effective virtual time of the two threads at the time of each context switch. At that moment, the actual virtual time is incremented by an amount equal to Δ if the thread was allowed to run for its time allocation. The scheduler compares the effective virtual time of the threads and runs first the one with the minimum effective virtual time.

Fig. 9.5 displays the effective virtual time and the real time of the threads a and b . When a thread is running, its effective virtual time increases as the real time increase; a running thread appears as a diagonal line. When a thread is runnable, but not running, its effective virtual time is constant; a runnable period is displayed as a horizontal line. We see that the two threads are allocated equal amounts of virtual time, but the thread a , with a larger weight, consumes its real time more slowly.

**FIGURE 9.5**

Example 1. The effective virtual time and the real time of the threads a (solid line) and b (dotted line) with weights $w_a = 2w_b$ when the actual virtual time is incremented in steps of $\Delta = 90$ mcu. The real time the two threads are allowed to use the CPU is proportional with their weights; the virtual times are equal but thread a consumes it more slowly. There is no time warp, the threads are dispatched based on their actual virtual time.

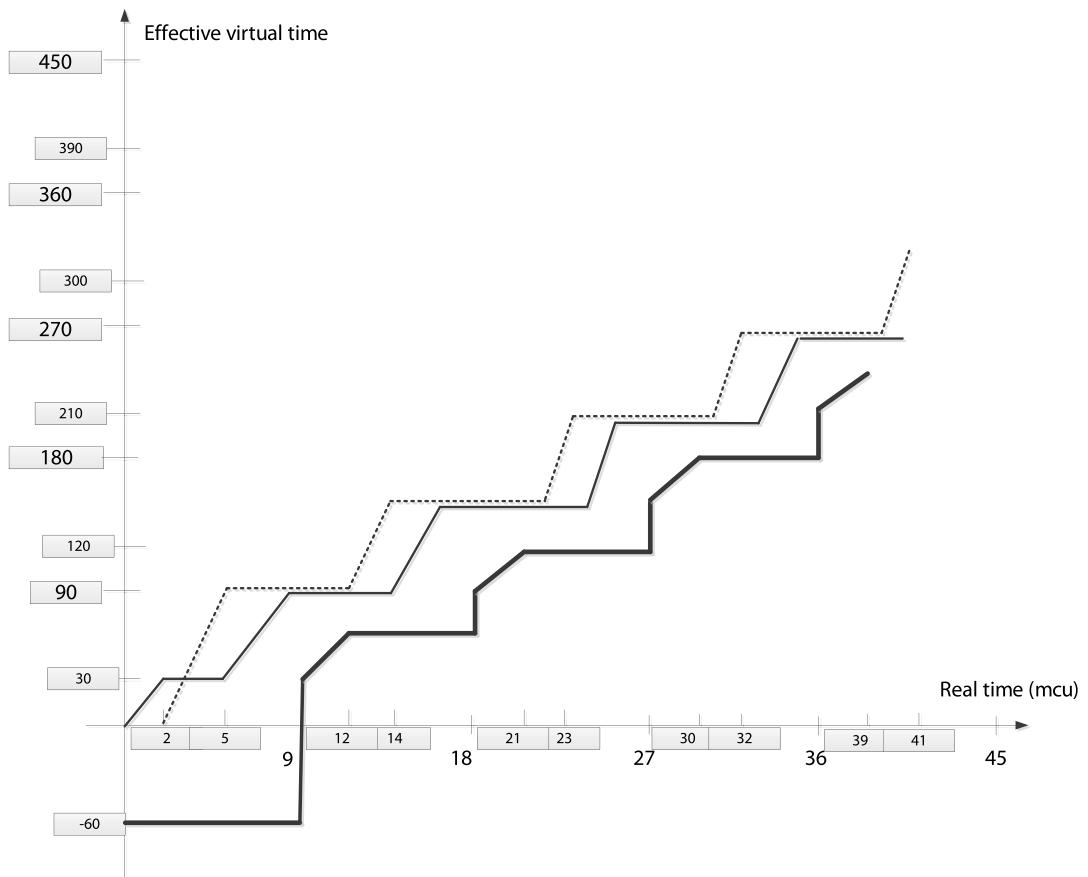
Example 2. Next, we consider the previous example but, this time, there is an additional thread, c , with real-time constraints, which wakes up at time $t = 9$ and then periodically at times $t = 18, 27, 36, \dots$ for 3 units of time.

Table 9.3 summarizes the evolution of the system when the real-time application thread c competes with the two best-effort threads a and b . Context switches occur now at real times $t = 2, 5, 9, 12, 14, 18, 21, 23, 27, 30, 32, 36, 39, 41, \dots$. Context switches at times $t = 9, 18, 27, 36, \dots$ are triggered by the waking up of the thread c which preempts the currently running thread. At $t = 9$ the time warp $W_c = -60$ gives priority to thread c . Indeed, $E_c(9) = A_c(9) - W_c = 0 - 60 = -60$ compared with $E_a(9) = 90$ and $E_b(9) = 90$. The same conditions occur every time the real-time thread wakes up. The best-effort application threads have the same effective virtual time when the real-time

Table 9.3 A real time thread c with a time warp $W_c = -60$ is waking up periodically at times $t = 18, 27, 36, \dots$ for 3 units of time and is competing with the two best-effort threads a and b . The real time and the effective virtual time of the three threads of each context switch are shown.

Context switch	Real time	Running thread	Effective virtual time of the running thread
1	$t = 2$	a	$E_a(2) = A_a(2) = A_a(0) + \Delta/3 = 0 + 90/3 = 30$
2	$t = 5$	b	$E_b^1 = A_b^1 = A_b^0 + \Delta = 0 + 90 = 90 \Rightarrow a \text{ runs next}$
3	$t = 9$	a	$c \text{ wakes up}$ $E_a^1 = A_a^1 + 2\Delta/3 = 30 + (-60) = 90$ $[E_a(9), E_b(9), E_c(9)] = (90, 90, -60) \Rightarrow c \text{ runs next}$
4	$t = 12$	c	$SVT(12) = \min(90, 90)$ $E_c^s(12) = SVT(12) + W_c = 90 + (-60) = 30$ $E_c(12) = E_c^s(12) + \Delta/3 = 30 + 30 = 60 \Rightarrow b \text{ runs next}$
5	$t = 14$	b	$E_b^2 = A_b^2 = A_b^1 + 2\Delta/3 = 90 + 60 = 150 \Rightarrow a \text{ runs next}$
6	$t = 18$	a	$c \text{ wakes up}$ $E_a^3 = A_a^3 = A_a^2 + 2\Delta/3 = 90 + 60 = 150$ $[E_a(18), E_b(18), E_c(18)] = (150, 150, 60) \Rightarrow c \text{ runs next}$
7	$t = 21$	c	$SVT = \min(150, 150)$ $E_c^s(21) = SVT + W_c = 150 + (-60) = 90$ $E_c(21) = E_c^s(21) + \Delta/3 = 90 + 30 = 120 \Rightarrow b \text{ runs next}$
8	$t = 23$	b	$E_b^3 = A_b^3 = A_b^2 + 2\Delta/3 = 150 + 60 = 210 \Rightarrow a \text{ runs next}$
9	$t = 27$	a	$c \text{ wakes up}$ $E_a^4 = A_a^4 = A_a^3 + 2\Delta/3 = 150 + 60 = 210$ $[E_a(27), E_b(27), E_c(27)] = (210, 210, 120) \Rightarrow c \text{ runs next}$
10	$t = 30$	c	$SVT = \min(210, 210)$ $E_c^s(30) = SVT + W_c = 210 + (-60) = 150$ $E_c(30) = E_c^s(30) + \Delta/3 = 150 + 30 = 180 \Rightarrow b \text{ runs next}$
11	$t = 32$	b	$E_b^4 = A_b^4 = A_b^3 + 2\Delta/3 = 210 + 60 = 270 \Rightarrow a \text{ runs next}$
10	$t = 36$	a	$c \text{ wakes up}$ $E_a^5 = A_a^5 = A_a^4 + 2\Delta/3 = 210 + 60 = 270$ $[E_a(36), E_b(36), E_c(36)] = (270, 270, 180) \Rightarrow c \text{ runs next}$
12	$t = 39$	c	$SVT = \min(270, 270)$ $E_c^s(39) = SVT + W_c = 270 + (-60) = 210$ $E_c(39) = E_c^s(39) + \Delta/3 = 210 + 30 = 240 \Rightarrow b \text{ runs next}$
13	$t = 41$	b	$E_b^5 = A_b^5 = A_b^4 + 2\Delta/3 = 270 + 60 = 330 \Rightarrow a \text{ runs next}$

application thread finishes and the scheduler chooses b to be dispatched first. We should also notice that the ratio of real times used by a and b is the same, as $w_a = 2w_b$. Fig. 9.6 shows the effective virtual times for the three threads a , b , and c . Every time when thread c wakes up, it preempts the current running thread, and it is immediately scheduled to run.

**FIGURE 9.6**

Example 2. The effective virtual time and the real time of the threads a (solid line), b (dotted line), and the c with real-time constraints (thick solid line). c wakes up periodically at times $t = 9, 18, 27, 36, \dots$, is active for 3 units of time and has a time warp of 60 mcu .

9.8 Cloud scheduling subject to deadlines (R)

Often, a service level agreement specifies the time when the results of computations done on the cloud should be available. This motivates us to examine cloud scheduling subject to deadlines, a topic drawing from a vast body of literature devoted to real-time applications.

Task characterization and deadlines. Real-time applications involve periodic or aperiodic tasks with deadlines. A task is characterized by a tuple (A_i, σ_i, D_i) , where A_i is the arrival time, $\sigma_i > 0$ is the data size of the task, and D_i is the *relative deadline*. Instances of a *periodic task*, Π_i^q , with period q are identical, $\Pi_i^q \equiv \Pi^q$, and arrive at times $A_0, A_1, \dots, A_i, \dots$, with $A_{i+1} - A_i = q$.

The deadlines satisfy the constraint $D_i \leq A_{i+1}$ and, generally, the data size is the same, $\sigma_i = \sigma$. The individual instances of *aperiodic tasks*, Π_i , are different, their arrival times A_i are generally uncorrelated, and the amount of data σ_i is different for different instances. The *absolute deadline* for the aperiodic task Π_i is $(A_i + D_i)$.

We distinguish *hard deadlines* from *soft deadlines*. In the first case, if the task is not completed by the deadline other tasks which depend on it may be affected and there are penalties. A hard deadline is strict and expressed precisely as milliseconds, or possibly seconds.

Soft deadlines are more of a guideline, and, in general, there are no penalties; soft deadlines can be missed by fractions of the units used to express them, e.g., minutes if the deadline is expressed in hours, or hours if the deadlines are expressed in days. The scheduling of tasks on a cloud is generally subject to soft deadlines, though, occasionally, applications with hard deadlines may be encountered.

System model. We consider only aperiodic tasks with arbitrarily divisible workloads. The application runs on a partition of a cloud, a virtual cloud with a *head node* called S_0 and n *worker nodes* S_1, S_2, \dots, S_n . The system is homogeneous, all workers are identical, and the communication time from the head node to every worker node is the same. The head node distributes the workload to worker nodes and this distribution is done sequentially. In this context there are two important problems:

1. The order of execution of the tasks Π_i .
2. The workload partitioning and the task mapping to worker nodes.

Scheduling policies. The most common scheduling policies used to determine the order of execution of the tasks are:

- FIFO—First-In-First-Out, the tasks are scheduled for execution in order of their arrival.
- EDF—Earliest Deadline First, the task with the earliest deadline is scheduled first.
- MWF—Maximum Workload Derivative First.

The *workload derivative* $DC_i(n^{min})$ of a task Π_i when n^{min} nodes are assigned to the application is defined as

$$DC_i(n^{min}) = W_i(n_i^{min} + 1) - W_i(n_i^{min}), \quad (9.36)$$

with $W_i(n)$ the workload allocated to task Π_i when n nodes of the cloud are available. The MWF policy requires that:

1. The tasks are scheduled in the order of their derivatives, the one with the highest derivative DC_i first.
2. The number n of nodes assigned to the application is kept to a minimum, n_i^{min} .

Two workload partitioning and task mappings to worker nodes are discussed next, the optimal partitioning and the equal partitioning. In our discussion, we use the derivations and some of the notations in [308]. These notations are summarized in Table 9.4.

Optimal Partitioning Rule (OPR). The workload is partitioned to ensure the earliest possible completion time. The optimality of OPR scheduling requires all tasks to finish execution at the same time.

Head node, S_0 , distributes data sequentially to individual worker nodes. The workload assigned to worker node S_i is $\alpha_i \sigma$. The time to deliver input data to worker node S_i is $\Gamma_i = (\alpha_i \times \sigma) \times \tau$, $1 \leq i \leq n$.

Table 9.4 The parameters used for scheduling with deadlines.

Name	Description
Π_i	the aperiodic tasks with arbitrary divisible load of an application \mathcal{A}
A_i	arrival time of task Π_i
D_i	the relative deadline of task Π_i
σ_i	the workload allocated to task Π_i
S_0	head node of the virtual cloud allocated to \mathcal{A}
S_i	worker nodes $1 \leq i \leq n$ of the virtual cloud allocated to \mathcal{A}
σ	total workload for application \mathcal{A}
n	number of nodes of the virtual cloud allocated to application \mathcal{A}
n^{min}	minimum number of nodes of the virtual cloud allocated to application \mathcal{A}
$\mathcal{E}(n, \sigma)$	the execution time required by n worker nodes to process the workload σ
τ	time for transferring a unit of workload from the head node S_0 to worker S_i
ρ	time for processing a unit of workload
α	the load distribution vector $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)$
$\alpha_i \times \sigma$	the fraction of the workload allocated to worker node S_i
Γ_i	time to transfer the data to worker S_i , $\Gamma_i = \alpha_i \times \sigma \times \tau$, $1 \leq i \leq n$
Δ_i	time the worker S_i needs to process a unit of data, $\Delta_i = \alpha_i \times \sigma \times \rho$, $1 \leq i \leq n$
t_0	start time of the application \mathcal{A}
A	arrival time of the application \mathcal{A}
D	deadline of application \mathcal{A}
$C(n)$	completion time of application \mathcal{A}

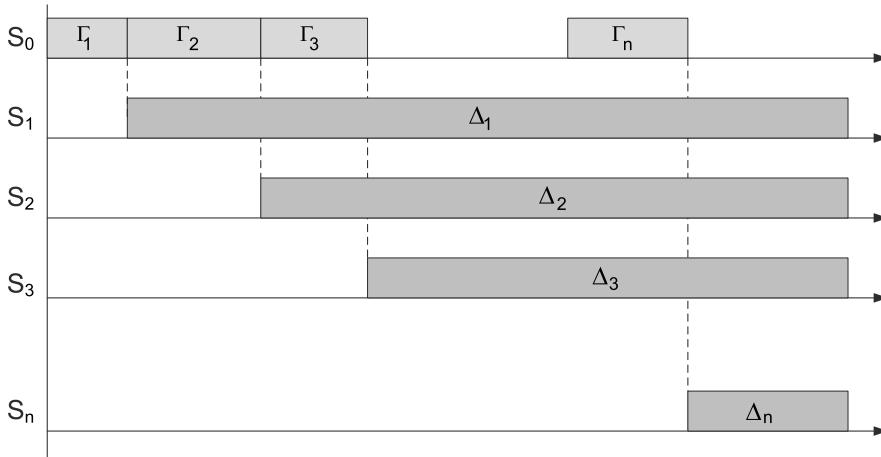
Worker node S_i starts processing the data as soon as the transfer is complete. The processing time of worker node S_i is $\Delta_i = (\alpha_i \times \sigma) \times \rho$, $1 \leq i \leq n$.

The timing diagram in Fig. 9.7 allows us to determine the execution time $\mathcal{E}(n, \sigma)$ for the OPR as

$$\begin{aligned}
 \mathcal{E}(1, \sigma) &= \Gamma_1 + \Delta_1 \\
 \mathcal{E}(2, \sigma) &= \Gamma_1 + \Gamma_2 + \Delta_2 \\
 \mathcal{E}(3, \sigma) &= \Gamma_1 + \Gamma_2 + \Gamma_3 + \Delta_3 \\
 &\vdots \\
 \mathcal{E}(n, \sigma) &= \Gamma_1 + \Gamma_2 + \Gamma_3 + \dots + \Gamma_n + \Delta_n.
 \end{aligned} \tag{9.37}$$

We substitute the expressions of Γ_i , Δ_i , $1 \leq i \leq n$ and rewrite these equations as

$$\begin{aligned}
 \mathcal{E}(1, \sigma) &= \alpha_1 \times \sigma \times \tau + \alpha_1 \times \sigma \times \rho \\
 \mathcal{E}(2, \sigma) &= \alpha_1 \times \sigma \times \tau + \alpha_2 \times \sigma \times \tau + \alpha_2 \times \sigma \times \rho \\
 \mathcal{E}(3, \sigma) &= \alpha_1 \times \sigma \times \tau + \alpha_2 \times \sigma \times \tau + \alpha_3 \times \sigma \times \tau + \alpha_3 \times \sigma \times \rho \\
 &\vdots \\
 \mathcal{E}(n, \sigma) &= \alpha_1 \times \sigma \times \tau + \alpha_2 \times \sigma \times \tau + \alpha_3 \times \sigma \times \tau + \dots + \alpha_n \times \sigma \times \tau + \alpha_n \times \sigma \times \rho.
 \end{aligned} \tag{9.38}$$

**FIGURE 9.7**

OPR timing diagram. The algorithm requires worker nodes to complete execution at the same time.

From the first two equations, we find the relation between α_1 and α_2 as

$$\alpha_1 = \frac{\alpha_2}{\beta} \quad \text{with} \quad \beta = \frac{\rho}{\tau + \rho}, \quad 0 \leq \beta \leq 1. \quad (9.39)$$

This implies that $\alpha_2 = \beta \times \alpha_1$; it is easy to see that in the general case

$$\alpha_i = \beta \times \alpha_{i-1} = \beta^{i-1} \times \alpha_1. \quad (9.40)$$

But α_i are the components of the load distribution vector; thus,

$$\sum_{i=1}^n \alpha_i = 1. \quad (9.41)$$

Next, we substitute the values of α_i and obtain the expression for α_1 :

$$\alpha_1 + \beta \times \alpha_1 + \beta^2 \times \alpha_1 + \beta^3 \times \alpha_1 \dots \beta^{n-1} \times \alpha_1 = 1 \quad \text{or} \quad \alpha_1 = \frac{1 - \beta}{1 - \beta^n}. \quad (9.42)$$

We have now determined the load distribution vector and we can now determine the execution time as

$$\mathcal{E}(n, \sigma) = \alpha_1 \times \sigma \times \tau + \alpha_1 \times \sigma \times \rho = \frac{1 - \beta}{1 - \beta^n} \sigma (\tau + \rho). \quad (9.43)$$

Call $C^A(n)$ the completion time of an application $A = (A, \sigma, D)$, which starts processing at time t_0 and runs on n worker nodes; then

$$C^A(n) = t_0 + \mathcal{E}(n, \sigma) = t_0 + \frac{1 - \beta}{1 - \beta^n} \sigma (\tau + \rho). \quad (9.44)$$

The application meets its deadline if and only if

$$C^A(n) \leq A + D, \quad (9.45)$$

or

$$t_0 + \mathcal{E}(n, \sigma) = t_0 + \frac{1 - \beta}{1 - \beta^n} \sigma(\tau + \rho) \leq A + D. \quad (9.46)$$

But $0 < \beta < 1$; thus, $1 - \beta^n > 0$, and it follows that

$$(1 - \beta)\sigma(\tau + \rho) \leq (1 - \beta^n)(A + D - t_0). \quad (9.47)$$

The application can meet its deadline only if $(A + D - t_0) > 0$, and under this condition, the inequality (9.47) becomes

$$\beta^n \leq \gamma \quad \text{with} \quad \gamma = 1 - \frac{\sigma \times \tau}{A + D - t_0}. \quad (9.48)$$

If $\gamma \leq 0$, there is not enough time even for data distribution and the application should be rejected. When $\gamma > 0$, then $n \geq \frac{\ln \gamma}{\ln \beta}$. Thus, the minimum number of nodes for the OPR strategy is

$$n^{min} = \lceil \frac{\ln \gamma}{\ln \beta} \rceil. \quad (9.49)$$

Equal Partitioning Rule (EPR). EPR assigns an equal workload to individual worker nodes, $\alpha_i = 1/n$. The workload allocated to worker node S_i is σ/n . The head node, S_0 , distributes sequentially the data to individual worker nodes. The time to deliver the input data to S_i is $\Gamma_i = (\sigma/n) \times \tau$, $1 \leq i \leq n$. Worker node S_i starts processing the data as soon as the transfer is complete. The processing time for node S_i is $\Delta_i = (\sigma/n) \times \rho$, $1 \leq i \leq n$.

From the diagram in Fig. 9.8, we see that

$$\mathcal{E}(n, \sigma) = \sum_{i=1}^n \Gamma_i + \Delta_n = n \times \frac{\sigma}{n} \times \tau + \frac{\sigma}{n} \times \rho = \sigma \times \tau + \frac{\sigma}{n} \times \rho. \quad (9.50)$$

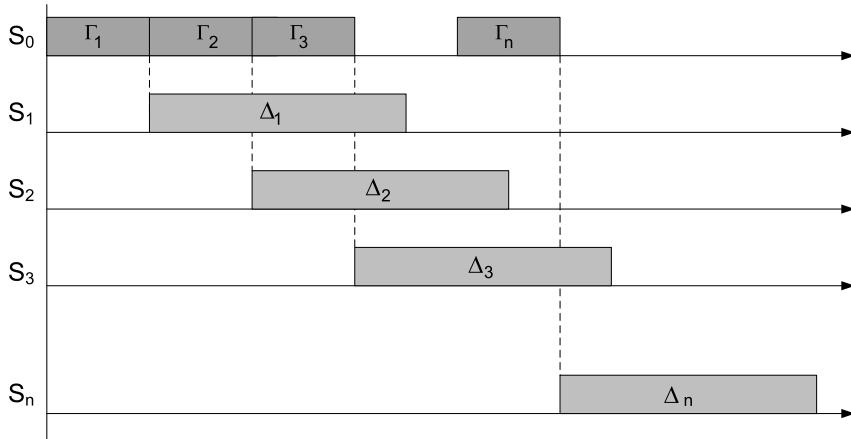
The condition for meeting the deadline, $C^A(n) \leq A + D$, leads to

$$t_0 + \sigma \times \tau + \frac{\sigma}{n} \times \rho \leq A + D \quad \text{or} \quad n \geq \frac{\sigma \times \rho}{A + D - t_0 - \sigma \times \tau}. \quad (9.51)$$

Thus,

$$n^{min} = \lceil \frac{\sigma \times \rho}{A + D - t_0 - \sigma \times \tau} \rceil. \quad (9.52)$$

The pseudocode for a general schedulability test for FIFO, EDF, and MWF scheduling policies, for two node allocation policies, MN (minimum number of nodes) and AN (all nodes), and for OPR and EPR partitioning rules is given in [308]. The same reference reports on a simulation study for ten algorithms.

**FIGURE 9.8**

EPR timing diagram.

The generic format of the names of the algorithms is *Sp-No-Pa* with Sp = FIFO/EDF/MWF, No = MN/AN, and Pa = OPR/EPR. For example, MWF-MN-OPR uses MWF scheduling, minimum number of nodes, and OPR partitioning. The relative performance of the algorithms depends on the relations between the unit cost of communication τ and the unit cost of computing ρ .

9.9 MapReduce application scheduling subject to deadlines (R)

We now discuss cloud scheduling of MapReduce applications subject to deadlines. Several options for scheduling Apache Hadoop, an open source implementation of the MapReduce algorithm are: FIFO, the Fair Scheduler [532], the Capacity Scheduler, and the Dynamic Proportional Scheduler [434].

A recent paper [267] applies the deadline scheduling framework discussed in Section 9.8 to Hadoop tasks. Table 9.5 summarizes the notations used for this analysis. The term *slots* is equivalent with *nodes* and means the number of instances.

We make two assumptions for our initial derivation:

- The system is homogeneous, ρ_m and ρ_r , the cost of processing a unit data by the Map and the Reduce task, respectively, are the same for all servers.
- Load equipartition.

Under these conditions, the duration of the job J with input of size σ is

$$\mathcal{E}(n_m, n_r, \sigma) = \sigma \left[\frac{\rho_m}{n_m} + \phi \left(\frac{\rho_r}{n_r} + \tau \right) \right]. \quad (9.53)$$

Table 9.5 Parameters used for Hadoop scheduling with deadlines.

Name	Description
Q	the query $Q = (A, \sigma, D)$
A	arrival time of query Q
D	deadline of query Q
Π_m^i	a map task, $1 \leq i \leq u$
Π_r^j	a reduce task, $1 \leq j \leq v$
J	the job to perform the query $Q = (A, \sigma, D)$, $J = (\Pi_m^1, \Pi_m^2, \dots, \Pi_m^u, \Pi_r^1, \Pi_r^2, \dots, \Pi_r^v)$
τ	cost for transferring a data unit
ρ_m	map task time for processing a unit data
ρ_r	reduce task time for processing a unit data
n_m	number of map slots
n_r	number of reduce slots
n_m^{min}	minimum number of slots for the map task
n	total number of slots, $n = n_m + n_r$
t_m^0	start time of the map task
t_r^{max}	maximum value for the start time of the reduce task
α	map distribution vector; the EPR strategy is used and, $\alpha_i = 1/u$
ϕ	filter ratio, the fraction of the input produced as output by the map process

Thus, the condition that the query $Q = (A, \sigma, D)$ with arrival time A meets the deadline D can be expressed as

$$t_m^0 + \sigma \left[\frac{\rho_m}{n_m} + \phi \left(\frac{\rho_r}{n_r} + \tau \right) \right] \leq A + D. \quad (9.54)$$

It follows immediately that the maximum value for the startup time of the *reduce* task is

$$t_r^{max} = A + D - \sigma \phi \left(\frac{\rho_r}{n_r} + \tau \right). \quad (9.55)$$

We now plug the expression of the maximum value for the startup time of the *reduce* task in the condition to meet the deadline

$$t_m^0 + \sigma \frac{\rho_m}{n_m} \leq t_r^{max}. \quad (9.56)$$

It follows immediately that n_m^{min} , the minimum number of slots for the *map* task, satisfies the condition

$$n_m^{min} \geq \frac{\sigma \rho_m}{t_r^{max} - t_m^0}, \quad \text{thus, } n_m^{min} = \lceil \frac{\sigma \rho_m}{t_r^{max} - t_m^0} \rceil. \quad (9.57)$$

The assumption of homogeneity of the servers can be relaxed and assume that individual servers have different costs for processing a unit workload $\rho_m^i \neq \rho_m^j$ and $\rho_r^i \neq \rho_r^j$. In this case, we can use the minimum values $\rho_m = \min \rho_m^i$ and $\rho_r = \min \rho_r^i$ in the expression we derived.

A Constraints Scheduler based on this analysis and an evaluation of the effectiveness of this scheduler are presented in [267].

9.10 Resource bundling; combinatorial auctions for cloud resources

Resources in a cloud are allocated in *bundles*; users get maximum benefit from a specific combination of resources. Indeed, along with CPU cycles, an application needs specific amounts of main memory, disk space, network bandwidth, and so on. Resource bundling complicates traditional resource allocation models and has generated an interest in economic models and, in particular, in auction algorithms. In the context of cloud computing, an auction is the allocation of resources to the highest bidder.

Combinatorial auctions. Auctions in which participants can bid on combinations of items or *packages* are called *combinatorial auctions* [118]; such auctions provide a relatively simple, scalable, and tractable solution to cloud resource allocation. Two recent combinatorial auction algorithms are the *Simultaneous Clock Auction* [36] and the *Clock Proxy Auction* [37]; the algorithm discussed in this section and introduced in [456] is called *Ascending Clock Auction (ASCA)*. In all these algorithms, the current price for each resource is represented by a “clock” seen by all participants at the auction.

We consider a strategy when prices and allocation are set as a result of an auction; in this auction, users provide bids for desirable bundles and the price they are willing to pay. We assume a population of U users, $u = \{1, 2, \dots, U\}$, and R resources, $r = \{1, 2, \dots, R\}$. The bid of user u is $\mathcal{B}_u = \{\mathcal{Q}_u, \pi_u\}$ with $\mathcal{Q}_u = (q_u^1, q_u^2, q_u^3, \dots)$ an R -component vector.

Each element of this vector, q_u^i , represents a bundle of resources a user u would accept and, in return, pay the total price π_u . Each vector component q_u^i is either a positive quantity and encodes the quantity of a resource desired, or if negative, it encodes the quantity of the resource offered. A user expresses her desires as an *indifference set* $\mathcal{I} = (q_u^1 \text{ XOR } q_u^2 \text{ XOR } q_u^3 \text{ XOR } \dots)$.

The final auction prices for individual resources are given by the vector $p = (p^1, p^2, \dots, p^R)$ and the amounts of resources allocated to user u are $x_u = (x_u^1, x_u^2, \dots, x_u^R)$. Thus, the expression $[(x_u)^T p]$ represents the total price paid by user u for the bundle of resources if the bid is successful at time T . The scalar $[\min_{q \in \mathcal{Q}_u} (q^T p)]$ is the final price established through the bidding process.

The bidding process aims to optimize an *objective function* $f(x, p)$. This function could be tailored to measure the net value of all resources traded, or it can measure the *total surplus*, the difference between the maximum amount the users are willing to pay minus the amount they pay. Other optimization function could be considered for a specific system, e.g., the minimization of energy consumption, or of the security risks.

Pricing and allocation algorithms. A pricing and allocation algorithm partitions the set of users in two disjoint sets, winners and losers, denoted as \mathcal{W} and \mathcal{L} , respectively; the algorithm should:

1. Be computationally tractable. Traditional combinatorial auction algorithms such as Vickrey-Clarke-Groves (VCG) fail this criteria; they are not computationally tractable.
2. Scale well. Given the scale of the system and the number of requests for service, scalability is a necessary condition.
3. Be objective; partitioning in winners and losers should only be based on the price π_u of a user’s bid; if the price exceeds the threshold then the user is a winner; otherwise, the user is a loser.

Table 9.6 The constraints for a combinatorial auction algorithm.

$x_u \in \{0 \cup \mathcal{Q}_u\}, \forall u$	a user gets all resources or nothing
$\sum_u x_u \leq 0$	final allocation leads to a net surplus of resources
$\pi_u \geq (x_u)^T p, \forall u \in \mathcal{W}$	auction winners are willing to pay the final price
$(x_u)^T p = \min_{q \in \mathcal{Q}_u} (q^T p), \forall u \in \mathcal{W}$	winners get the cheapest bundle in \mathcal{I}
$\pi_u < \min_{q \in \mathcal{Q}_u} (q^T p), \forall u \in \mathcal{L}$	the bids of the losers are below the final price
$p \geq 0$	prices must be non-negative

4. Be fair. Make sure that the prices are *uniform*; all winners within a given resource pool pay the same price.
5. Indicate clearly at the end of the auction the unit prices for each resource pool.
6. Indicate clearly to all participants the relationship between the supply and the demand in the system.

The function to be maximized is $\max_{x, p} f(x, p)$. The constraints in Table 9.6 correspond to our intuition: (a) the first one states that a user either gets one of the bundles it has opted for or nothing, no partial allocation is acceptable; (b) the second one expresses the fact that the system awards only available resources, only offered resources can be allocated; (c) the third one is that the bid of the winners exceeds the final price; (d) the fourth one states that the winners get the least expensive bundles in their indifference set; (e) the fifth one states that losers bid below the final price; (f) finally, the last one states that all prices are positive numbers.

Ascending Clock Auction combinatorial auction algorithm. Informally, the participants at the auction based on the ASCA algorithm [456] specify the resource and the quantities of that resource offered or desired at the price listed for that time slot. Then the *excess vector* is computed as follows:

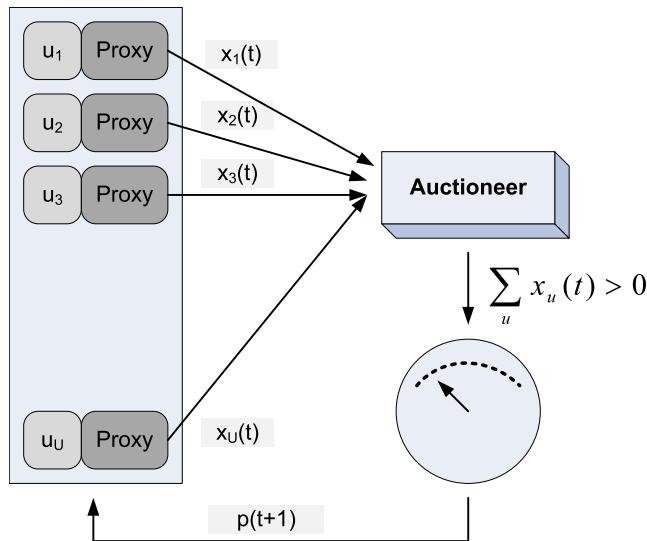
$$z(t) = \sum_u x_u(t) \quad (9.58)$$

If all its components are negative, then auction stops; negative components mean that the demand does not exceed the offer. If the demand is larger than the offer, $z(t) \geq 0$, then the auctioneer increases the price for items with a positive excess demand and solicits bids at the new price. The algorithm satisfies conditions 1 through 6. All users discover the price at the same time and pay or receive a “fair” payment relative to uniform resource prices, the computation is tractable, and the execution time is linear in the number of participants at the auction and the number of resources. The computation is robust, generates plausible results, regardless of the initial parameters of the system.

There is a slight complication as the algorithm involves user bidding in multiple rounds. To address this problem the user proxies automatically adjust their demands on behalf of the actual bidders, as shown in Fig. 9.9. These proxies can be modeled as functions which compute the “best bundle” from each \mathcal{Q}_u set given the current price

$$\mathcal{Q}_u = \begin{cases} \hat{q}_u & \text{if } \hat{q}_u^T p \leq \pi_u \text{ with } \hat{q}_u \in \arg \min(q_u^T p) \\ 0 & \text{otherwise} \end{cases} \quad (9.59)$$

In this algorithm, $g(x(t), p(t))$ is the function for setting the price increase. This function can be correlated with the excess demand $z(t)$ as in $g(x(t), p(t)) = \alpha z(t)^+$ (the notation x^+ means $\max(x, 0)$)

**FIGURE 9.9**

The schematics of the ASCA algorithm; to allow for a single-round auction, users are represented by proxies which place the bids $x_u(t)$. The auctioneer determines if there is an excess demand and, in that case, it raises the price of resources for which the demand exceeds the supply and requests new bids.

with α a positive number. An alternative is to ensure that the price does not increase by an amount larger than δ ; in that case $g(x(t), p(t)) = \min(\alpha z(t)^+, \delta e)$ with $e = (1, 1, \dots, 1)$ is an R -dimensional vector and minimization is done componentwise.

The input to the ASCA algorithm: U users, R resources, \bar{p} the starting price, and the update increment function, $g : (x, p) \mapsto \mathbb{R}^R$. The pseudo code of the algorithm is:

Pseudo code for the ASCA algorithm

```

1 set  $t = 0$ ,  $p(0) = \bar{p}$ 
2   loop
3     collect bids  $x_u(t) = \mathcal{G}_u(p(t)) \ \forall u$ 
4     calculate excess demand  $z(t) = \sum_u x_u(t)$ 
5     if  $z(t) < 0$  then
6       break
7     else
8       update prices  $p(t + 1) = p(t) + g(x(t), p(t))$ 
9        $t \leftarrow t + 1$ 
10      end if
11    end loop

```

The convergence of the optimization problem is guaranteed *only if* all participants at the auction are either providers of resources or consumers of resources, but not both providers and consumers at the same time. Nevertheless, the clock algorithm only finds a feasible solution; it does not guarantee its optimality.

The authors of [456] have implemented the algorithm and allowed internal use of it within Google; their preliminary experiments show that the system led to substantial improvements. One of the most interesting side effects of the new resource allocation policy is that the users were encouraged to make their applications more flexible and mobile to take advantage of the flexibility of the system controlled by the ASCA algorithm.

Auctioning algorithms are very appealing because they support resource bundling and do not require a model of the system. At the same time, a practical implementation of such algorithms is challenging. First, requests for service arrive at random times, while in an auction, all participants must react to a bid at the same time. Periodic auctions must then be organized but this adds to the delay of the response time. Second, there is an incompatibility between cloud elasticity, which guarantees that the demand for resources of an existing application will be satisfied immediately, and the idea of periodic auctions.

9.11 Cloud resource utilization and energy efficiency

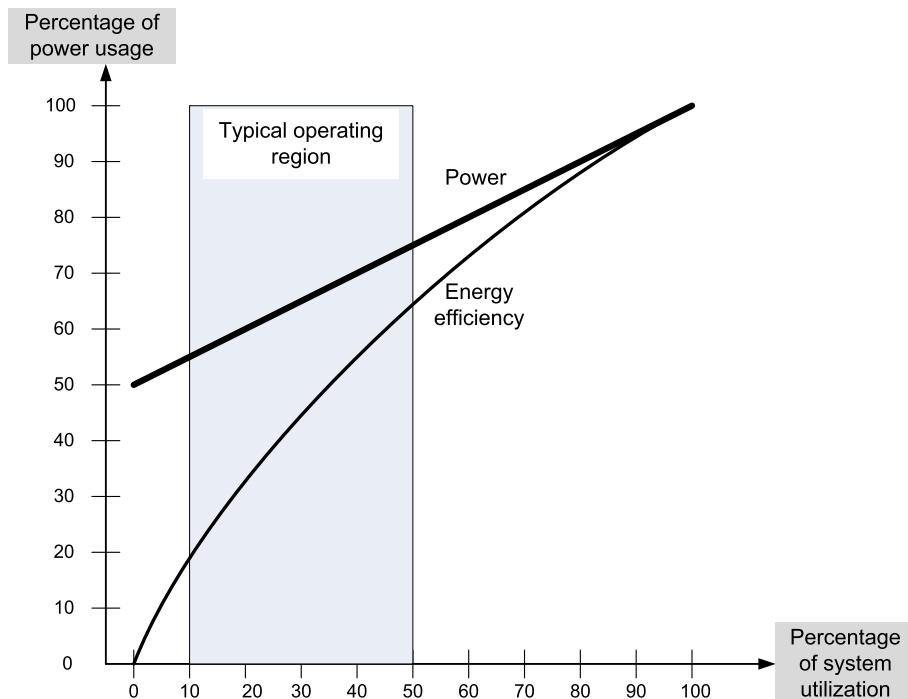
According to Moore's Law, the number of transistors on a chip, thus, the computing power of microprocessors doubles approximately every 1.5 years, and electrical efficiency of computing devices doubles also about every 1.5 years [280]. Thus, performance growth rate and improvements in electrical efficiency almost cancel out. It follows that the energy use for computing scales linearly with the number of computing devices. The number of computing devices continues to grow and many are now housed in large cloud data centers.

The energy consumption of cloud data centers is growing and has a significant ecological impact. It also affects the cost of cloud services. The energy costs are passed on to the users of cloud services and differ from one country to another and from one region to another. For example, the costs for two AWS regions, US East and South America are: upfront for a year \$2 604 versus \$5 632 and hourly \$0.412 versus \$0.724, respectively. Higher energy and communication costs are responsible for the significant difference in this example; the energy costs for the two regions differ by about 40%.

All these facts justify the need to take a closer look at cloud energy consumption, a complex subject extensively discussed in the literature [9,30,45,51,52,331,494,499]. The topics to be covered are: how to define the energy efficiency, how energy-efficient are the processors, the storage devices, the networks, and the other physical elements of the cloud infrastructure, and what are the constraints and how well are these resources managed?

Cloud elasticity and over provisioning. One of the main appeals of utility computing is elasticity. *Elasticity* means that additional resources are guaranteed to be allocated when an application needs them and these resources will be released when they are no longer needed. The user ends up paying only for the resources it has actually used.

Over-provisioning means that a cloud service provider has to invest in a larger infrastructure than the *typical* cloud workload warrants. It follows that the average cloud server utilization is low [9,64,331]; see Fig. 9.10. The low server utilization affects negatively the common measure of energy efficiency,

**FIGURE 9.10**

Even when power requirements scale linearly with the load, the energy efficiency of a computing system is not a linear function of the load; even when idle, a system may use 50% of the power corresponding to the full load. Data collected over a long period of time shows that the typical operating region for the servers at a data center is from about 10% to 50% of the load [51].

the performance per watt of power and the ecological impact of cloud computing. Overprovisioning is not economically sustainable [95].

Elasticity is based on overprovisioning and on the assumption that there is an effective admission control mechanism. Another assumption is that the likelihood of all running applications dramatically increasing their resource consumption at the same time is extremely low. This assumption is realistic, though we have seen cases when a system is overloaded due to concurrent access by large crowds, e.g., the phone system in case of a catastrophic event such as an earthquake. A possible solution is to ask cloud users to specify in their service request the type of workloads and to pay for access accordingly, e.g., a low rate for slow varying and a high rate for workloads with sudden peaks.

Energy efficiency and energy-proportional systems. An energy-proportional system consumes no power when idle, very little power under a light load and, gradually, more power as the load increases. By definition, an ideal energy-proportional system is always operating at 100% efficiency. Humans are a good approximation of an ideal energy proportional system; the human energy consumption is

about 70 W at rest, 120 W on average on a daily basis, and can go as high as 1 000–2 000 W during a strenuous, short-time effort [51].

In real life, even systems whose power requirements scale linearly, when idle, use more than half the power consumed at full load [9]. Indeed, a 2.5 GHz Intel E5200 dual-core desktop processor with 2 GB of RAM consumes 70 W when idle and 110 W when fully loaded; a 2.4 GHz Intel Q6600 processor with 4 GB of RAM consumes 110 W when idle and 175 W when fully loaded [45].

Different subsystems of a computing system behave differently in terms of energy efficiency; while many processors have relatively good energy-proportional profiles, significant improvements in memory and disk subsystems are necessary. The processors used in servers consume less than one-third of their peak power at very low load and have a dynamic range of more than 70% of peak power; the processors used in mobile and/or embedded applications are better in this respect.

The dynamic power range⁸ of other system components is narrower, less than 50% for DRAM, 25% for disk drives, and 15% for networking switches [51]. Power consumption of such devices is: 4.9 KW for a 604.8 TB, HP 8100 EVA storage server, 3.8 KW for the 320 Gbps Cisco 6509 switch, 5.1 KW for the 660 Gbps Juniper MX-960 gateway router [45].

The alternative to the wasteful resource management policy when the servers are *always on*, regardless of their load, is to develop *energy-aware load balancing and scaling* policies. Such policies combine *dynamic power management* with load balancing and attempt to identify servers operating outside their optimal energy regime and decide if and when they should be switched to a sleep state or what other actions should be taken to optimize the energy consumption.

Energy saving. The effort to reduce the energy use is focused on the computing, networking, and storage activities of a data center. A 2010 report shows that a typical Google cluster spends most of its time within the 10–50% CPU utilization range; there is a mismatch between server workload profile and server energy efficiency [9]. A similar behavior is also seen in the data center networks; these networks operate in a very narrow dynamic range; the power consumed when the network is idle is significant compared to the power consumed when the network is fully utilized.

A strategy to reduce energy consumption is to concentrate the workload on a small number of disks and allow the others to operate in a low-power mode. One of the techniques to accomplish this is based on replication. A replication strategy based on a sliding window is reported in [499]; measurement results indicate that it performs better than LRU, MRU, and LFU⁹ policies for a range of file sizes, file availability, and number of client nodes and the power requirements are reduced by as much as 31%.

Another technique is based on data migration. The system in [229] uses data storage in virtual nodes managed with a distributed hash table; the migration is controlled by two algorithms: a short-term optimization algorithm used for gathering or spreading virtual nodes according to the daily variation of the workload so that the number of active physical nodes is reduced to a minimum, and a long-term optimization algorithm, used for coping with changes in the popularity of data over a longer period, e.g., a week.

A number of proposals have emerged for *energy proportional* networks [10]; the energy consumed by such networks is proportional with the communication load. For example, a data center interconnec-

⁸ The dynamic range in this context is determined by the lower and the upper limit of the device power consumption. A large dynamic range means that the device is better, it is able to operate at a lower fraction of its peak power when its load is low.

⁹ LRU (Least Recently Used), MRU (Most Recently Used), and LFU (Least Frequently Used) are replacement policies used by memory hierarchies for caching and paging.

tion network based on a flattened butterfly topology is more energy and cost efficient according to [9]. High-speed channels typically consist of multiple serial lanes with the same data rate; a physical unit is stripped across all the active lanes. Channels commonly operate plesiochronously¹⁰ and are always on, regardless of the load, because they must still send idle packets to maintain byte and lane alignment across the multiple lanes. An example of an energy proportional network is *InfiniBand*, discussed in Section 6.7.

Many proposals argue that dynamic resource provisioning is necessary to minimize power consumption. Two main issues are critical for energy saving: the amount of resources allocated to each application and the placement of individual workloads. For example, a resource management framework combining a utility-based dynamic Virtual Machine provisioning manager with a dynamic VM placement manager to minimize power consumption and reduce Service Level Agreement violations is presented in [490].

Energy optimization is an important policy for cloud resource management but it cannot be considered in isolation, it has to be coupled with admission control, capacity allocation, load balancing, and quality of service. Existing mechanisms cannot support concurrent optimization of all the policies. Mechanisms based on a solid foundation such as control theory are too complex and do not scale well, those based on machine learning are not fully developed, and the others require a model of a system with a dynamic configuration operating in a fast-changing environment.

9.12 Resource management and dynamic application scaling

The demand for resources can be a function of the time of the day, can monotonically increase or decrease in time, or can experience predictable or unpredictable peaks. For example, a new web service will experience a low request rate at the very beginning, and the load will exponentially increase if the service is successful. A service for income tax processing will experience a peak around the tax filling deadline, while access to a service provided by FEMA (Federal Emergency Management Agency) will increase dramatically after a natural disaster.

The elasticity of a public cloud, the fact that it can supply to an application precisely the amount of resources it needs and that one pays only for the resources it consumes are serious incentives to migrate to a public cloud. The question we address is how scaling can be actually implemented in a cloud when a very large number of applications exhibit this often unpredictable behavior [83,335,489]. To make matters worse, in addition to an unpredictable external load, the cloud resource management has to deal with resource reallocation due to server failures.

We distinguish two scaling strategies, vertical and horizontal. *Vertical scaling* keeps the number of VMs of an application constant, but increases the amount of resources allocated to each one of them. This can be done either by migrating the VMs to more powerful servers, or by keeping the VMs on the same servers, but increasing their share of the CPU time. The first alternative involves additional overhead; the VM is stopped, a snapshot of it is taken, the file is transported to a more powerful server, and, finally, the VM is restated at the new site.

¹⁰ Different parts of the system are almost, but not quite perfectly, synchronized; in this case, the core logic in the router operates at a frequency different from that of the I/O channels.

Horizontal scaling is the most common scaling strategy on a cloud; it is done by increasing the number of VMs as the load increases and reducing this number when the load decreases. Often, this leads to an increase of communication bandwidth consumed by the application. Load balancing among the running VMs is critical for this mode of operation. For a very large application, multiple load balancers may need to cooperate with one another. In some instances the load balancing is done by a front-end server, which distributes incoming requests of a transaction-oriented system to back-end servers.

An application should be designed to support scaling. Workload partitioning of a *modularly divisible* application is static as we have seen in Section 11.5. Static workload partitioning is decided a priori and cannot be changed thus, the only alternative is vertical scaling. The workload of an *arbitrarily divisible* application can be partitioned dynamically. As the load increases, the system can allocate additional VMs to process the additional workload. Most cloud applications belong to this class, and this justifies our statement that horizontal scaling is the most common scaling strategy.

Mapping a computation means to assign suitable physical servers to the application. A very important first step in application processing is to identify the type of application and map it accordingly. For example, a communication-intensive application should be mapped to a powerful server to minimize the network traffic. Such a mapping may increase the cost per unit of CPU usage, but it will decrease the computing time and, probably, reduce the overall cost for the user. At the same time, it will reduce network traffic, a highly desirable effect from the perspective of an *arbitrarily divisible* application.

To scale up or down a compute-intensive application a good strategy is to increase/decrease the number of VMs or instances. As the load is relatively stable, the overhead of starting up or terminating an instance does not increase significantly the computing time or the cost.

Several strategies to support scaling exist. *Automatic VM scaling* uses predefined metrics, e.g., CPU utilization to make scaling decisions. Automatic scaling requires *sensors* to monitor the state of the VMs and servers and *controllers* to make decisions based on the information about the state of the cloud. Controllers often use a state machine model for decision making. Amazon and Rightscale offer automatic scaling. AWS *CloudWatch* service supports applications monitoring and allows a user to set up conditions for automatic migrations. AWS *Elastic Load Balancing* service automatically distributes incoming application traffic across multiple EC2; *Elastic Beanstalk* allows dynamic scaling between a low and a high number of instances specified by the user; see Section 2.2. A cloud user usually has to pay for the more sophisticated scaling services such as *Elastic Beanstalk*.

Nonscalable or single load balancers are also used for horizontal scaling. Google uses machine learning to optimize resource utilization and minimize the risk of killing a task or performance degradation due to CPU throttling. The system called *Autoscaling* [426] configures automatically resources allocated to a job using horizontal and vertical scaling in an attempt to reduce the *slack*, i.e., the difference between user specified limits for CPU cores and RAM and the actual resource usage.

9.13 Control theory and optimal resource management (R)

Control theory has been used to design adaptive resource management for many classes of applications including power management [268], task scheduling [315], QoS adaptation in web servers [4], and load balancing [355,395]. The classical feedback control methods are used in all these cases to regulate the key operating parameters of the system based on measurement of the system output. The feedback con-

trol for these methods assumes a linear time-invariant system model, and a closed-loop controller. This controller is based on an open-loop system transfer function, which satisfies stability and sensitivity constraints.

A technique to design self-managing systems based on concepts from control theory is discussed in [504]. This technique allows multiple QoS objectives and operating constraints to be expressed as a cost function. The technique can be applied to stand-alone or distributed web servers, database servers, high performance application servers, and embedded systems.

The following discussion considers a single processor serving a stream of input requests with the goal of minimizing a cost function reflecting the response time and the power consumption. The goal is to illustrate the methodology for optimal resource management based on control theory concepts. The analysis is intricate and cannot be easily extended to a collection of servers.

Control theory principles. An overview of control theory principles used for optimal resource allocation is presented next. Optimal control generates a sequence of control inputs over a look-ahead horizon, while estimating changes in operating conditions. A convex cost function has as arguments $x(k)$, the state at step k , and $u(k)$, the control vector; this cost function is minimized subject to the constraints imposed by the system dynamics. The discrete-time optimal control problem is to determine the sequence of control variables $u(i), u(i+1), \dots, u(n-1)$ minimizing the expression

$$J(i) = \Phi(n, x(n)) + \sum_{k=i}^{n-1} L^k(x(k), u(k)) \quad (9.60)$$

where $\Phi(n, x(n))$ is the cost function of the final step, n , and $L^k(x(k), u(k))$ is a time-varying cost function at the intermediate step k over the horizon $[i, n]$. The minimization is subject to the constraints $x(k+1) = f^k(x(k), u(k))$, where $x(k+1)$, the system state at time $k+1$, is a function of $x(k)$, the state at time k , and of $u(k)$, the input at time k ; in general, the function f^k is time-varying, thus its superscript.

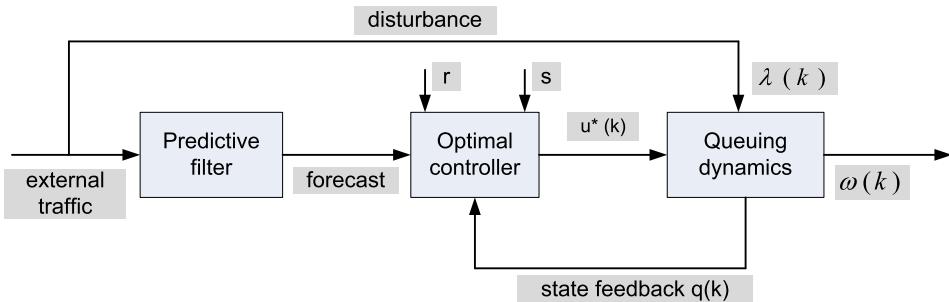
One of the techniques to solve this problem is based on the *Lagrange multiplier* method of finding the extremes (minima or maxima) of a function subject to constraints; more precisely, if we wish to maximize the function $g(x, y)$ subject to the constraint $h(x, y) = k$ we introduce a Lagrange multiplier λ . Then we study the function

$$\Lambda(x, y, \lambda) = g(x, y) + \lambda \times [h(x, y) - k]. \quad (9.61)$$

A necessary condition for the optimality is that (x, y, λ) is a stationary point for $\Lambda(x, y, \lambda)$, in other words

$$\nabla_{x,y,\lambda} \Lambda(x, y, \lambda) = 0 \quad \text{or} \quad \left(\frac{\partial \Lambda(x, y, \lambda)}{\partial x}, \frac{\partial \Lambda(x, y, \lambda)}{\partial y}, \frac{\partial \Lambda(x, y, \lambda)}{\partial \lambda} \right) = 0. \quad (9.62)$$

The Lagrange multiplier at time step k is $\lambda(k)$ and we solve Eq. (9.62) as an unconstrained optimization problem. We define an adjoint cost function which includes the original state constraints as the Hamiltonian function H , then we construct the adjoint system consisting of the original state equation

**FIGURE 9.11**

The structure of an optimal controller described in [504]; the controller uses the feedback regarding the current state as well as the estimation of the future disturbance due to environment to compute the optimal inputs over a finite horizon. The two parameters r and s are the weighting factors of the performance index.

and the *costate equation*¹¹ governing the Lagrange multiplier. Thus, we define a two-point boundary problem¹²; the state x_k develops forward in time, while the costate occurs backward in time.

A model capturing both QoS and energy consumption for a single server system. We now turn our attention to the case of a single processor serving a stream of input requests. To compute the optimal inputs over a finite horizon, the controller in Fig. 9.11 uses the feedback regarding the current state and the estimation of the future disturbance due to the environment. The control task is solved as a state regulation problem updating the initial and final states of the control horizon.

We use a simple queuing model to estimate the response time; requests for service at processor P are processed on an FCFS basis. We do not assume a priori distributions of the arrival process and of the service process; instead, we use the estimate $\hat{\Lambda}(k)$ of the arrival rate $\Lambda(k)$ at time k . We also assume that the processor can operate at frequencies $u(k)$ in the range $u(k) \in [u_{min}, u_{max}]$ and call $\hat{c}(k)$ the time to process a request at time k when the processor operates at the highest frequency in the range, u_{max} . Then, we define the scaling factor $\alpha(k) = u(k)/u_{max}$ and we express an estimate of the processing rate $N(k)$ as $\alpha(k)/\hat{c}(k)$.

The behavior of a single processor is modeled as a non-linear, time-varying, discrete-time state equation. If T_s is the sampling period, defined as the time difference between two consecutive observations of the system, e.g., the one at time $(k + 1)$ and the one at time k , then the size of the queue at time $(k + 1)$ is

$$q(k + 1) = \max \left\{ \left[q(k) + \left(\hat{\Lambda}(k) - \frac{u(k)}{\hat{c}(k) \times u_{max}} \right) \times T_s \right], 0 \right\} \quad (9.63)$$

¹¹ The costate equation is related to the state equation in optimal control.

¹² A boundary value problem has conditions specified at the extremes of the independent variable while an initial value problem has all of the conditions specified at the same value of the independent variable in the equation. The common case is when boundary conditions are supposed to be satisfied at two points—usually the starting and ending values of the integration.

The first term, $q(k)$, is the size of the input queue at time k , and the second one is the difference between the number of requests arriving during the sampling period, T_s , and those processed during the same interval. The response time $\omega(k)$ is the sum of the waiting time and the processing time of the requests $\omega(k) = (1 + q(k)) \times \hat{c}(k)$. Indeed, the total number of requests in the system is $(1 + q(k))$, and the departure rate is $1/\hat{c}(k)$.

We wish to capture both the QoS and the energy consumption, as both affect the cost of providing the service. A utility function, such as the one depicted in Fig. 9.14, captures the rewards, as well as the penalties specified by the SLA for the response time. In the queuing model, the utility is a function of the size of the queue and can be expressed as a quadratic function of the response time

$$S(q(k)) = 1/2 \left(s \times (\omega(k) - \omega_0)^2 \right) \quad (9.64)$$

with ω_0 , the response time set point and $q(0) = q_0$, the initial value of the queue length. The energy consumption is a quadratic function of the frequency

$$R(u(k)) = 1/2 \left(r \times u(k)^2 \right). \quad (9.65)$$

The two parameters s and r are weights for the two components of the cost, derived from the utility function and from the energy consumption, respectively. We have to pay a penalty for the requests left in the queue at the end of the control horizon, a quadratic function of queue length

$$\Phi(q(N)) = 1/2 \left(v \times q(N)^2 \right). \quad (9.66)$$

The performance measure of interest is a cost expressed as

$$J = \Phi(q(N)) + \sum_{k=1}^{N-1} [S(q(k)) + R(q(k))]. \quad (9.67)$$

The problem is to find the optimal control u^* and the finite time horizon $[0, N]$ such that the trajectory of the system subject to optimal control is q^* , and the cost J in Eq. (9.67) is minimized subject to the following constraints

$$q(k+1) = \left[q(k) + \left(\hat{\Lambda}(k) - \frac{u(k)}{\hat{c}(k) \times u_{max}} \right) \times T_s \right], \quad q(k) \geq 0, \quad \text{and} \quad u_{min} \leq u(k) \leq u_{max}. \quad (9.68)$$

When the state trajectory $q(\cdot)$ corresponding to the control $u(\cdot)$ satisfies the constraints

$$\Gamma 1 : q(k) > 0, \quad \Gamma 2 : u(k) \geq u_{min}, \quad \Gamma 3 : u(k) \leq u_{max}, \quad (9.69)$$

then the pair $[q(\cdot), u(\cdot)]$ is called a *feasible state*. If the pair minimizes the Eq. (9.67) then the pair is *optimal*.

The Hamiltonian H in our example is

$$H = S(q(k)) + R(u(k)) + \lambda(k+1) \times \left[q(k) + \left(\Lambda(k) - \frac{u(k)}{c \times u_{max}} \right) T_s \right] \\ + \mu_1(k) \times (-q(k)) + \mu_2(k) \times (-u(k) + u_{min}) + \mu_3(k) \times (u(k) - u_{max}). \quad (9.70)$$

According to Pontryagin's minimum principle,¹³ the necessary condition for a sequence of feasible pairs to be optimal pairs is the existence of a sequence of costates λ and a Lagrange multiplier $\mu = [\mu_1(k), \mu_2(k), \mu_3(k)]$ such that

$$H(k, q^*, u^*, \lambda^*, \mu^*) \leq H(k, q, u^*, \lambda^*, \mu^*), \quad \forall q \geq 0 \quad (9.71)$$

where the Lagrange multipliers, $\mu_1(k), \mu_2(k), \mu_3(k)$, reflect the sensitivity of the cost function to the queue length at time k and the boundary constraints and satisfy several conditions

$$\mu_1(k) \geq 0, \quad \mu_1(k)(-q(k)) = 0, \quad (9.72)$$

$$\mu_2(k) \geq 0, \quad \mu_2(k)(-u(k) + u_{min}) = 0, \quad (9.73)$$

$$\mu_3(k) \geq 0, \quad \mu_3(k)(u(k) - u_{max}) = 0. \quad (9.74)$$

A detailed analysis of the methods to solve this problem and the analysis of the stability conditions is beyond the scope of our discussion and can be found in [504].

The extension of the techniques for optimal resource management from a single system to a cloud with a very large number of servers is a rather challenging area of research. The problem is even harder when, instead of transaction-based processing, the cloud applications require the implementation of a complex workflow.

9.14 Stability of two-level resource allocation strategy (R)

The discussion in Section 9.13 shows that a server can be assimilated with a closed-loop control system and that we can apply theoretical control principles to resource allocation. We now discuss a two-level resource allocation architecture based on control theory concepts for the entire cloud; see Fig. 9.12. The automatic resource management is based on two levels of controllers, one for the service provider and one for the application.

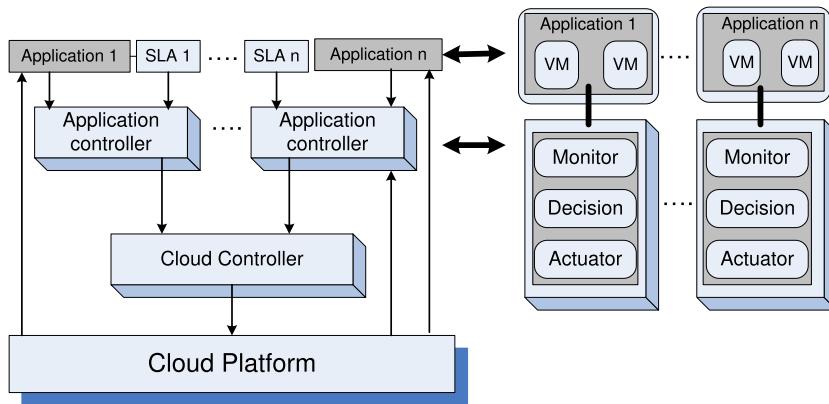
The main components of a control system are: the inputs, the control system components, and the outputs. The inputs in such models are: the offered workload and the policies for admission control, the capacity allocation, the load balancing, the energy optimization, and the QoS guarantees in the cloud. The system components are *sensors* used to estimate relevant measures of performance and *controllers* which implement various policies; the output is the resource allocations to the individual applications.

The controllers use the feedback provided by sensors to stabilize the system; stability is related to the change in the output. If the change is too large, then the system may become unstable. In our context, the system could experience thrashing, the amount of useful time dedicated to the execution of applications becomes increasingly smaller and most of the system resources are occupied by management functions.

There are three main sources of instability in any control system:

1. The delay in getting the system reaction after a control action.

¹³ Pontryagin's principle is used in the optimal control theory to find the best possible control which leads a dynamic system transition from one state to another, subject to a set of constraints.

**FIGURE 9.12**

A two-level control architecture; application controllers and cloud controllers work in concert.

2. The granularity of the control, the fact that a small change enacted by the controllers leads to very large changes in the output.
3. Oscillations, when the changes in the input are too large and the control is too weak, such that the changes in the input propagate directly to the output.

Two types of policies are used in autonomic systems: (i) threshold-based policies and (ii) sequential decision policies based on Markovian decision models. In the first case, upper and lower bounds on performance trigger adaptation through resource reallocation; such policies are simple and intuitive but require setting per-application thresholds.

Lessons learned from the experiments with two levels of controllers and the two types of policies are discussed in [156]. A first observation is that the actions of the control system should be carried out in a rhythm that does not lead to instability; adjustments should only be carried out after the performance of the system has stabilized. The controller should measure the time for an application to stabilize and adapt to the manner in which the controlled system reacts.

If an upper and a lower threshold are set, then instability occurs when they are too close to one another if the variations of the workload are large enough and the time required to adapt does not allow the system to stabilize. The actions consist of allocation/deallocation of one or more VMs; sometimes, allocation/deallocation of a single VM required by one of the thresholds may cause crossing the other threshold, another source of instability.

9.15 Feedback control based on dynamic thresholds (R)

The elements involved in a control system are sensors, monitors, and actuators. The *sensors* measure the parameter(s) of interest, then transmit the measured values to a *monitor* which determines if the system behavior must be changed, and, if so, it requests the *actuators* to carry out the necessary actions. Often,

the parameter used for admission control policy is the current system load; when a threshold, e.g., 80%, is reached, the cloud stops accepting additional load.

In practice, the implementation of such a policy is challenging, or outright infeasible. First, due to the very large number of servers and to the fact that the load changes rapidly in time, the estimation of the current system load is likely to be inaccurate. Second, the ratio of average to maximal resource requirements of individual users specified in a service-level agreement is typically very high. Once an agreement is in place, user demands must be satisfied; user requests for additional resources within the SLA limits cannot be denied.

Thresholds. A threshold is the value of a parameter related to the state of a system that triggers a change in the system behavior. Thresholds are used in control theory to keep critical parameters of a system in a predefined range. The threshold could be *static*, defined once and for all, or could be *dynamic*. A dynamic threshold could be based on an average of measurements carried out over a time interval, a so called *integral control*; the dynamic threshold could also be a function of the values of multiple parameters at a given time, or a mixture of the two.

To maintain the system parameters in a given range, a *high* and a *low* threshold are often defined. The two thresholds determine different actions; for example, a high threshold could force the system to limit its activities and a low threshold could encourage additional activities. *Control granularity* refers to the level of detail of the information used to control the system. *Fine control* means that very detailed information about the parameters controlling the system state is used, while *coarse control* means that the accuracy of these parameters is traded off for the efficiency of implementation.

Proportional thresholding. Application of these ideas to cloud computing, in particular to the IaaS delivery model, and a strategy for resource management called *proportional thresholding* are discussed in [306]. The questions addressed are:

- Is it beneficial to have two types of controllers: (1) *application controllers*, which determine if additional resources are needed, and (2) *cloud controllers*, which arbitrate requests for resources and allocate the physical resources?
- Is it feasible to consider *fine control*? Is *coarse control* more adequate in a cloud computing environment?
- Are dynamic thresholds based on time averages better than static ones?
- Is it better to have a high and a low threshold, or is it sufficient to define only a high threshold?

The first two questions are interrelated. It seems more appropriate to have two controllers, one with knowledge of the application and one aware of the state of the cloud. In this case, a coarse control is more adequate for many reasons. As mentioned earlier, the cloud controller can only have a very rough approximation of the cloud state. Moreover, the service provider may wish to hide some of the information they have to simplify its resource management policies. For example, CSPs may not allow a VM to access information available to hypervisor-level sensors and actuators.

To answer the last two questions one has to define a measure of “goodness.” In the experiments reported in [306], the parameter measured is the average CPU utilization, and a strategy is better than another if it reduces the number of requests made by the application controllers to add or remove VMs to the pool of those available to the application.

A control theoretical approach to address these questions is challenging. A pragmatic approach, qualitative arguments, and simulation results using a synthetic workload for a transaction-oriented application on a web server are reported in [306].

The essence of the proportional thresholding is captured by the following algorithm:

1. Compute the integral value of the high and the low threshold as averages of the maximum and, respectively, the minimum of the processor utilization over the process history.
2. Request additional VMs when the average value of the CPU utilization over the current time slice exceeds the high threshold.
3. Release a VM when the average value of the CPU utilization over the current time slice falls below the low threshold.

The conclusions reached based on experiments with 3 VMs are: (a) dynamic thresholds perform better than the static ones and (b) two thresholds are better than one. While confirming our intuition, such results have to be justified by experiments in a realistic environment. Moreover, convincing results cannot be based on empirical values for some of the parameters required by integral control equations.

9.16 Coordination of autonomic performance managers (R)

Can specialized autonomic performance managers cooperate to optimize power consumption and, at the same time, satisfy the requirements of SLAs? This is the question examined in a paper reporting on experiments carried out on a set of blades mounted on a chassis [268]. The experimental setup is shown in Fig. 9.13. Extending the techniques discussed in this report to a large-scale farm of servers poses significant problems; computational complexity is just one of them.

Virtually all modern processors support Dynamic Voltage Scaling as a mechanism for energy saving; indeed, the energy dissipation scales quadratically with the supply voltage. The power management controls the CPU frequency, thus, the rate of instruction execution. For some compute-intensive workloads the performance decreases linearly with the CPU clock frequency, while for others the effect of lower clock frequency is less noticeable or non-existent. The clock frequency of individual blades/servers is controlled by a power manager typically implemented in the firmware; it adjusts the clock frequency several times a second.

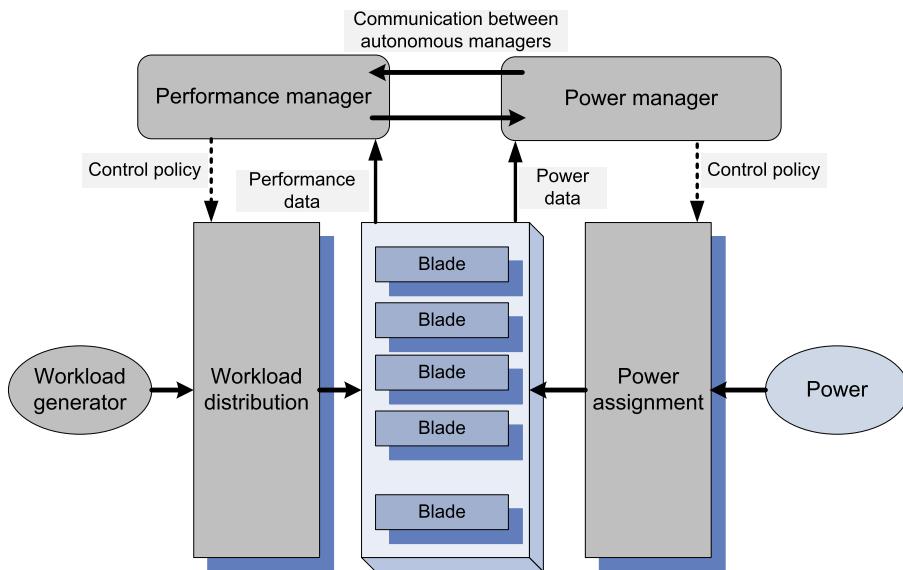
The approach to coordinating power and performance management in [268] is based on several ideas:

- Use a joint utility function for power and performance. The joint performance-power utility function, $U_{pp}(R, P)$, is a function of the response time, R , and the power, P , and it can be of the form

$$U_{pp}(R, P) = U(R) - \epsilon \times P \quad \text{or} \quad U_{pp}(R, P) = \frac{U(R)}{P}, \quad (9.75)$$

with $U(R)$ being the utility function based on response time only and ϵ being a parameter to weigh the influence of the two factors, response time and power.

- Identify a minimal set of parameters to be exchanged between the two managers.
- Set up a power cap for individual systems based on the utility-optimized power management policy.
- Use a standard performance manager modified only to accept input from the power manager regarding the frequency determined according to the power management policy. The power manager consists of TCL and C programs to compute the per-server (per-blade) power caps and send them

**FIGURE 9.13**

Autonomous performance and power managers cooperate to ensure SLA prescribed performance and energy optimization; they are fed with performance and power data and implement the performance and power management policies, respectively.

via IPMI¹⁴ to the firmware controlling the blade power. The power manager and the performance manager interact, but no negotiation between the two agents is involved.

- Use standard software systems. For example, use the WXD (WebSphere Extended Deployment), a middleware that supports setting performance targets for individual web applications and for the monitor response time, and periodically recompute the resource allocation parameters to meet the targets set. Use the Wide-Spectrum Stress Tool from IBM Web Services Toolkit as a workload generator.

For practical reasons, the utility function is expressed in terms of n_c , the number of clients, and p_κ , the power cap, as in

$$U'(p_\kappa, n_c) = U_{pp}(R(p_\kappa, n_c), P(p_\kappa, n_c)). \quad (9.76)$$

The optimal power cap, p_κ^{opt} is a function of the workload intensity expressed by the number of clients, n_c ,

$$p_\kappa^{opt}(n_c) = \arg \max U'(p_\kappa, n_c). \quad (9.77)$$

¹⁴ Intelligent Platform Management Interface (IPMI) is a standardized computer system interface developed by Intel and used by system administrators to manage a computer system and monitor its operation.

The hardware used for these experiments were the Goldensbridge blade with an Intel Xeon processor running at 3 GHz with 1 GB of level 2 cache and 2 GB of DRAM and with hyper threading enabled. A blade could serve 30 to 40 clients with a response time at or better than the 1000 msec limit. When p_k is lower than 80 Watts, the processor runs at its lowest frequency, 375 MHz, while for p_k at or larger than 110 Watts, the processor runs at its highest frequency, 3 GHz.

Three types of experiments were conducted: (i) with the power management turned off; (ii) when the dependence of the power consumption and the response time were determined through a set of exhaustive experiments; and (iii) when the dependency of the power cap p_k on n_c is derived via reinforcement-learning models.

The second type of experiments lead to the conclusion that both the response time and the power consumed are nonlinear functions of the power cap, p_k , and the number of clients, n_c . More specifically, the conclusions of these experiments are: (i) at a low load, the response time is well below the target of 1000 msec; (ii) at medium and high loads, the response time decreases rapidly when p_k increases from 80 to 110 watts; and (iii) the consumed power increases rapidly as the load increase for a given value of the power cap.

The machine learning algorithm used for the third type of experiments was based on the Hybrid Reinforcement Learning algorithm described in [474]. In the experiments using the machine learning model, the power cap required to achieve a response time lower than 1000 msec for a given number of clients was the lowest when $\epsilon = 0.05$ and the first utility function given by Eq. (9.75) was used; for example, when $n_c = 50$ then $p_k = 109$ Watts when $\epsilon = 0.05$, while $p_k = 120$ when $\epsilon = 0.01$.

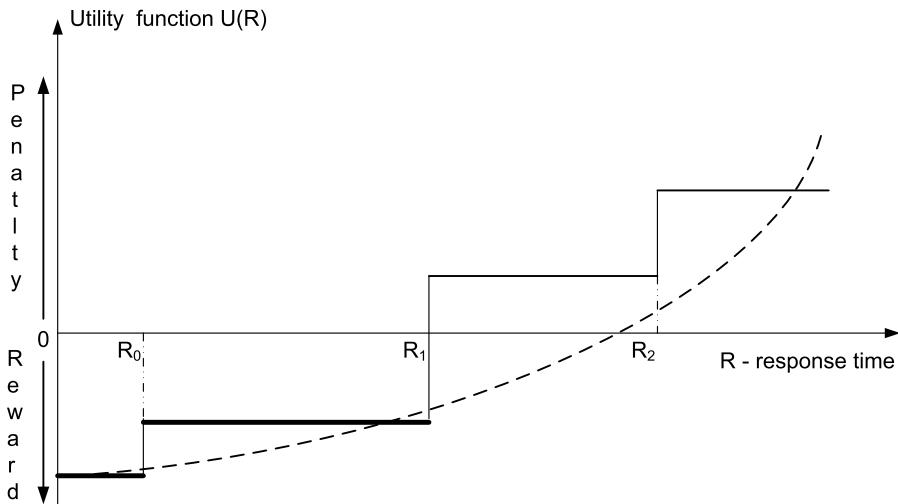
9.17 A utility model for cloud-based web services (R)

A *utility function* relates the “benefits” of a service to the “cost” to provide the service. For example, the benefit could be revenue and the cost could be the power consumption. In this section, we discuss a utility-based approach for autonomic management. The goal is to maximize the total profit computed as the difference between the revenue guaranteed by an SLA and the total cost to provide the services.

An SLA often specifies the rewards as well as penalties associated with specific performance metrics. Sometimes, the quality of service translates into average response time; this is the case of cloud-based web services when the SLA often specifies explicitly this requirement. For example, Fig. 9.14 shows the case when the performance metric is R , the response time. The largest reward can be obtained when $R \leq R_0$; a slightly lower reward corresponds to $R_0 < R \leq R_1$; when $R_1 < R \leq R_2$, instead of gaining a reward, the provider of service pays a small penalty; the penalty increases when $R > R_2$. A utility function, $U(R)$, which captures this behavior is a sequence of step functions; the utility function is sometimes approximated by a quadratic curve as discussed in Section 9.13.

Formulated as an optimization problem, the solution discussed in [11] addresses multiple policies, including QoS. The cloud model for this optimization is quite complex and requires a fair number of parameters. We assume a cloud providing $|K|$ different classes of service, each class k involving N_k applications. For each class $k \in K$ call v_k the revenue (or the penalty) associated with a response time r_k and assume a linear dependency for this utility function of the form $v_k = v_k^{\max} (1 - r_k/r_k^{\max})$; see Fig. 9.15(a); call $m_k = -v_k^{\max}/r_k^{\max}$ the slope of the utility function.

The system is modeled as a network of queues with multiqueues for each server and with a delay center, which models the think time of the user after the completion of service at one server and the start

**FIGURE 9.14**

A utility function $U(R)$ is a series of step functions with jumps corresponding to the response time, $R = R_0|R_1|R_2$, when the reward and the penalty levels change according to the SLA. Dotted lines show a quadratic approximation of the utility function.

of processing at the next server; see Fig. 9.15(b). Upon completion, a class k request either completes with probability $1 - \sum_{k' \in K} \pi_{k,k'}$ or returns to the system as a class k' request with transition probability $\pi_{k,k'}$. Call λ_k the external arrival rate of class k requests and Λ_k the aggregate rate for class k , $\Lambda_k = \lambda_k + \sum_{k' \in K} \Lambda_{k'} \pi_{k,k'}$.

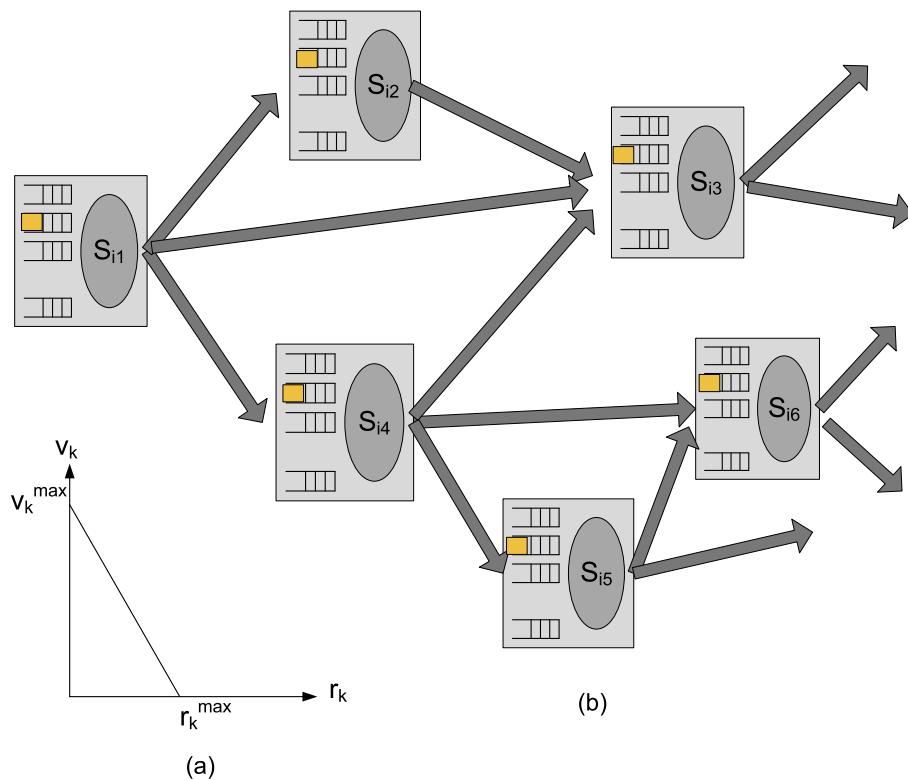
Typically, CPU and memory are considered as representative for resource allocation. For simplicity, we assume a single CPU that runs at a discrete set of clock frequencies and a discrete set of supply voltages according to a DVFS model; the power consumption on a server is a function of the clock frequency. The scheduling of a server is work-conserving¹⁵ and is modeled as a Generalized Processor Sharing (GPS) scheduling [536]. Analytical models [6], [384] are too complex for large systems.

The optimization problem formulated in [11] involves five terms: A and B reflect revenues, C is the cost for servers in a low power, stand-by mode, D is the cost of active servers given their operating frequency, E is the cost for switching servers from low-power, stand-by mode, to active state, and F is the cost for migrating VMs from one server to another. There are 9 constraints $\Gamma_1, \Gamma_2, \dots, \Gamma_9$ for this mixed integer non-linear programming problem. The decision variables for this optimization problem are listed in Table 9.7 and the parameters used are shown in Table 9.8.

The expression to be maximized is:

$$(A + B) - (C + D + E + F) \quad (9.78)$$

¹⁵ A scheduling policy is work conserving if the server cannot be idle while there is work to be done.

**FIGURE 9.15**

(a) The utility function, v_k the revenue (or the penalty) associated with a response time r_k for a request of class $k \in K$; the slope of the utility function is $m_k = -v_k^{\max}/r_k^{\max}$. (b) At each server S_i , there are $|K|$ queues for each one of $k \in K$ classes of requests. A tier consists of all requests of class $k \in K$ at all servers S_{ij} , $i \in I$, $1 \leq j \leq 6$.

Table 9.7 Decision variables for the optimization problem.

Name	Description
x_i	$x_i = 1$ if server $i \in I$ is running, $x_i = 0$ otherwise
$y_{i,h}$	$y_{i,h} = 1$ if server i is running at frequency h , $y_{i,h} = 0$ otherwise
$z_{i,k,j}$	$z_{i,k,j} = 1$ if application tier j of a class k request runs on server i , $z_{i,k,j} = 0$ otherwise
$w_{i,k}$	$w_{i,k} = 1$ if at least one class k request is assigned to server i , $w_{i,k} = 0$ otherwise
$\lambda_{i,k,j}$	rate of execution of applications tier j of class k requests on server i
$\phi_{i,k,j}$	fraction of capacity of server i assigned to tier j of class k requests

with

$$A = \max \sum_{k \in K} \left(-m_k \sum_{i \in I, j \in N_k} \frac{\lambda_{i,k,j}}{\sum_{h \in H_i} (C_{i,h} \times y_{i,h}) \mu_{k,j} \times \phi_{i,k,j} - \lambda_{i,k,j}} \right), \quad B = \sum_{k \in K} u_k \times \Lambda_k, \quad (9.79)$$

Table 9.8 The parameters used for the A , B , C , D , E , and F terms and the constraints Γ_i of the optimization problem.

Name	Description
I	the set of servers
K	the set of classes
Λ_k	the aggregate rate for class $k \in K$, $\Lambda_k = \lambda_k + \sum_{k' \in K} \Lambda_{k'} \pi_{k,k'}$
a_i	the availability of server $i \in I$
A_k	minimum level of availability for request class $k \in K$ specified by the SLA
m_k	the slope of the utility function for a class $k \in K$ application
N_k	number of applications in class $k \in K$
H_i	the range of frequencies of server $i \in I$
$C_{i,h}$	capacity of server $i \in I$ running at frequency $h \in H_i$
$c_{i,h}$	cost for server $i \in I$ running at frequency $h \in H_i$
\bar{c}_i	average cost of running server i
$\mu_{k,j}$	maximum service rate for a unit capacity server for tier j of a class k request
cm	the cost of moving a VM from one server to another
cs_i	the cost for switching server i from the stand-by mode to an active state
$RAM_{k,j}$	the amount of main memory for tier j of class k request
\overline{RAM}_i	the amount of memory available on server i

$$C = \sum_{i \in I} \bar{c}_i, \quad D = \sum_{i \in I, h \in H_i} c_{i,h} \times y_{i,h}, \quad E = \sum_{i \in I} cs_i \max(0, x_i - \bar{x}_i), \quad (9.80)$$

and

$$F = \sum_{i \in I, k \in K, j \in N_j} cm \max(0, z_{i,j,k} - \bar{z}_{i,j,k}). \quad (9.81)$$

The nine constraints are:

- (Γ_1) $\sum_{i \in I} \lambda_{i,k,j} = \Lambda_k, \forall k \in K, j \in N_k, \Rightarrow$ the traffic assigned to all servers for class k requests equals the predicted load for the class.
- (Γ_2) $\sum_{k \in K, j \in N_k} \phi_{i,k,j} \leq 1 \forall i \in I, \Rightarrow$ server i cannot be allocated an workload more than its capacity.
- (Γ_3) $\sum_{h \in H_i} y_{i,h} = x_i, \forall i \in I, \Rightarrow$ if server $i \in I$ is active; it runs at one frequency in the set H_i , only one $y_{i,h}$ is nonzero.
- (Γ_4) $z_{i,k,j} \leq x_i, \forall i \in I, k \in K, j \in N_k \Rightarrow$ requests can only be assigned to active servers.
- (Γ_5) $\lambda_{i,k,j} \leq \Lambda_k \times z_{i,k,j}, \forall i \in I, k \in K, j \in N_k \Rightarrow$ requests may run on server $i \in I$ only if the corresponding application tier has been assigned to server i .
- (Γ_6) $\lambda_{i,k,j} \leq \left(\sum_{h \in H_i} C_{i,h} \times y_{i,h} \right) \mu_{k,j} \times \phi_{i,k,j}, \forall i \in I, k \in K, j \in N_k \Rightarrow$ resources cannot be saturated.
- (Γ_7) $RAM_{k,j} \times z_{i,k,j} \leq \overline{RAM}_i \forall i \in I, k \in K \Rightarrow$ the memory on server i is sufficient to support all applications running on it.

(Γ_8) $\prod_{j=1}^{N_k} (1 - \prod_{i=1}^M (1 - a_i^{w_{i,k}})) \geq A_k, \forall k \in K \Rightarrow$ the availability of all servers assigned to class k request should be at least equal to the minimum required by the SLA.

(Γ_9) $\sum_{j=1}^{N_k} z_{i,k,j} \geq N_k \times w_{i,k}, \forall i \in I, k \in K$
 $\lambda_{i,j,k}, \phi_{i,j,k} \geq 0, \forall i \in I, k \in K, j \in N_k$
 $x_i, y_{i,h}, z_{i,k,j}, w_{i,k} \in \{0, 1\}, \forall i \in I, k \in K, j \in N_k \Rightarrow$ constraints and relations among decision variables.

Clearly, this approach is not scalable to clouds with a very large number of servers. Moreover, the large number of decision variables and parameters of the model make this approach unfeasible for a realistic cloud computing resource management strategy.

9.18 Cloud self-organization

Computer clouds are complex systems and should be analyzed in the context of the environment they operate in. The more diverse the environment, the more challenging is the cloud resource management. Cloud self-organization offers a glimpse of hope [332].

Two most important concepts for understanding complex systems are *emergence* and *self-organization*. *Emergence* lacks a clear and widely accepted definition, but it is generally understood as *a property of a system that is not predictable from the properties of individual system components*. There is a continuum of emergence spanning multiple scales of organization. Simple emergence occurs in systems at, or near thermodynamic equilibrium while complex emergence occurs only in nonlinear systems driven far from equilibrium [224].

Physical phenomena which do not manifest themselves at microscopic scales but occur at macroscopic scale are manifestations of emergence, e.g., temperature is a manifestation of microscopic behavior of large ensembles of particles. For such systems at equilibrium, the temperature is proportional to the average kinetic energy per degree of freedom. This is not true for ensembles of a small number of particles. Even the laws of classical mechanics can be viewed as limiting cases of quantum mechanics applied to large masses.

Emergence is critical for complex systems such as financial systems, the air-traffic system, and the power grid. The May 6, 2010, event when the Dow Jones Industrial Average dropped 600 points in a short period of time was a manifestation of emergence. The failure of the trading systems is attributed to interactions of trading systems developed independently and owned by organizations which work together but are motivated by self interest.

A recent paper [453] points out that dynamic coalitions of software-intensive systems used for financial activities pose serious challenges because there is no central authority, and there are no means to control the behavior of individual trading systems. The failures of the power grid (for example, the northeast blackout of 2003) can also be attributed to emergence. Indeed, during the first few hours of this event, the cause of the failure could not be identified due to the large number of independent systems involved. It was established only later that multiple causes, including the deregulation of the electricity market and the inadequacy of transmission lines of the power grid, contributed to this failure.

Informally, self-organization means synergistic activities of elements when no single element acts as a coordinator and the global patterns of behavior are distributed [191, 441]. The intuitive meaning

Table 9.9 Attributes associated with self-organization and complexity.

Simple systems; no self-organization	Complex systems; self-organization
Mostly linear	Non-linear
Close to equilibrium	Far from equilibrium
Tractable at component level	Intractable at component level
One or few scales of organization	Many scales of organization
Similar patterns at different scales	Different patterns at different scales
Do not require a long history	Require a long history
Simple emergence	Complex emergence
Unaffected by phase transitions	Affected by phase transitions
Limited scalability	Scale-free

of self-organization is captured by the observation of Alan Turing [482]: “global order can arise from local interactions.”

Self-organization is prevalent in nature; for example, in chemistry this process is responsible for molecular self-assembly, for self-assembly of monolayers, for the formation of liquid and colloidal crystals, and in many other instances. Spontaneous folding of proteins and other biomacromolecules, formation of lipid bilayer membranes, the flocking behavior of different species, creation of structures by social animals, are all manifestation of self-organization of biological systems. Inspired by biological systems, self-organization was proposed for organization of computing and communication systems [243], including sensor networks [328], for space exploration [239], and economical systems [284].

Generic attributes of complex systems exhibiting self-organization are summarized in Table 9.9. Nonlinearity of physical systems used to build computing and communication systems has countless manifestations and consequences. For example, when the clock rate of a microprocessor doubles, the power dissipation increases from $2^2 = 4$ to $2^3 = 8$ times, depending of the solid state technology used. This means that the heat removal system of much faster microprocessors has to use a different technology when we double the speed.

Nonlinearity is ultimately the reason why recently we have seen the clock rate of general-purpose microprocessors increasing only slightly. The increased number of transistors postulated by Moore’s law is now used for multicore processors. This example illustrates also the so called *incommensurate scaling*, another attribute of complex systems. Incommensurate scaling means that when the size of the system, or when one of its important attributes such as speed increases, different system components are subject to different scaling rules.

The fact that computing and communication systems operate far from equilibrium is clearly illustrated by the traffic carried out by the Internet; there are patterns of traffic specific to the time of the day, but there is no steady-state. The many scales of the organization and the fact that there are different patterns at different scales is also clear in the Internet which is a collection of networks where, in turn, each network is also a collection of smaller networks, each one with its own specific traffic patterns.

The concept of *phase transition* comes from thermodynamics and describes the transformation, often discontinuous, of a system from one phase/state to another, as a result of a change in the environment. Examples of phase transitions are: *freezing*, transition from liquid to solid, and its reverse, *melting*; *deposition*, transition from gas to solid, and its reverse, *sublimation*; *ionization*, transition from gas to plasma, and its reverse, *recombination*.

Phase transitions can occur in computing and communication systems due to avalanche phenomena, when the process designed to eliminate the cause of an undesirable behavior leads to a further deterioration of the systems state. A typical example is thrashing due to competition among several memory-intensive processes that lead to excessive page faults.

Another example is an acute congestion which can cause a total collapse of a network; the routers start dropping packets, and, unless congestion avoidance and congestion control means are in place and operate effectively, the load increases as senders retransmit packets and the congestion increases. To prevent such phenomena some form of *negative feedback* has to be built into the system.

Scalability, the ability of the system to grow without affecting its global function(s), is a defining attribute of *self-organization*. Complex systems encountered in nature, or man-made, exhibit an intriguing property, *scale-free organization* [46,47]. This property reflects one of the few attributes of self-organization that can be precisely quantified.

Scale-free organization can be best explained in terms of the network model of the system, a random graph with vertices representing the entities and the links representing the relationships among them. In a scale-free organization the probability $P(m)$ that a vertex interacts with m other vertices decays as a power law $P(m) \approx m^{-\gamma}$ with γ a real number, regardless of the type and function of the system, the identity of its constituents and the relationships between them [67].

Empirical data available for social networks, power grids, the web, or the citation of scientific papers, confirm this trend. As an example of a social network, consider the collaborative graph of movie actors where links are present if two actors were ever cast in the same movie; in this case $\gamma \approx 2.3$. The power grid of the Western US has some 5 000 vertices representing power generating stations and in this case $\gamma \approx 4$.

The exponent of the World Wide Web scale-free network is $\gamma \approx 2.1$. This means that the probability that m pages point to one page is $P(m) \approx m^{-2.1}$ [47]. Recent studies indicate that $\gamma \approx 3$ for the citation of scientific papers. The larger the network, the closer a power law with $\gamma \approx 3$ approximates the distribution [46].

9.19 Cloud interoperability

Vendor lock-in is a concern, therefore cloud interoperability is a topic of great interest for the cloud community [314,337]. This section addresses several questions: What are the challenges? What is realistic to expect now? What could be done in the future for cloud interoperability? It makes only sense to discuss interoperability of PaaS and IaaS cloud delivery models; the expectation that SaaS services offered by one CSP will be offered by others, e.g., that Google's Gmail will be supported by Amazon, is not realistic.

Now, a workload can migrate from one server to another server in the same data center or among data centers of the same CSP. Migrating a workload to a different CSP is not feasible at this time. To use multiple clouds, data must be replicated and application binaries must be created for all targeted clouds. This is a costly proposition, thus, unfeasible in practice.

There is already a significant body of work on cloud standardization carried out at NIST, but it may take some time before the standards are adopted. First, cloud computing is a fast-changing field and early standardization would hinder progress and slowdown or stifle innovation. It is also likely that the CSPs will resist the standardization efforts.

CSPs are adamant to share information about internal specifications of the software stack, their policies, the mechanisms implementing these policies, and data formats. Each CSP is confident that such information gives it an advantage over the competition. There are also technical reasons why cloud interoperability poses a fair number of challenges, some insurmountable due to current limitations of computing and communication technologies.

So why a complex system such as the Internet is so successful while the development of an Intercloud, a worldwide organization allowing CSPs to share load is so challenging? The Internet is a network of networks, and its architecture is based on two simple ideas: (i) every communicating entity must be identified by an address, thus, a host at the periphery of the Internet, or a router at its core must have one or more IP address; (ii) data sent should be able to reach its destination in this maze of networks, therefore each network should route packets using the same protocol, the IP.

The function of a digital network, regardless of its physical substrate used for communication, is to transport bits of data regardless of what their provenance, music, voice, images, data collected by a sensor, text, or any other conceivable type of information. To make matters even easier, these bits can be packaged together in blocks of small, large, medium, or of any desirable size and can be repackaged whenever the need arises. All that matters is to deliver these bits from a source to a destination in a finite amount of time or in a very short time if the application so requires.

Things cannot be more different for a SuperCloud. First of all, most applications running on clouds need a large volume of data as input to produce results. Transferring say 1 TB of data over a 10 Gbps network takes 8×10^5 seconds, slightly less than a day. Increasing the network speed by an order of magnitude will still require a few hours to transfer this relatively modest volume of data for most cloud applications. Often, we need to transfer more than 1 TB and it is unlikely that Internet speeds of 100 Gbps will soon be available to connect data centers to one another.

What we expect from a SuperCloud is different from what is expected from the Internet where the only function required is to transport data and where all routers, regardless of their architecture, run software implementing the IP protocol. On the other hand, the spectrum of computations done on a cloud is extremely broad, and the heterogeneity of the cloud infrastructure cannot be ignored. Clouds use processors with different architectures, different configurations of cache, memory, and secondary storage, support different operating systems, and use different hypervisors.

The architecture of the server matters; one can only execute code on a server with the same ISA as the one the code was compiled for. The operating system running on a server matters because user code makes system calls to carry out privileged operations and a binary running under one OS cannot be migrated to another OS. The hypervisor running on the server matters because each hypervisor supports only a set of operating systems.

Fortunately, there is a glimpse of hope. A VM including the OS and the application can be migrated to a system with similar architecture and the same hypervisor. Nested virtualization, discussed in Section 5.10, allows a hypervisor to run another hypervisor and this idea discussed later in the section adds to the degrees of freedom for VM migration. Container technologies, such as Docker and LXC, are useful but one cannot move a Docker container from one host to another. What can be done to preserve data that an application has created inside the container is to commit the changes in the container to an image using Docker *commit*, move the image to a new host, and then start a new container with Docker *run*. Moreover, a Docker container is intended to run a single application. There are Docker containers to run an application such as MySQL. A new back-end Docker engine the *libcontainer* can run any application. LXC containers run an application under an instance of Linux.

9.20 Further readings

Cloud resource management poses new and extremely challenging problems, so there should be no surprise that this is a very active area of research. A fair number of papers, including [96], [432], and [27], are dedicated to different resource management policies. Several papers are concerned with QoS [171] and Service Level Agreements [199]. SLA-driven capacity management and SLA-based resource allocation policies are covered in [6] and [29], respectively. [169] analyzes performance monitoring for SLAs. Dynamic request scheduling of applications subject to SLA requirements is presented in [68]. Clouds QoS is analyzed in [68]. Semantic resource allocation using a multiagent system is discussed in [161].

The autonomic computing era is presented in [187]. Energy-aware resource allocation in autonomic computing is covered in [30]. Policies for autonomic computing based on utility functions are analyzed in [269]. Coordination of multiple autonomic managers and power-performance tradeoffs are dissected in [268]. Autonomic management of cloud services subject to availability guarantees is presented in [11]. The use of self-organizing agents for service composition in cloud computing is the subject of [219].

An authoritative reference on fault-tolerance is [38]; applications of control theory to resource allocation, discussed in [156] and [91], cover resource multiplexing in data centers. Admission control policies are discussed in [216]. Optimal control problems are analyzed in [226], and system testing is covered in [231]. Verification of performance assertion on the cloud is the subject of [278]. Power and performance management are the subject of [287] and performance management for cluster based web services is covered in [384]. Autonomic management of heterogeneous workloads is discussed in [467] and application placement controllers is the topic of [470]. Applications of pattern matching for forecasting on-demand resources needs is discussed in [88]. Economic models for resource allocations are covered in [327, 330], and [431].

Scheduling and resource allocation are also covered by numerous papers: a batch queuing system on clouds with Hadoop and HBase is presented in [538]; data flow-driven scheduling for business applications is covered in [151]. Scalable thread scheduling is the topic of [516]. Reservation-based scheduling is discussed in [126]. Anti-caching in database management is the subject of [132]. Scheduling of real time services in cloud computing is presented in [310]. The OGF (Open Grid Forum) OCCI (Open Cloud Computing Interface) is involved in the definition of virtualization formats and APIs for IaaS. Flexible memory exchange, bag of tasks scheduling and distributed low latency scheduling are covered in [370], [377], and [382], respectively. Reference [421] analyzes capacity management for pools of resources.

9.21 Exercises and problems

Problem 1. Analyze the benefits and the problems posed by resource management implementations based on control theory, machine learning, utility-based, and market-oriented policies.

Problem 2. Can optimal strategies for the five classes of policies, admission control, capacity allocation, load balancing, energy optimization, and QoS guarantees be actually implemented in a cloud? The term “optimal” is used in the sense of control theory. Support your an-

swer with solid arguments. Optimal strategies for one may be in conflict with optimal strategies for one or more of the other classes. Identify and analyze such cases.

- Problem 3.** Analyze the relationship between the scale of a system and policies and mechanisms for resource management. Consider also the geographic scale of the system.
- Problem 4.** What are the limitations of the control theoretic approach discussed in Section 9.13? Do the approaches discussed in Sections 9.14 and 9.15 remove or relax some of these limitations? Justify your answers.
- Problem 5.** Multiple controllers are probably necessary due to the scale of the cloud. Is it beneficial to have system and application controllers? Should the controllers be specialized? For example, some to monitor performance, others to monitor power consumption. Should all the functions we want to base the resource management policies on be integrated in a single controller and one such controller be assigned to a given number of servers, or to a geographic region? Justify your answers.
- Problem 6.** In a scale-free network, the degrees of the nodes have an exponential distribution. A scale-free network could be used as a virtual network infrastructure for cloud computing. *Controllers* represent a dedicated class of nodes tasked with resource management; in a scale-free network nodes with a high connectivity can be designated as controllers. Analyze the potential benefit of such a strategy.
- Problem 7.** Use the start-time fair queuing (SFQ) scheduling algorithm to compute the virtual startup and the virtual finish time for two threads a and b with weights $w_a = 1$ and $w_b = 5$ when the time quantum is $q = 15$ and thread b blocks at time $t = 24$ and wakes up at time $t = 60$. Plot the virtual time of the scheduler function of the real time.
- Problem 8.** Apply the borrowed virtual time (BVT) scheduling algorithm to the problem in Example 2 of Section 9.7 but with a time warp of $W_c = -30$.
- Problem 9.** In Section 9.11, we introduced the concept of energy-proportional systems, and we saw that different system components have different dynamic ranges. Sketch a strategy to reduce the power consumption in a lightly loaded cloud and discuss the steps for placing a computational server in a standby mode and then for bringing it back up to an active mode.
- Problem 10.** Overprovisioning is the reliance on extra capacity to satisfy the needs of a large community of users when the average-to-peak resource demand ratio is very high. Give an example of a large-scale system using overprovisioning and discuss if overprovisioning is sustainable in that case and what are the limitations of it. Is cloud elasticity based on overprovisioning sustainable? Give the arguments to support your answer.

Concurrency and cloud computing

10

Concurrency is at the heart of cloud computing, large workloads generated by many cloud applications run concurrently on multiple instances, taking advantage of ample resources provided by a cloud infrastructure. The practical motivations for concurrent execution are to: (i) overcome physical limitations of any single system by distributing the workload among several servers and (ii) significantly reduce the completion time of a computation.

Leslie Lamport began his 2013 Turing Lecture dedicated to Edsger Dijkstra [294] with the observation that concurrency has been known by several names: “I don’t know if concurrency is a science, but it is a field of computer science. What I call concurrency has gone by many names, including parallel computing, concurrent programming, and multiprogramming. I regard distributed computing to be part of the more general topic of concurrency.”

Is there a distinction between *concurrency* and *parallel processing*? According to some, concurrency describes the necessity that multiple computations are executed at the same time, while parallel processing implies a solution, there are multiple physical systems capable of carrying out the computations required by these activities at the same time, i.e., concurrently. Concurrency emphasizes cooperation and interference among activities, while parallel processing aims to shorten the completion time of the set of activities, and it is hindered by cooperation and activity interference.

Execution of multiple activities in parallel can proceed either quasi-independently, or tightly coordinated with an explicit communication pattern. In either case, some form of communication is necessary for coordination of concurrent activities. Coordination complicates the description of a complex activity because it has to characterize the work done by individual entities working in concert, as well as the interactions among them.

Communication affects overall efficiency of concurrent activities and could significantly increase the completion time of an application with a large number of concurrent tasks communicating frequently. Furthermore, communication requires prior agreement on the communication discipline described by a communication protocol. Measures to ensure that anomalies, e.g., deadlocks, do not affect the overall orchestration required by the application are also necessary.

The chapter starts with a brief discussion of concurrency’s enduring challenges in Section 10.1 and continues with an overview of concurrent execution of communicating processes in Section 10.2, followed by a discussion of computational models in Section 10.3. BSP, a bridging hardware-software model designed to avoid logarithmic losses of efficiency in parallel processing, and its version for a multicore computational model are covered in Sections 10.4 and 10.5. Petri nets, discussed in Section 10.6, are intuitive models able to describe concurrency and conflict.

The concept of process state, critical for understanding concurrency, is covered in Section 10.7. Many functions of a computer cloud require information about process state. For example, controllers

for cloud resource management discussed in Chapter 9 require accurate state information. Process coordination is analyzed in Section 10.8, while Section 10.9 presents logical clocks and message delivery rules in an attempt to bridge the gap between the abstractions used to analyze concurrency and the physical systems.

The concept of consistent cuts and distributed snapshots are at the heart of *checkpoint–restart* procedures for long-lasting computations. Checkpoints are taken periodically in anticipation of the need to restart a software process when one or more systems fail; when a failure occurs, the computation is restarted from the last checkpoint, rather than from the beginning. These concepts are discussed in Sections 10.10 and 10.11. Atomic actions, consensus protocols, and load balancing are covered in Sections 10.12, 10.13, and 10.14, respectively. Section 10.15 presents multithreading and concurrency in Java, FlumeJava, and Apache Crunch.

This chapter reviews theoretical foundations of important algorithms at the heart of system and application software. The concepts introduced in the next sections help us better understand cloud resource management policies and the mechanisms implementing these policies. For a deeper understanding of the many subtle concepts related to concurrency, the reader should consult the classical references discussed at the end of the chapter.

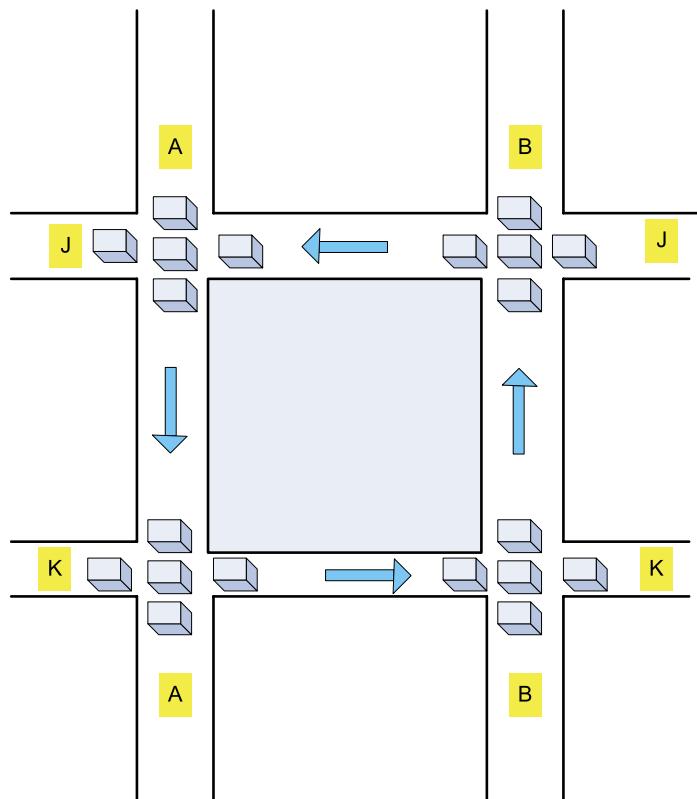
10.1 Enduring challenges

Concurrency is a very broad subject, and in this chapter we restrict the discussion to topics closely related to cloud computing. We start with a gentle introduction to some of the enduring challenges posed by concurrency and coordination. Coordination starts with resource allocation to the units carrying out individual tasks and workload distribution among these units. The initial phase is followed by communication during execution of the tasks, and finally, by the assembly of individual results. Coordination is ubiquitous in our daily life, and the lack of coordination has implications for the results. For example, Fig. 10.1 shows that lack of coordination and disregard of the rules regarding the management of shared resources lead to traffic deadlock, an unfortunate phenomenon we often experience.

*Synchronization*¹ is an important aspect of concurrency illustrated by the dining philosophers problem, see Fig. 10.2. Five philosophers sitting at a table alternately think and eat. To eat a philosopher needs the two chopsticks placed left and right of her plate. After finishing eating she must place the chopsticks back on the table to give a chance to her left and right neighbors to eat. The problem is nontrivial, the naive solution when each philosopher picks up the chopstick to the left and waits for the one to the right to become available, or vice versa, fails because it allows the system to reach a deadlock when no progress is possible. Deadlock would lead to philosopher starvation, a state to be avoided.

This problem captures critical aspects of concurrency, such as mutual exclusion and resource starvation discussed in this chapter. Edsger Dijkstra proposed the following solution formulated in terms of resources in general:

¹ Synchronization refers to one of two distinct concepts: process synchronization and data synchronization. Process synchronization, discussed in this section, means that multiple processes handshake at a certain point to reach an agreement or to commit to a sequence of actions. Data synchronization means keeping multiple copies of a dataset in coherence with one another, or maintaining data integrity.

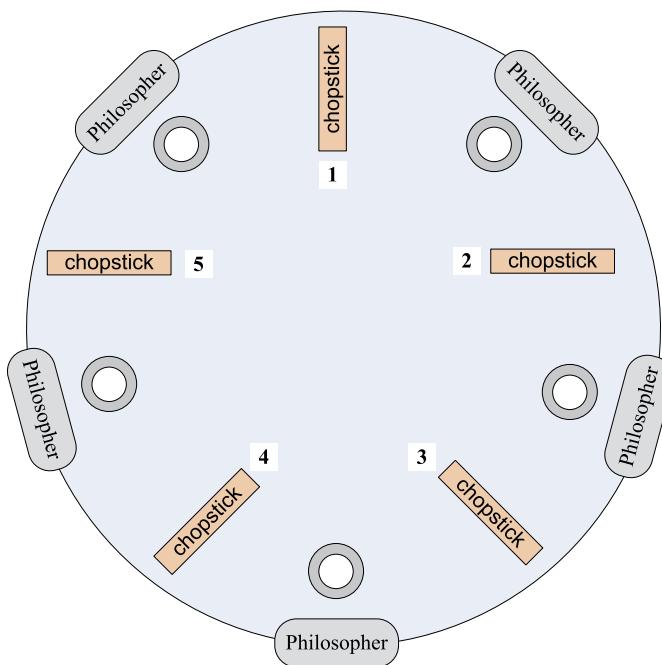
**FIGURE 10.1**

The consequences of the lack of coordination. The absence of traffic lights or signs in intersections causes a traffic jam. Blocking the intersections, A-J, A-K, B-J, and B-K, shared resources for North-South and East-West traffic flows, creates deadlocks.

- i. Assign a partial order to the resources.
- ii. Impose the rule that resources must be requested in order.
- iii. Impose the rule that no two resources unrelated by order will ever be used by a single unit of work at the same time.

In the dining philosopher problem, the resources are the chopsticks, numbered 1 through 5, and each unit of work, i.e., philosopher, will always pick up the lower-numbered chopstick first and then the higher-numbered one next to her. The order in which each philosopher puts down the chopsticks does not matter.

This solution avoids starvation. If four of the five philosophers pick up their lower-numbered chopstick at the same time, only the highest-numbered chopstick will be left on the table. Therefore, the fifth philosopher will not be able to pick up a chopstick. Only one philosopher will have access to that

**FIGURE 10.2**

Dining philosophers problem. To avoid deadlock Dijkstra's solution requires numbering of chopsticks and two additional rules.

highest-numbered chopstick, so she will be able to eat using two chopsticks. Reference [309] presents a solution of the dining philosopher problem based on a Petri Net model.

The division of work comes naturally when some activities of a complex task require special competence and can only be assigned to uniquely qualified entities. In other cases, all entities have the same competence and the work should be assigned based on the individual ability to carry out a task more efficiently. Balancing the workload could be difficult in all cases because some activities may be more intense than others.

Though concurrency reduces the time to completion, it can negatively affect the efficiency of individual entities involved. Sometimes, a complex task consists of multiple stages and transitions from one stage to the next can only occur when all concurrent activities in one stage have finished their assigned work. In this case, the entities finishing early have to wait for the others to complete, an effect called *barrier synchronization*.

This discussion shines some light on the numerous challenges inherent to concurrency. Many computational problems are rather complex, and concurrency has the potential to greatly affect our ability to compute efficiently. This motivates our interest in concurrency and its applications to cloud computing.

Parallel and distributed computing exploit concurrency and have been investigated since mid-1960s. Parallel processing refers to concurrent execution on a system with a large number of processors, while

distributed computing means concurrent execution on multiple systems, often at different locations. Communication latency is considerably lower in the first case, while distributed computing could only be efficient for coarse-grained parallel applications when concurrent activities seldom communicate with one another. Metrics such as execution time, speedup, and processor utilization discussed in Chapter 3 characterize how efficiently a parallel or distributed system can process a particular application.

Topics discussed in this chapter, such as computational models, checkpointing, atomic actions, and consensus algorithms, are relevant to both parallel and distributed computing. Multithreading is more relevant to parallel processing, while load balancing is particularly important to distributed systems.

This distinction is blurred for computer clouds because their infrastructure consists of millions of servers at one data center, possibly linked by high-speed networks with computers at another data center of the same cloud service provider. Nevertheless, communication latency is a matter of concern in cloud computing as we can see in Chapters 4, 9, and 11.

10.2 Communication and concurrency

Concurrency implies cooperative execution of multiple processes/threads in parallel. Concurrency can be exploited for minimizing the completion time of a task, while maximizing efficiency of the computing substrate. Computing substrate is a generic term for the physical systems used for application processing. To analyze the benefits of concurrency, we consider the decomposition of a computation into virtual tasks and relate them to the physical processing elements of the computing substrate.

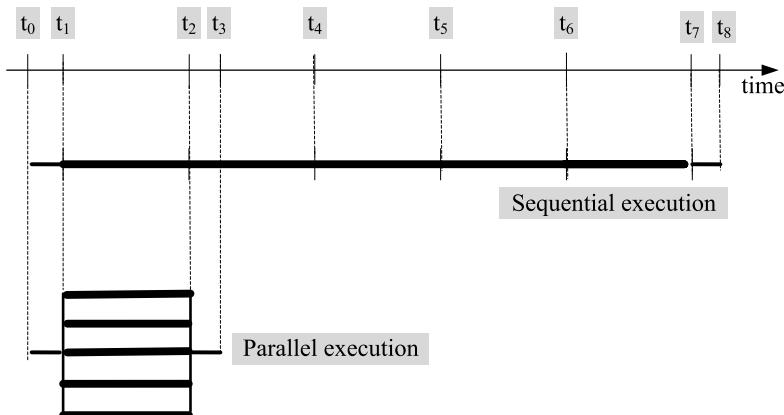
The larger the cardinality of the set of virtual tasks, the higher the degree of concurrency, thus the higher the potential speedup. A parallel algorithm can then be implemented by a parallel program able to run on a system with multiple execution units. A *process* is a program in execution and requires an *address space* hosting the code and the data of an application. A *thread* is a lightweight execution unit running in the address space of a process.

The *speedup* of concurrent execution of an application quantifies the effect of concurrent execution, and it is defined as the ratio of the completion time of sequential execution of the task versus the concurrent execution completion time. For example, Fig. 10.3 illustrates concurrent execution of an application in which the workload is partitioned and assigned to five different processors or cores running concurrently. The application is the conversion of 5×10^6 images from one format to another. This is an *embarrassingly parallel* application because the five threads running on five cores, each processing 10^6 images, do not communicate with one another. The speedup S for the example in Fig. 10.3 is

$$S = \frac{t_8 - t_0}{t_3 - t_0} \approx 5. \quad (10.1)$$

There are two sides of concurrency, algorithmic or logical concurrency, discussed in this chapter, and physical concurrency, discovered and exploited by the software and the hardware of the computing substrate. For example, a compiler can unroll loops and optimize sequential program execution, and a processor core may execute multiple program instructions concurrently, as discussed in Chapter 3.

Concurrency is intrinsically tied to communication; concurrent entities must communicate with one another to accomplish the common task. The corollary of this statement is that communication and

**FIGURE 10.3**

Sequential versus parallel execution of an application. The sequential execution starts at time t_0 , goes through a brief initialization phase until time t_1 and then starts the actual image processing. When all images have been processed, it enters a brief termination phase at time t_7 , and finishes at time t_8 . The concurrent execution has its own brief initialization and termination phases; the actual image processing starts at time t_1 and ends at time t_2 . The results are available at time $t_3 \ll t_8$. The speedup is close to the maximum speedup.

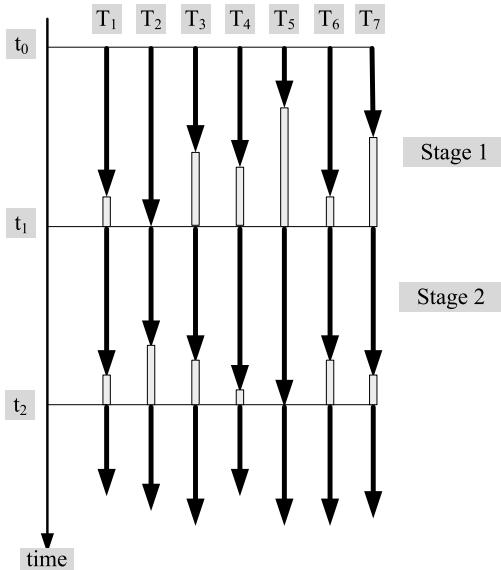
computing are deeply intertwined. Explicit communication is built into the blueprint of concurrent activities; we shall call it *algorithmic communication* for lack of a better term. A multithreaded application consists of one or more *thread groups*, sets of threads that work in concert and have to coordinate with each other.

Sometimes, a multithreaded computation consists of multiple stages when concurrently running threads of a thread group cannot continue to the next stage until all of them have finished the current one. This leads to inefficiencies, as shown in Fig. 10.4, that illustrates the effects of *barrier synchronization*. Such thread groups (or process groups) have to be scheduled together to minimize the blocking time, a process called *co-scheduling* [34].

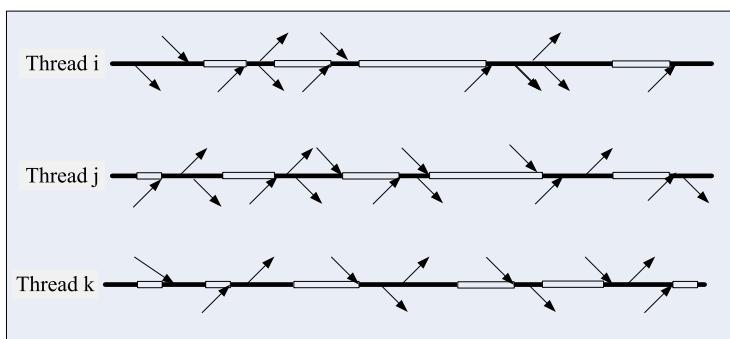
Communication is a more intricate process than the execution of a computational step on a machine that rigorously follows a sequence of instructions from its own repertoire. Besides the sender and the receiver(s), communication involves a third party, a communication channel that may, or may not be reliable. Therefore, the two or more communicating entities have to agree on a communication protocol consisting of multiple messages. *Communication complexity* reflects the amount of communication that the participants of a communication system need to exchange to be able to perform certain tasks [525].

Fig. 10.5 illustrates the case when short bursts of computations alternate with relatively long periods when a thread is blocked, waiting for messages from other threads, the so-called *fine-grained* parallelism. The opposite is *coarse-grained* parallelism when little or no communication among concurrent threads takes place, as can be seen in Fig. 10.3.

Communication speed is considerably slower than computation speed. A processor could execute billions of instructions during the time it is blocked waiting to receive a message. It is thus expected that an application exhibiting fine-grained parallelism, i.e., communicating frequently, experiences long

**FIGURE 10.4**

Barrier synchronization. The computation involving a group of seven threads running concurrently consists of multiple stages; all threads in the thread group must finish execution of one stage before proceeding to the next. All threads start execution of Stage 1 at time t_0 . Threads T_1 , T_3 , T_4 , T_5 , T_6 , and T_7 finish early and have to wait for thread T_2 before proceeding to Stage 2 at time t_1 . Similarly, threads T_1 , T_2 , T_3 , T_4 , T_6 , and T_7 have to wait for thread T_5 , before proceeding to the next stage at time t_2 . Black arrows indicate running threads, while white bars represent blocked tasks, waiting to proceed to the next stage.

**FIGURE 10.5**

Fine-grained parallelism. Short bursts of computations of three concurrent threads are interspersed with blocked periods when a thread is waiting for messages from other threads. Solid black bars represent running threads, while white bars represent blocked threads waiting for messages. Arrows represent messages sent or received by a thread.

blocking periods. Intensive communication can slow down considerably a group of concurrent threads of an application, especially in case of multistage computations with frequent barrier synchronizations.

The word *message* should be understood in an information theoretical sense, rather than the more narrow meaning used in computer networks context. Embarrassingly parallel activities can proceed without message exchanges and enjoy linear or even superlinear speedup, while in all other cases, the completion time of communicating activities is increased because the communication bandwidth is significantly lower than the processor bandwidth.

Nonalgorithmic communication is required by the organization of a computing substrate consisting of various components that need to act in concert to carry out the required task. For example, in a system consisting of multiple processors and memories, a thread running on one processor may require data stored in the memory of another one. The nonalgorithmic communication, unavoidable in a distributed systems, occurs as a side effect of thread scheduling and data distribution strategies and can dramatically reduce the benefits of concurrency.

Spatial and temporal locality affect the efficiency of a parallel computation. A sequence of instructions or data references enjoy *spatial locality* if the items referenced in a narrow window of time are close in space, e.g., are stored in nearby memory addresses, or nearby sectors on a disk. A sequence of items enjoy *temporal locality* if accesses to the same item are clustered in time.

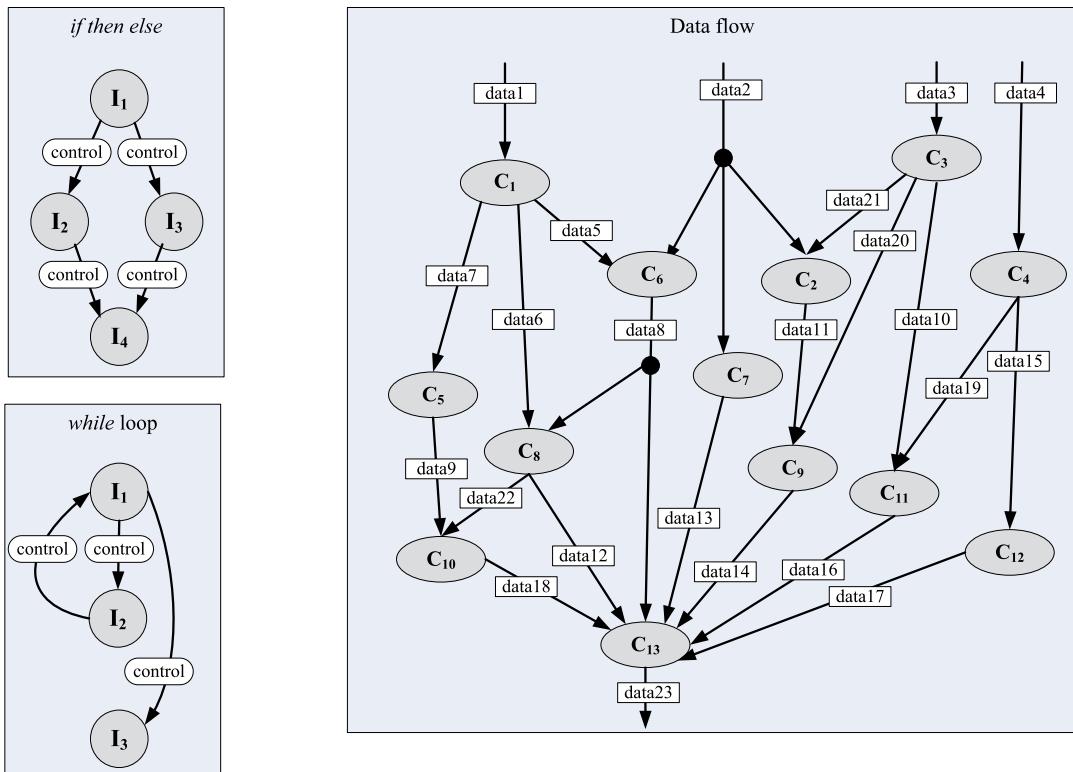
Nonalgorithmic communication complexity may decrease due to locality; it is more likely that a virtual task may find the data it needs in the memory of the physical processor where it runs. Also, the efficiency of the computing substrate may increase because at any given time the scheduler may find sufficient read-to-run virtual tasks to keep the physical processing elements busy. The scale of a system amplifies the negative effects of both algorithmic and nonalgorithmic communication. The interaction between the virtual and physical aspects of concurrency gives us a glimpse at the challenges faced by a computational model of concurrent activities.

Today's computing systems implement two computational models, one based on *control flow* and the other one on *data flow*. The ubiquitous von Neumann model for sequential execution implements the former; at each step, the algorithm specifies the step to be executed next. In this model, concurrency can only be exploited through the development of a parallel algorithm reflecting the decomposition of a computational task into processes/threads that can run concurrently, while exchanging messages to coordinate their progress. In the case of control flow, the *if then else* construct is shown by the top graph and the *while* loop on the bottom graph in Fig. 10.6. Instructions I_2 and I_3 of the *if then else* construct can run concurrently only after instruction I_1 finishes. Instruction I_4 can only be executed when I_2 and I_3 finish execution.

The *data flow* execution model is based on an implicit communication model, where a thread starts execution as soon as its input data become available. The advantages of this execution model are self-evident; it effortlessly extracts maximum amount of parallelism from any computation.

The data flow model example in Fig. 10.6 shows a maze of computations C_1, \dots, C_{13} with complex dependencies. The data flow model allows all computations to run as soon as their input data become available. For example, C_1 , C_3 , and C_4 start execution concurrently with input *data1*, *data3*, and *data4*, respectively, while C_2 and C_6 wait for completion of C_3 and C_1 , respectively. Finally, C_{13} can only start when *data18*, *data12*, *data8*, *data13*, *data14*, *data16*, and *data17* are produced by C_{10} , C_8 , C_6 , C_7 , C_9 , C_{11} , and C_{12} respectively.

The time required by computations C_i $1 \leq i \leq 13$ to finish execution and deliver data to the ones waiting depends upon the size of the input data which is not known a priori. To capture the dynamics of

**FIGURE 10.6**

Control flow versus data flow models. Left graphs show the control flow from one instruction I_i , $1 \leq i \leq 4$ to the others for *if then else* and the *while* loop constructs. Either I_2 or I_3 will ever run in the *if then else* construct. The right graph illustrates the data flow model; now the arrival of the input data triggers the execution of each one of the computations C_i , $1 \leq i \leq 13$. For example, C_9 execution is triggered by the arrival of *data11* produced by C_2 and *data20* produced by C_3 .

a computational task, the control flow model would require individual computations to send and receive messages, in addition to sending the data.

There are only a few data flow computer systems in today's landscape, but it is not beyond the realm of possibilities to see some added to the cloud computing infrastructure in the next future. Moreover, some of the frameworks for data processing discussed in Chapters 4 and 11 attempt to mimic the data flow execution model to optimize their performance.

Petri Nets models discussed in Section 10.6 are powerful enough to express either control flow or data flow. In these bipartite graphs one type of vertices, *places*, model system state, while the other type of vertices, *transactions*, model actions. Tokens flowing out of places trigger transactions and end up in other places, indicating a change of system state. Tokens may represent either control or data.

For several decades, concurrency was of interest mostly for systems programming and for high-performance computing in science and engineering. The majority of other application developers were content with sequential processing and expected increased computing power due to faster clock rates. Concurrency is now mainstream due to the disruptive effects of the multicore processor technology, the limitation imposed on the processor clock speed, combined with the insatiable appetite for computing power encapsulated in tiny and energy-frugal computing devices. Writing and debugging concurrent software is considerably more challenging than developing sequential code; it requires a different frame of mind and effective development tools.

The next three sections discuss computational models, abstractions needed to gain insight into fundamental aspects of computing and concurrency.

10.3 Computational models; communicating sequential processes

Several computational models are used in the analysis of algorithms. Each one of the two broad classes of models, *abstract machines* and *decision trees*, specify the set of primitive operations allowed. Turing machines, circuit models, lambda calculus, finite-state machines, cellular automaton, stack machines, accumulator machines, and random access machines are abstract machines used in proofs of computability and for establishing upper bounds on computational complexity of algorithms [439]. Decision tree models are used in proofs of lower bounds on computational complexity.

Computational models. Turing machines are *uniform models of computation*; the same computational device is used for all possible input lengths. Boolean circuits are *non-uniform models of computation*, inputs of different lengths are processed by different circuits. A computational problem \mathcal{P} is associated with the family of Boolean circuits $\mathcal{C} = \{C_1, C_2, \dots, C_n, \dots\}$, where each C_n is a Boolean circuit handling inputs of n bits. A family of Boolean circuits \mathcal{C}_n , $n \in \mathbb{N}$, is *polynomial-time uniform* if there exists a deterministic Turing machine TM, such that TM runs in polynomial time and $\forall n \in \mathbb{N}$ it outputs a description of \mathcal{C}_n on input 1^n .

A finite-state machine (FSM) consists of a logical circuit \mathcal{L} and a memory \mathcal{M} . An execution step with the external input $L^{in} \in \Sigma$ takes the current state $S \in \mathcal{S}$ and uses the logic circuit \mathcal{L} for producing a successor state $S^{new} \in \mathcal{S}$ and an output letter from the same alphabet, Σ , $L^{out} \in \Sigma$. FSMs are used for the implementation of sophisticated processing scheme such as TCP (Transmission Control Protocol) covered in Chapter 6 and used by the Internet protocol stack and the Zookeeper coordination model discussed in Section 11.4.

The operation of a serial computer using two synchronous interconnected FSMs, a central processing unit (CPU) with a small number of registers and a random-access memory, is modeled by a *random-access machine* (RAM). CPU operates on data stored in its registers. *Parallel random-access machine* (PRAM) is an abstract programming model consisting of a bounded set of RAM processors and a common memory of a potentially unlimited size. Each RAM has its own program and program counter and a unique Id. During a PRAM execution step, the RAMs execute synchronously three operations: read from the common memory, perform a local computation, and write to the common memory.

The von Neumann model, based on the von Neumann machine architecture, has been an enduring model of sequential computations. It has continued to allow “a diverse and chaotic software to run efficiently on a diverse and chaotic world of hardware” [485]. The model has endured the rapid pace

of technological changes since 1947 when ENIAC performed the first Monte Carlo simulations for the Manhattan Project [221]. A brilliant feature of the von Neumann model is its clairvoyance, the ability to remain consistent in the face of later concepts such as memory hierarchies, certainly not available in the mid-1940s. This model has been the *zeitgeist*² of computing for more than half a century.

Communicating Sequential Processes (CSP). CSP is a formal language for describing patterns of interaction in concurrent systems proposed by Tony Hoare in 1978 [241]. CSP treats input and output as fundamental programming primitives and includes a simple form of parallel composition based on synchronized communication. CSP programs are a parallel composition of a fixed number of sequential named processes communicating with each other only through synchronous message-passing. Each message is specified by the name of the intended sending or receiving process. The concept of a *process* is a natural abstraction of the way parallel activities behave.

Brookes, Hoare, and Roscoe refined the theory of CSP into a process algebraic form [74]. CSP process algebra includes two classes of primitives, *initial events*, indivisible and instantaneous primitives representing interactions, and *primitive processes*, representing actions/behaviors. Several algebraic operators combine events and processes:

prefix—creates a new process by combining an event and a process, e.g., $a \rightarrow P$ describes a process like P willing to use primitive event a to communicate with the environment.

hiding—makes event a unobservable to primitive process P , $(a \rightarrow P) \setminus \{a\}$.

interleaving—a new process that behaves like primitive process P or Q simultaneously; the operator \parallel represents independent concurrent activities.

deterministic choice—a new process is defined as a choice between two processes P and Q chosen by the environment which communicates its choice via two primitive events a and b , $(a \rightarrow P) \square (b \rightarrow Q)$.

nondeterministic choice—a new process is defined as a choice between two processes P and Q , but the environment has no control over the choice, $(a \rightarrow P) \sqcap (b \rightarrow Q)$. The new process can behave as $(a \rightarrow P)$ or as $(b \rightarrow Q)$.

A fair number of tools to analyze systems described by CSP exist. For example, Failures/Divergence Refinement (FDR) is a family of model checker that converts two CSP process expressions into Labelled Transition Systems and then determines whether one of the processes is a refinement of the other within some specified semantic model; see <http://www.fsel.com/>.

CSP was used in 1998 to analyze the software, consisting of some 22 000 lines of code, for a fault-tolerant computer used to control the assembly, reboot³ operations and flight control, and data management for experiments carried out by the International Space Station. The analysis confirmed that the design was deadlock-free [82].

² The German word “Zeitgeist” literally translated as “time spirit” is used to identify the dominant set of ideas and concepts of the society in a particular field and at a particular time.

³ Reboot is the process of boosting the altitude of an artificial satellite, to increase the time until its orbit will decay and it reenters the atmosphere.

10.4 The bulk synchronous parallel model

Leslie Valiant developed in the early 1990s a bridging hardware-software model designed to avoid logarithmic losses of efficiency in parallel processing [485]. Valiant argues that *parallel and distributed computing should be based on a model emphasizing portability and the algorithms should be parameter-aware and designed to run efficiently on a broad range of systems with a wide variation of parameters*. The algorithms have to be written in a language that can be compiled effectively on a computer implementing the BSP model.

BSP is an unambiguous computational model, including parameters quantifying the physical constraints of the computing substrate. It also includes a nonlocal primitive, the barrier synchronization. The BSP model aims to free the programmer from managing memory, communication, and low-level synchronization, provided that programs are written with sufficient *parallel slackness*. Parallel slackness means hiding communication latency by providing each processor with a large pool of ready-to-run tasks, while other tasks are waiting for either a message or the completion of another operation.

BSP programs are written for v virtual parallel processors running on $p \leq v$ physical processors. This choice permits high-level language compilers to create executables sharing a virtual address space and to schedule and pipeline computation and communication efficiently. The BSP model includes:

- i. Processing elements and memory components.
- ii. A router involved in message exchanges between pairs of components; the throughput of the router is \bar{g} and s is the startup cost.
- iii. Synchronization mechanisms acting every L units of time.

The computation involves *supersteps* and *tasks*. A superstep is an execution unit allocated L units of time and consisting of multiple tasks. Each task is a combination of local computations and message exchanges. The computation proceeds as follows:

1. At the beginning of a superstep, each component is assigned a task.
2. At the end of the time allotted to the superstep, after L units of time since its start, a global check determines if the superset has been completed:
 - 2.1. If so, the next superstep is initiated;
 - 2.2. Else, another period of L units of time is allocated allowing the superset to continue its execution.

Local operations do not automatically slow down other processors. A processor can proceed without waiting for the completion of processes in the router or in other components when the synchronization is switched off. The model does not make any assumption about communication latency. The value of the periodicity parameter L may be controlled, and its lower bound is determined by the hardware, the faster the hardware the lower L , while its upper bound is controlled by the granularity of the parallelism, hence by the software.

The router of the BSP computing engine supports arbitrary *h-relations*, which are supersteps of duration $\bar{g} \times h + s$ when each component sends and receives at most h messages. It is assumed that h is large, and $\bar{g} \times h$ is comparable to s . When $g = 2\bar{g}$ and $\bar{g} \times h \geq s$, an h-relation will require at most $g \times h$ units of time.

Hash functions are used by the model to distribute memory accesses to memory units of all BSP components. The mapping of logical or symbolic addresses to physical ones must be efficient and dis-

tribute the references as uniformly as possible, and this can be done by a pseudo-random mapping. When a superstep requires p random accesses to p components, then a component will get, with high probability, $\log p / \log \log p$ accesses. If a superstep requires $p \log p$ memory accesses, then each component will get not more than about $3 \log p$ accessed, and these can be sustained by a router in the optimal bound $\mathcal{O}(\log p)$.

The simulation of a parallel program with $v \geq p \log p$ virtual processors on a BSP computer with p components is discussed next. Each BSP computer component consists of a processor and a memory. Each one of the p components of the BSP computer is assigned v/p virtual processors. The BSP will then simulate one step of a virtual processor in one superstep, and each one of the p memory modules will get v/p references if the memory references are evenly spread. The BSP will execute a superstep in an optimal time of $\mathcal{O}(v/p)$.

The multiplication of two $n \times n$ matrices A and B using a standard algorithm is a good example of simulation on a BSP machine with $p \leq n^2$ processors. Each processor is assigned an $n/\sqrt{p} \times n/\sqrt{p}$ submatrix and receives $n/\sqrt{p} \times n/\sqrt{p}$ rows of A and $n/\sqrt{p} \times n/\sqrt{p}$ columns of B . Thus, each processor will carry out $2n^3/p$ additions and multiplications and send $2n^2/\sqrt{p} \leq 2n^3/p$ messages. If each processor sends $2n^2/\sqrt{p}$ messages, then the running time is affected only by a constant factor.

Concurrency can be implemented efficiently by replicating data when h is small. In the matrix multiplication example, $2n^2/p$ elements of matrices A and B can be distributed to each one of the p processors, and each processor replicates each of its elements \sqrt{p} times and sends them to the \sqrt{p} processors that need the entries. The number of messages sent by each processor is $2n^2\sqrt{p}$. When $g = \mathcal{O}(n\sqrt{p})$ and $L = \mathcal{O}(n^3/p)$, we have the optimal running time of $\mathcal{O}(n^3/p)$. In addition to matrix multiplication, it is shown that several other important algorithms can be implemented efficiently on a BSP computer.

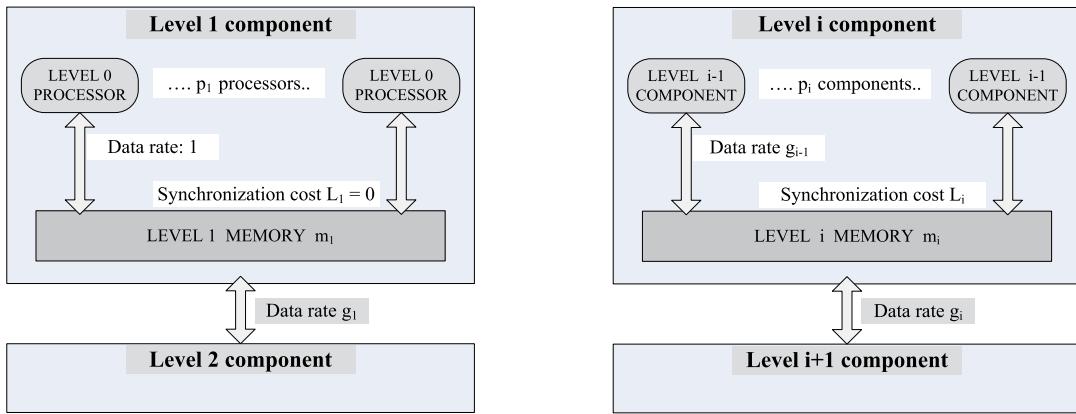
Valiant concludes that the BSP model helps programmability for computations with sufficient slack because the memory and communication management required to implement a virtual shared memory can be achieved with only a constant factor loss in processor utilization [485]. Moreover, the BSP model can be implemented efficiently for different communication technologies and interconnection network topologies, e.g., for a hypercube interconnection network.

10.5 A model for multicore computing

Extracting an optimal multicore processor performance is very challenging. Most of these challenges are intrinsically related to the complexity and diversity of the computational engines because the performance of an application on one system may not translate on high performance on another. Developing parallel algorithms is in itself nontrivial for many applications of interest and the application performance may not scale with the problem size for various computation substrates.

The Multi-BSP [486] is a hierarchical multicore computing model with an arbitrary number of levels. The computing substrate of the model includes multiple levels of cache, as well as the size of the memory. A *depth-d* model is a tree of depth d with caches and memory as the internal nodes of the tree and the processors at the leaves.

A depth- d model requires $4d$ parameters. Level i , $1 \leq i \leq d$, has four parameters (p_i, g_i, L_i, m_i) that quantify the number of subcomponents, the communication bandwidth, the synchronization cost, and the memory/cache size, respectively. The model captures the unavoidable communication cost due

**FIGURE 10.7**

Multi-BSP model. A $depth-d$ model is a tree of depth d with caches and memory as the internal nodes of the tree and the processors at the leaves. p_i —the number of $(i - 1)$ level components inside an i -th level component; g_i —the communication bandwidth; L_i —the cost for barrier synchronization for a level- i superstep; m_i —the number of words of memory inside an i -th level component.

to latency L_i and bandwidth g_i as seen in Fig. 10.7. An in-depth description of level i parameters follows:

- p_i —the number of $(i - 1)$ level components inside an i -th level component. The 1st level components consist of p_1 raw processors; a computation step on a word in level 1 memory represents one unit of time.
- g_i —the communication bandwidth, the ratio of the number of operations of a processor to the number of words transmitted between the memories of a component at level i and its parent component at level $(i + 1)$, in a unit of time. A word represents the data the processor operates on. Level 1 memories can keep up with the processors; thus their data rate, g_0 , is one.
- L_i —the cost for barrier synchronization for a level i superstep. The barrier synchronization is between the subcomponents of a component; there is no synchronization across separated branches in the component hierarchy.
- m_i —the number of words of memory inside an i -th level component that is not inside any $(i - 1)$ level component.

The number of processors in a level- i component is $P_i = p_1 \cdot p_2 \cdot \dots \cdot p_i$. The number of level- i components in a level- j component is $Q_{i,j} = p_{i+1} \cdot p_{i+2} \cdot \dots \cdot p_j$, and the number in the whole system is $Q_{i,d} = Q_i = p_{i+1} \cdot p_{i+2} \cdot \dots \cdot p_d$. The total memory at level- i component is $M_i = m_i + p_i \cdot m_{i-1} + p_{i-1} \cdot p_i \cdot m_{i-2} + \dots + p_2 \cdot \dots \cdot p_{i-1} \cdot p_i m_1$. The cost of communication from level 1 to outside level i is $G_i = g_i + g_{i-1} + \dots + g_1$.

A level- i superstep is a construct within a level- i component that allows each of p_i -level $(i - 1)$ components to execute independently until they reach a barrier. Only after reaching the barrier can all exchange information with the m_i memory of the level- i component at a communication cost of $(g_i -$

1). The communication cost is $m_{i-1}g_{i-1}$ where m_i is the number of words communicated between the memory at the i -th level and its level $(i - 1)$ subcomponents. The model tolerates constant factors k_{comp} , k_{comm} , and k_{synch} , but for each depth d these constants are independent of the (p, g, L, m) parameters.

The question is whether a model with such a large number of parameters could be useful. It is shown that Multi-BSP algorithms for problems, such as matrix multiplication, FFT (Fast Fourier Transform), and comparison sorting, can be expressed as parameter-free [486]. A *parameter-free* version of an optimal Multi-BSP algorithm with respect to a given algorithm \mathcal{A} means that it is optimal in:

1. Parallel computation steps to within multiplicative constant factors and in total computation steps to within additive lower order terms.
2. Parallel communication costs to within constant multiplicative factors among Multi-BSP algorithms.
3. Synchronization costs to within constant multiplicative factors among Multi-BSP algorithms.

The proofs of optimality of communication and synchronization in [486] are based on previous work on lower bounds on the number of communication steps of distributed algorithms.

10.6 Modeling concurrency with Petri nets

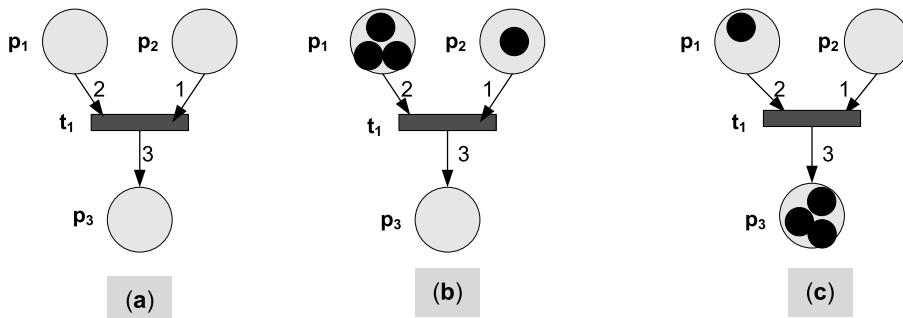
In 1962, Carl Adam Petri introduced a family of graphs for modeling concurrent systems, the Petri nets (PNs) [398]. PNs are used to model the dynamic, rather than the static system behavior, e.g., detect synchronization anomalies. PN models have been extended to study the performance of concurrent systems. More than 8 000 articles on the properties of different flavors of PNs and their applications have been published.

PNs are bipartite graphs populated with tokens flowing through the graph. A *bipartite graph* is one with two classes of vertices; arcs always connect a vertex in one class with one or more vertices in the other class. The two classes of PN vertices are *places* and *transitions*, thus the name Place-Transition (P/T) Nets often used for this class of bipartite graphs; arcs connect either one place with one or more transitions or a transition with one or more places.

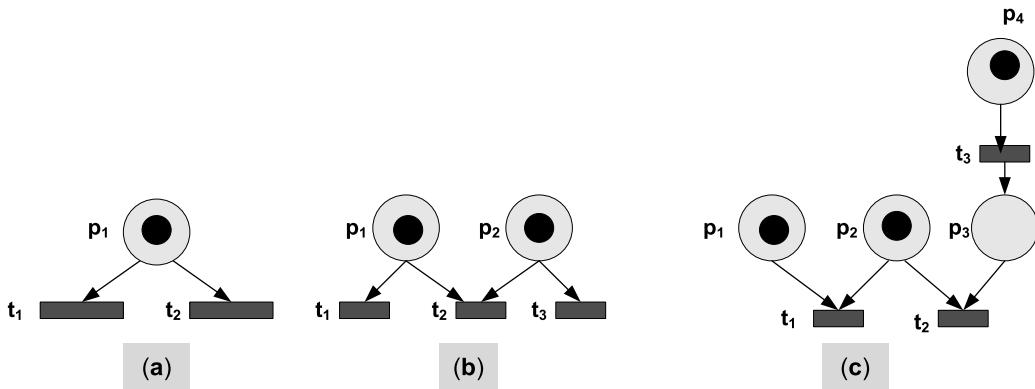
Petri nets model the dynamic behavior of systems. The places of a PN contain tokens; firing of transitions removes tokens from the *input places* of the transition and adds them to its *output places*; see Fig. 10.8. PNs can model different activities in a distributed system. A *transition* may model the occurrence of an event, the execution of a computational task, the transmission of a packet, a logic statement, etc.

The *input places* of a transition model the preconditions of an event, the input data for the computational task, the presence of data in an input buffer, the preconditions of a logic statement. The *output places* of a transition model the postconditions associated with an event, the results of the computational task, the presence of data in an output buffer, or the conclusions of a logic statement.

The distribution of tokens in the places of a PN at a given time is called the *marking* of the net and reflects the state of the system being modeled. PNs are very powerful abstractions and can express both concurrency and choice as we can see in Fig. 10.9.

**FIGURE 10.8**

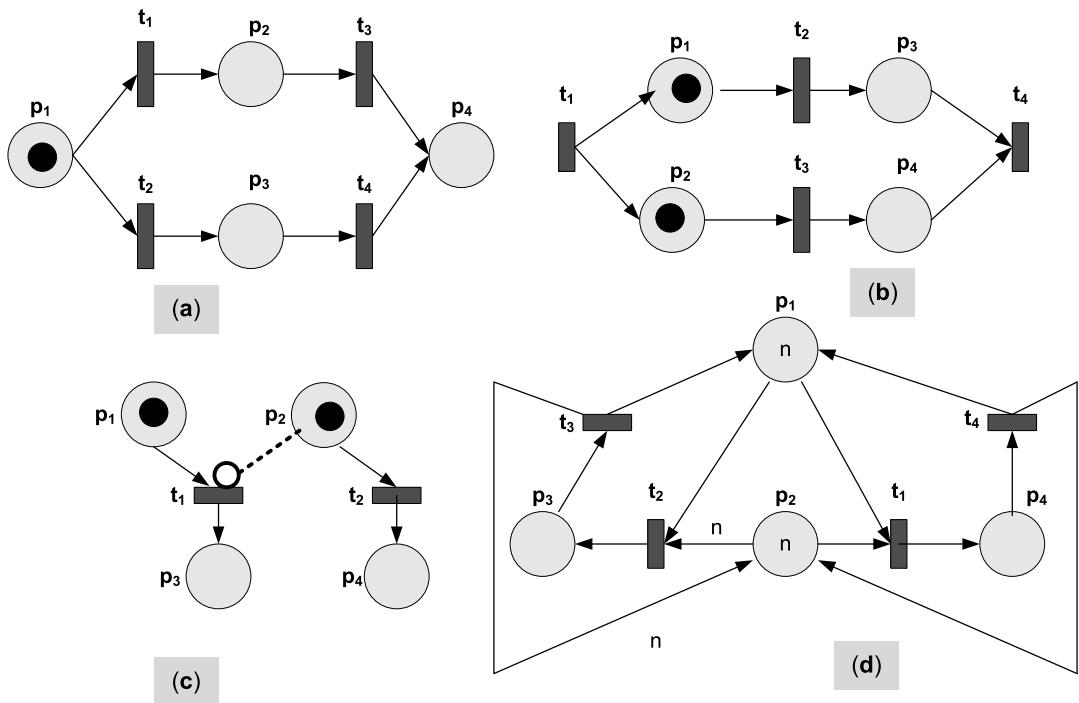
Petri nets firing rules. (a) An unmarked net with one transition t_1 with two input places, p_1 and p_2 , and one output place, p_3 . (b) The marked net, the net with places populated by tokens; the net before firing the enabled transition t_1 . (c) The marked net after firing transition t_1 , two tokens from place p_1 and one from place p_2 are removed and transported to place p_3 .

**FIGURE 10.9**

Petri nets modeling. (a) Choice; only one of transitions t_1 or t_2 may fire. (b) Symmetric confusion; transitions t_1 and t_3 are concurrent, and, at the same time, they are in conflict with t_2 . If t_2 fires, then t_1 and/or t_3 is disabled. (c) Asymmetric confusion; transition t_1 is concurrent with t_3 , and it is in conflict with t_2 if t_3 fires before t_1 .

PNs can model concurrent activities. For example, the net in Fig. 10.9(a) models conflict or choice; only one of the transitions t_1 and t_2 may fire, but not both. Two transitions are said to be *concurrent* if they are causally independent. Concurrent transitions may fire before, after, or in parallel with each other; examples of concurrent transitions are t_1 and t_3 in Figs. 10.9(b) and (c).

When choice and concurrency are mixed, we end up with a situation called *confusion*. *Symmetric confusion* means that two or more transitions are concurrent and, at the same time, they are in conflict with another one. For example, transitions t_1 and t_3 in Fig. 10.9(b) are concurrent and, at the same time, they are in conflict with t_2 . If t_2 fires, either one or both of them will be disabled. *Asymmetric confusion*

**FIGURE 10.10**

(a) A state machine; there is the choice of firing t_1 , or t_2 ; only one transition fires at any given time; concurrency is not possible. (b) A *marked graph* can model concurrency but not choice; transitions t_2 and t_3 are concurrent; there is no causal relationship between them. (c) An extended net used to model priorities; the arc from p_2 to t_1 is an inhibitor arc. The process modeled by transition t_1 is activated only after the process modeled by transition t_2 is activated. (d) Modeling exclusion; transitions t_1 and t_2 model writing and, respectively, reading with n processes to a shared memory. At any given time only one process may write, but any subset of the n processes may read at the same time, provided that no process writes.

occurs when a transition t_1 is concurrent with another transition t_3 and will be in conflict with t_2 if t_3 fires before t_1 as shown in Fig. 10.9(c).

Concurrent transitions t_2 and t_3 in Fig. 10.10(a) model concurrent execution of two processes. A *marked graph* can model concurrency, but not choice; transitions t_2 and t_3 in Fig. 10.10(b) are concurrent; there is no causal relationship between them. Transition t_4 and its input places p_3 and p_4 in Fig. 10.10(b) model synchronization; t_4 fire only if conditions associated with p_3 and p_4 are satisfied.

PNs can be used to model *priorities*. The net in Fig. 10.10(c) models a system with two processes modeled by transitions t_1 and t_2 ; the process modeled by t_2 has a higher priority than the one modeled by t_1 . If both processes are ready to run, places p_1 and p_2 hold tokens. When the two processes are ready, transition t_2 will fire first, modeling the activation of the second process. Only after t_2 is activated, transition t_1 , modeling the activation of the first process, will fire.

PNs are able to model *exclusion*, for example, the net in Fig. 10.10(d) models a group of n concurrent processes in a shared-memory environment. At any given time only one process may write, but any subset of the n processes may read at the same time, provided that no process writes. Place p_3 models the process allowed to write, p_4 the ones allowed to read, p_2 the ones ready to access the shared memory, and p_1 the running tasks. Transition t_2 models the initialization/selection of the process allowed to write and t_1 of the processes allowed to read, whereas t_3 models the completion of a write and t_4 the completion of a read. Indeed, p_3 may have at most one token, while p_4 may have at most n . If all n processes are ready to access the shared memory, all n tokens in p_2 are consumed when transition t_1 fires. However, place p_4 may contain n tokens obtained by successive firings of transition t_2 .

Formal definitions. After this informal discussion of PNs, we switch to a more formal presentation and give several definitions.

Labeled Petri Net: a tuple $N = (p, t, f, l)$ such that:

- $p \subseteq U$ is a finite set of *places*,
- $t \subseteq U$ is a finite set of *transitions*,
- $f \subseteq (p \times t) \cup (t \times p)$ a set of directed arcs, called *flow relations*,
- $l : t \rightarrow L$ a labeling or a weight function

with U a universe of identifiers and L a set of labels. The weight function describes the number of tokens necessary to enable a transition. Labeled PNs describe a static structure; places may contain *tokens* and the distribution of tokens over places defines the state, or the markings of the PN. The dynamic behavior of a PN is described by the structure, together with the markings of the net.

Marked Petri Net: a pair (N, s) where $N = (p, t, f, l)$ is a labeled PN and s is a bag⁴ over p denoting the markings of the net.

Preset and Postset of Transitions and Places. The preset of transition t_i , denoted as $\bullet t_i$, is the set of input places of t_i and the postset, denoted by $t_i \bullet$, is the set of the output places of t_i . The preset of place p_j , denoted as $\bullet p_j$, is the set of input transitions of p_j and the postset, denoted by $p_j \bullet$, is the set of the output transitions of p_j .

Fig. 10.8(a) shows a PN with three places, p_1 , p_2 , and p_3 , and one transition, t_1 . The weights of the arcs from p_1 and p_2 to t_1 are two and one, respectively; the weight of the arc from t_1 to p_3 is three. The preset of transition t_1 in Fig. 10.8(a) consists of two places, $\bullet t_1 = \{p_1, p_2\}$ and its postset consist of only one place, $t_1 \bullet = \{p_3\}$. The preset of place p_4 in Fig. 10.10(a) consists of transitions t_3 and t_4 , $\bullet p_4 = \{t_3, t_4\}$ and the postset of p_1 is $p_1 \bullet = \{t_1, t_2\}$.

Ordinary Net. A PN is ordinary if the weights of all arcs are 1. The nets in Figs. 10.10(a), (b), and (c) are ordinary nets; the weights of all arcs are 1.

Enabled Transition: a transition $t_i \in t$ of the ordinary PN (N, s) , with s the initial marking of N , is enabled if and only if each of its input places contain a token, $(N, s)[t_i] > \Leftrightarrow \bullet t_i \in s$. The notation

⁴ A bag $B(\mathcal{A})$ is a multiset of symbols from an alphabet, \mathcal{A} ; it is a function from \mathcal{A} to the set of natural numbers. For example, $[x^3, y^4, z^5, w^6 \mid P(x, y, z, w)]$ is a bag consisting of three elements x , four elements y , five elements z , and six elements w such that the $P(x, y, z, w)$ holds. P is a predicate on symbols from the alphabet. x is an element of a bag A denoted as $x \in A$ if $x \in \mathcal{A}$ and if $A(x) > 0$.

$(N, s)[t_i >$ means that t_i is enabled. The marking of a PN changes as a result of transition firing; a transition must be enabled to fire.

Firing Rule: the firing of the transition t_i of the ordinary net (N, s) means that a token is removed from each of its input places, and one token is added to each of its output places, its marking changes $s \mapsto (s - \bullet t_i + t_i \bullet)$. Thus, firing of transition t_i changes a marked net (N, s) into another marked net $(N, s - \bullet t_i + t_i \bullet)$.

Firing Sequence: a nonempty sequence of transitions $\sigma \in t^*$ of the marked net (N, s_0) with $N = (p, t, f, l)$ is called a *firing sequence* if and only if there exist markings $s_1, s_2, \dots, s_n \in \mathcal{B}(p)$ and transitions $t_1, t_2, \dots, t_n \in t$ such that $\sigma = t_1, t_2, \dots, t_n$ and for $i \in (0, n)$, $(N, s_i)[t_{i+1}] >$ and $s_{i+1} = s_i - \bullet t_i + t_i \bullet$. All firing sequences that can be initiated from marking s_0 are denoted as $\sigma(s_0)$.

Reachability is the problem of finding if marking s_n is reachable from the initial marking s_0 , with $s_n \in \sigma(s_0)$. Reachability is a fundamental concern for dynamic systems. The reachability problem is decidable; reachability algorithms require exponential time and space.

Liveness: a marked Petri net (N, s_0) is said to be *live* if it is possible to fire any transition infinitely often, starting from the initial marking, s_0 . The absence of deadlock in a system is guaranteed by the liveness of its net model.

Incidence Matrix: given a Petri net with n transitions and m places, the incidence matrix $F = [f_{i,j}]$ is an integer matrix with $f_{i,j} = w(i, j) - w(j, k)$. Here, $w(i, j)$ is the weight of the flow relation (arc) from transition t_i to its output place p_j , and $w(j, k)$ is the weight of the arc from the input place p_j to transition t_k . In this expression, $w(i, j)$ represents the number of tokens added to the output place p_j and $w(j, k)$ the number of tokens removed from the input place p_j when transition t_i fires. F^T is the transpose of the incidence matrix.

A marking s_k can be written as a $m \times 1$ column vector, and its j -th entry denotes the number of tokens in place j after some transition firing. The necessary and sufficient condition for transition t_k to be enabled at a marking s is that $w(j, k) \leq s(j) \quad \forall s_j \in \bullet t_i$, the weight of the arc from every input place of the transition, be smaller or equal to the number of tokens in the corresponding input place.

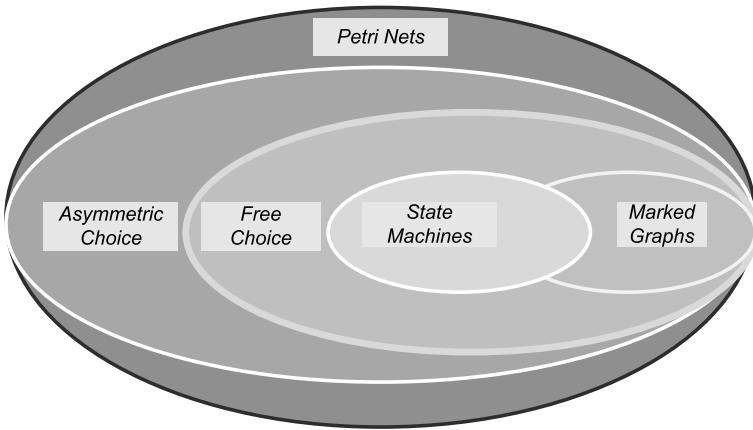
Extended Nets: PNs with inhibitor arcs; an *inhibitor arc* prevents the enabling of a transition. For example, the arc from p_2 to t_1 in the net in Fig. 10.10(a) is an inhibitor arc; the process modeled by transition t_1 can be activated only after the process modeled by transition t_2 is activated.

Modified Transition Enabling Rule for Extended Nets: a transition is not enabled if one of the places in its preset is connected with the transition with an inhibitor arc and if the place holds a token. For example, transition t_1 in the net in Fig. 10.10(c) is not enabled while place p_2 holds a token.

There are several classes of PNs distinguished from one another by their structural properties:

1. State Machines—are used to model finite state machines and cannot model concurrency and synchronization.
2. Marked Graphs—cannot model choice and conflict.
3. Free-choice Nets—cannot model confusion.
4. Extended Free-choice Nets—cannot model confusion but allow inhibitor arcs.
5. Asymmetric Choice Nets—can model asymmetric confusion but not symmetric ones.

This partitioning is based on the number of input and output flow relations from/to a transition or a place and by the manner in which transitions share input places, as indicated in Fig. 10.11.

**FIGURE 10.11**

Classes of Petri nets.

State Machine: a PN is a state machine if and only if

$$|\bullet t_i| = 1 \wedge |t_i \bullet| = 1, \forall t_i \in t. \quad (10.2)$$

All transitions of a state machine have exactly one incoming and one outgoing arc. This topological constraint limits the expressiveness of a state machine; no concurrency is possible. For example, the transitions t_1 , t_2 , t_3 , and t_4 of state machine in Fig. 10.10(a) have only one input and one output arc; the cardinality of their presets and postsets is one. No concurrency is possible; once a choice was made by firing either t_1 , or t_2 , the evolution of the system is entirely determined. This state machine has four places p_1 , p_2 , p_3 , and p_4 , and the marking is a 4-tuple (p_1, p_2, p_3, p_4) ; the possible markings of this net are $(1, 0, 0, 0)$, $(0, 1, 0, 0)$, $(0, 0, 1, 0)$, $(0, 0, 0, 1)$, with a token in places p_1 , p_2 , p_3 , or p_4 , respectively.

Marked Graph: a PN is a marked graph if and only if

$$|\bullet p_i| = 1 \wedge |p_i \bullet| = 1, \forall p_i \in p. \quad (10.3)$$

In a marked graph, each place has only one incoming and one outgoing flow relation; thus, marked graphs do not allow modeling of choice.

Free Choice, Extended Free Choice, and Asymmetric Choice Petri Nets: the marked net, (N, s_0) with $N = (p, t, f, l)$ is a free-choice net if and only if

$$(\bullet t_i) \cap (\bullet t_j) = \emptyset \Rightarrow |\bullet t_i| = |\bullet t_j|, \forall t_i, t_j \in t. \quad (10.4)$$

N is an extended free-choice net if $(\bullet t_i) \cap (\bullet t_j) = \emptyset \Rightarrow \bullet t_i = \bullet t_j, \forall t_i, t_j \in t$.

N is an asymmetric choice net if and only if $(\bullet t_i) \cap (\bullet t_j) \neq \emptyset \Rightarrow (\bullet t_i \subseteq \bullet t_j) \text{ or } (\bullet t_i \supseteq \bullet t_j), \forall t_i, t_j \in t$.

In an extended free-choice net, if two transitions share an input place, they must share all places in their presets. In an asymmetric choice net two transitions may share only a subset of their input places.

Several extensions of PNs have been proposed. For example, Colored Petri Nets (CPNs) allow tokens of different colors, thus increasing the expressivity of the PNs but not simplifying their analysis. Several extensions of PNs to support performance analysis by associating a random time with each transition have been proposed. In case of Stochastic Petri Nets (SPNs), a random time elapses between the time a transition is enabled and the moment it fires. This random time allows the model to capture the service time associated with the activity modeled by the transition.

Applications of SPNs to performance analysis of complex systems is generally limited by the explosion of the state space of the models. Stochastic High-Level Petri Nets (SHLPN) were introduced in 1988 [309]; the SHLPNs allow easy identification of classes of equivalent markings even when the corresponding aggregation of states in the Markov domain is not obvious. This aggregation could reduce the size of the state space by one or more orders of magnitude depending on the system being modeled.

Cloud applications often require a large number of tasks to run concurrently. The interdependencies among these tasks are quite intricate, and PNs can be very useful to intuitively present a high-level model of the interactions among the tasks. The five classes of systems mentioned earlier, including finite-state machines, as well as control flow and data flow systems, can be modeled by PNs. The workflow patterns common to many cloud applications discussed in Section 11.2 translate easily to PN models. Unfortunately, the size of the state space of detailed PN models of complex systems grows very fast, and this limits the usefulness of these models.

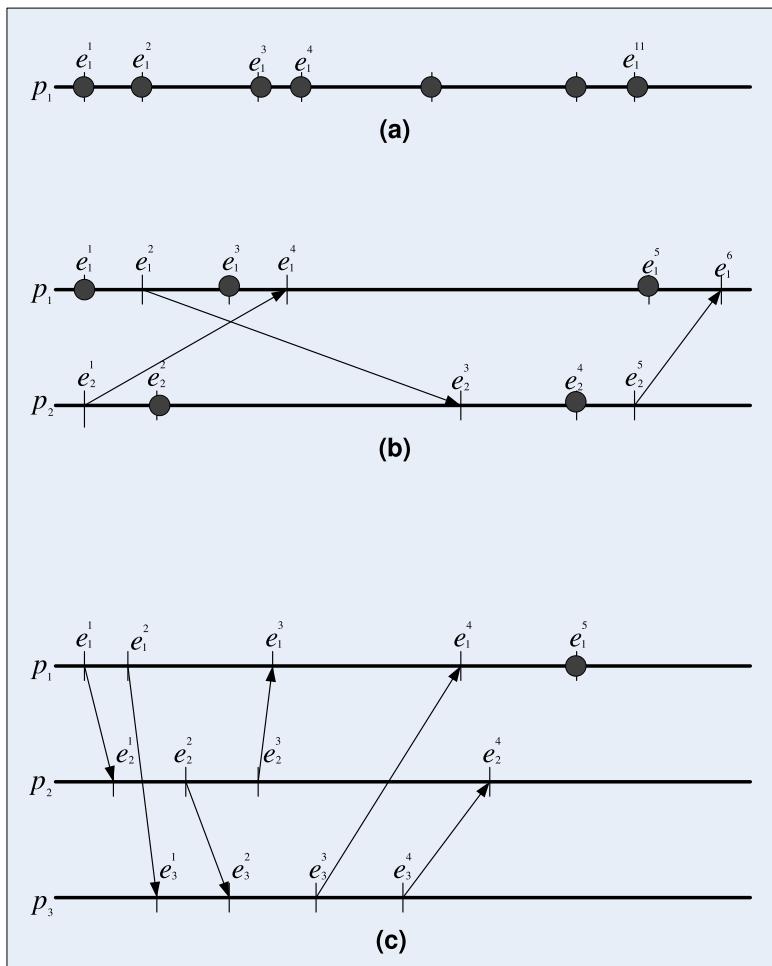
10.7 Process state; global state of a process or thread group

To understand the important properties of distributed systems, we use a model, an abstraction based on two critical components, processes/threads and communication channels. A *process* is a program in execution, and a *thread* is a light-weight process. A thread of execution is the smallest unit of processing that can be scheduled by an operating system.

A process or a thread is characterized by its *state*. The state is the ensemble of information that we need to restart a process or thread after it was suspended. An *event* is a change of state of a process or a thread. The events affecting the state of process p_i are numbered sequentially as $e_i^1, e_i^2, e_i^3, \dots$, as shown in the space-time diagram in Fig. 10.12(a). A process p_i is in state σ_i^j immediately after the occurrence of event e_i^j and remains in that state until the occurrence of the next event, e_i^{j+1} .

A *process or a thread group* is a collection of cooperating processes and threads; to reach a common goal, the processes work in concert and communicate with one another. For example, a parallel algorithm to solve a system of partial differential equations (PDEs) over a domain D may partition the data in several segments and assign each segment to one of the members of the group. The processes or the threads in the group must cooperate with one another and iterate until the common boundary values computed by one process agree with the common boundary values computed by another.

A *communication channel* provides the means for processes/threads to communicate with one another and coordinate their actions by exchanging messages. Without loss of generality, we assume that

**FIGURE 10.12**

Space-time diagrams display local and communication events during a process lifetime. Local events are small black circles. Communication events in different processes/threads are connected by lines originating at a *send* event and terminated by an arrow at the *receive* event. (a) All events in the case of a single process p_1 are local; the process is in state σ_1 immediately after the occurrence of event e_1^1 and remains in that state until the occurrence of event e_1^2 . (b) Two processes/threads p_1 and p_2 ; event e_1^2 is a communication event; p_1 sends a message to p_2 ; event e_2^3 is a communication event, process or thread p_2 receives the message sent by p_1 . (c) Three processes or threads interact by means of communication events.

communication among processes is done only by means of $send(m)$ and $receive(m)$ communication events, where m is a message. We use the term “message” for a structured unit of information that can be interpreted only in a semantic context by the sender and the receiver. The *state of a communication*

channel is defined as follows: given two processes p_i and p_j , the state of the channel, $\xi_{i,j}$, from p_i to p_j consists of messages sent by p_i but not yet received by p_j .

These two abstractions allow us to concentrate on critical properties of distributed systems without the need to discuss the detailed physical properties of the entities involved. The model presented is based on the assumption that a channel is a unidirectional bit pipe of infinite bandwidth and zero latency, but unreliable; messages sent through a channel may be lost or distorted or the channel may fail, losing its ability to deliver messages. We also assume that the time a process needs to traverse a set of states is of no concern and that processes may fail, or be aborted.

A *protocol* is a finite set of messages exchanged among processes and threads to help them coordinate their actions. Fig. 10.12(c) illustrates the case when communication events are dominant in the local history of processes, p_1 , p_2 , and p_3 . In this case, only e_1^5 is a local event; all others are communication events. The particular protocol illustrated in Fig. 10.12(c) requires processes p_2 and p_3 to send messages to the other processes in response to a message from process p_1 .

Informal definition of the state of a single process or a thread can be extended to collections of communicating processes/threads. The *global state of a distributed system* consisting of several processes and communication channels is the union of the states of individual processes and channels [40].

Call h_i^j the history of process p_i up to and including its j -th event, e_i^j , and call $\sigma_i^{j_i}$ the local state of process p_i following event e_i^j . Consider a system consisting of n processes, $p_1, p_2, \dots, p_i, \dots, p_n$ with $\sigma_i^{j_i}$ the local state of process p_i ; then, the global state of the system is an n -tuple of local states

$$\Sigma^{(j_1, j_2, \dots, j_n)} = (\sigma_1^{j_1}, \sigma_2^{j_2}, \dots, \sigma_i^{j_i}, \dots, \sigma_n^{j_n}). \quad (10.5)$$

The state of the channels does not appear explicitly in this definition of the global state because the state of the channels is encoded as part of the local state of the processes/threads communicating through the channels. The global states of a distributed computation with n processes form an n -dimensional lattice. The elements of this lattice are global states $\Sigma^{(j_1, j_2, \dots, j_n)}(\sigma_1^{j_1}, \sigma_2^{j_2}, \dots, \sigma_n^{j_n})$.

Fig. 10.13(a) shows the lattice of global states of the distributed computation in Fig. 10.12(b). This is a two-dimensional lattice because we have two processes, p_1 and p_2 . The lattice of global states for the distributed computation in Fig. 10.12(c) is a three-dimensional lattice; the computation consists of three concurrent processes, p_1 , p_2 , and p_3 .

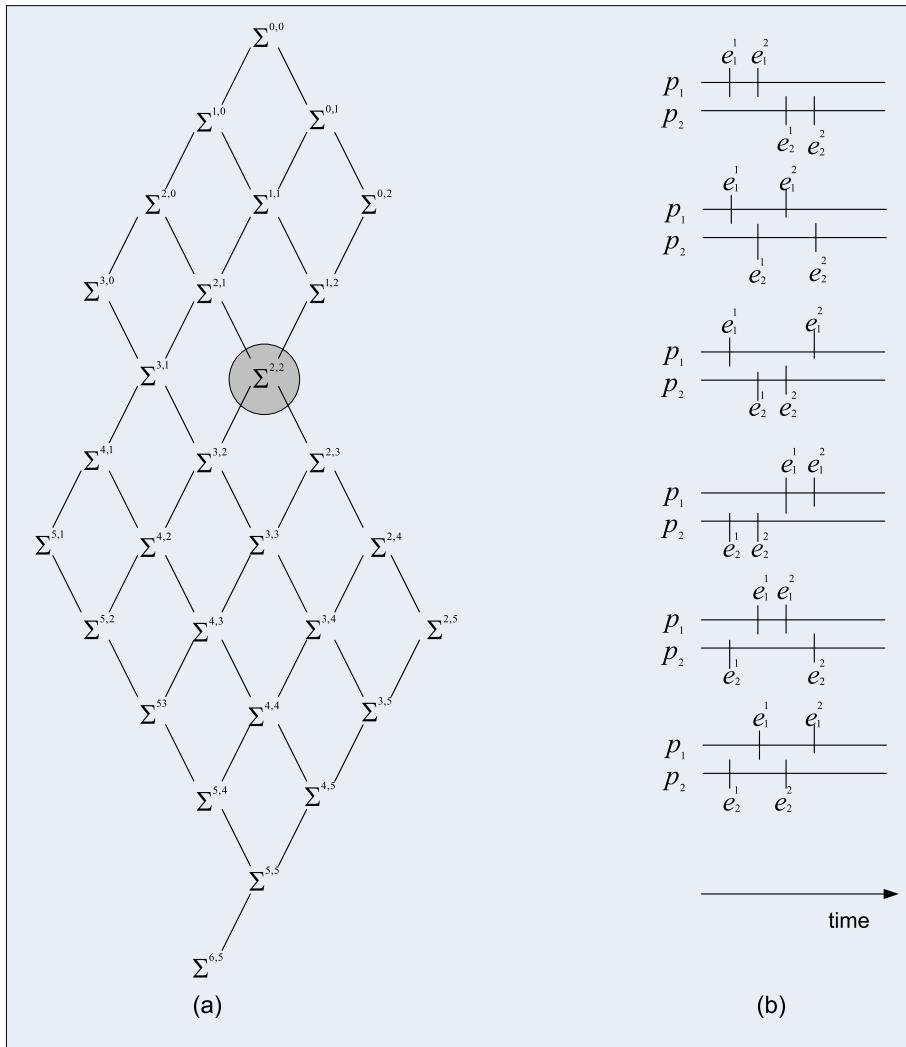
The initial state of the system in Fig. 10.13(b) is the state before the occurrence of any event, and it is denoted by $\Sigma^{(0,0)}$; the only global states reachable from $\Sigma^{(0,0)}$ are $\Sigma^{(1,0)}$, and $\Sigma^{(0,1)}$. The communication events limit the global states the system may reach; in this example the system cannot reach the state $\Sigma^{(4,0)}$ because process p_1 enters state σ_4 only after process p_2 has entered the state σ_1 . Fig. 10.13(b) shows the six possible sequences of events to reach the global state $\Sigma^{(2,2)}$:

$$(e_1^1, e_2^2, e_1^1, e_2^2), (e_1^1, e_2^1, e_1^2, e_2^2), (e_1^1, e_2^1, e_2^2, e_1^2), (e_2^1, e_1^2, e_1^1, e_2^2), (e_2^1, e_1^2, e_2^1, e_1^2), (e_2^1, e_1^1, e_2^1, e_2^2). \quad (10.6)$$

An interesting question is: How many paths needed to reach a global state exist? The more paths exist, the harder it is to identify the events leading to a state when we observe an undesirable behavior of the system. A large number of paths increases the difficulties to debug the system.

We conjecture that, in the case of two threads in Fig. 10.13(a), the number of paths from the global state $\Sigma^{(0,0)}$ to $\Sigma^{(m,n)}$ is

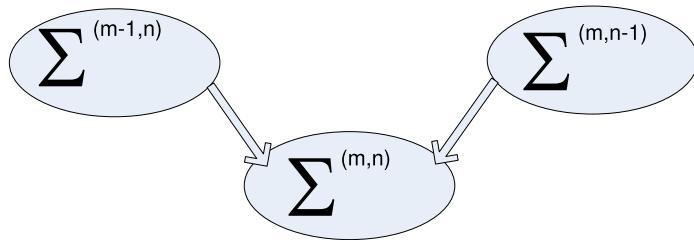
$$N_p^{(m,n)} = \frac{(m+n)!}{m!n!}. \quad (10.7)$$

**FIGURE 10.13**

(a) The lattice of the global states of two processes/threads with the space–time diagrams in Fig. 10.12(b). (b) The six possible sequences of events leading to the state $\Sigma^{(2,2)}$.

We have already seen that there are six paths leading to state $\Sigma^{(2,2)}$ and, indeed

$$N_p^{(2,2)} = \frac{(2+2)!}{2!2!} = \frac{24}{4} = 6. \quad (10.8)$$

**FIGURE 10.14**

In the two-dimensional case the global state $\Sigma^{(m,n)}$, $\forall(m, n) \geq 1$ can only be reached from two states, $\Sigma^{(m-1,n)}$ and $\Sigma^{(m,n-1)}$.

To prove Eq. (10.7), we use a method resembling induction; we notice first that the global state $\Sigma^{(1,1)}$ can only be reached from the states $\Sigma^{(1,0)}$ and $\Sigma^{(0,1)}$ and that $N_p^{(1,1)} = (2)!/1!1! = 2$; thus the formula is true for $m = n = 1$. Then, we show that, if the formula is true for the $(m - 1, n - 1)$ case, it will also be true for the (m, n) case. If our conjecture is true, then

$$N_p^{[(m-1),n]} = \frac{[(m-1)+n]!}{(m-1)!n!} \quad (10.9)$$

and

$$N_p^{[m,(n-1)]} = \frac{[(m+(n-1)]!}{m!(n-1)!}. \quad (10.10)$$

We observe that the global state $\Sigma^{(m,n)}$, $\forall(m, n) \geq 1$ can only be reached from two states, $\Sigma^{(m-1,n)}$ and $\Sigma^{(m,n-1)}$, as seen in Fig. 10.14, thus

$$N_p^{(m,n)} = N_p^{(m-1,n)} + N_p^{(m,n-1)}. \quad (10.11)$$

Indeed

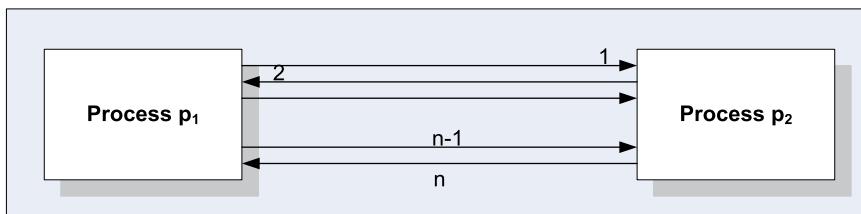
$$\frac{[(m-1)+n]!}{(m-1)!n!} + \frac{[m+(n-1)]!}{m!(n-1)!} = (m+n-1)! \left[\frac{1}{(m-1)!n!} + \frac{1}{m!(n-1)!} \right] = \frac{(m+n)!}{m!n!}. \quad (10.12)$$

This shows that our conjecture is true, and, thus, Eq. (10.7) gives the number of paths to reach the global state $\Sigma^{(m,n)}$ from $\Sigma^{(0,0)}$ when two threads are involved. This expression can be generalized for the case of q threads; using the same strategy, we see that the number of paths from the state $\Sigma^{(0,0,\dots,0)}$ to the global state $\Sigma^{(n_1,n_2,\dots,n_q)}$ is

$$N_p^{(n_1,n_2,\dots,n_q)} = \frac{(n_1+n_2+\dots+n_q)!}{n_1!n_2!\dots n_q!}. \quad (10.13)$$

Indeed,

$$N_p^{(n_1,n_2,\dots,n_q)} = N_p^{(n_1-1,n_2,\dots,n_q)} + N_p^{(n_1,n_2-1,\dots,n_q)} + \dots + N_p^{(n_1,n_2,\dots,n_q-1)}. \quad (10.14)$$

**FIGURE 10.15**

Process coordination in the presence of errors; each message may be lost with probability ϵ . If a protocol consisting of n messages exists, then the protocol should be able to function properly with $n - 1$ messages reaching their destination, one of them being lost.

Eq. (10.13) gives us an indication of how difficult it is to debug a system with a large number of concurrent threads.

Many problems in distributed systems are instances of the *global predicate evaluation problem* (GPE) where the goal is to evaluate a Boolean expression whose elements are functions of the global state of the system.

10.8 Communication protocols and process coordination

A major concern in any parallel and distributed system is communication in the presence of channel failures. There are multiple modes for a channel to fail, and some lead to messages being lost. In the general case, it is impossible to guarantee that two processes will reach an agreement in case of channel failures; see Fig. 10.15.

Statement. Given two processes p_1 and p_2 connected by a communication channel that can lose a message with probability $\epsilon > 0$, no protocol capable of guaranteeing that two processes will reach agreement exists, regardless of how small the probability ϵ is.

The proof of this statement is by contradiction. Assume that such a protocol exists and it consists of n messages. Recall that a protocol consists of a finite number of messages. Since any message might be lost with probability ϵ , the protocol should be able to function when only $n - 1$ messages reach their destination, the last one being lost. Induction on the number of messages proves that indeed no such protocol exists; indeed, the same reasoning leads us to conclude that the protocol should function correctly with $n - 2$ messages, and so on.

In practice, error detection and error correction codes allow processes to communicate reliably through noisy digital channels. The redundancy of a message is increased by more bits and packaging the message as a codeword; the recipient of the message is then able to decide if the sequence of bits received is a valid codeword, and, if the code satisfies some distance properties, then the recipient of the message is able to extract the original message from a bit string in error.

Communication protocols implement not only *error control* mechanisms, but also flow control and congestion control. *Flow control* provides feedback from the receiver; it forces the sender to transmit

only the amount of data the receiver is able to buffer and then process. *Congestion control* ensures that the offered load of the network does not exceed the network capacity. In store-and-forward networks, individual routers may drop packets when the network is congested and the sender is forced to retransmit. Based on the estimation of the RTT (Round-Trip-Time), the sender can detect congestion and reduce the transmission rate.

The implementation of these mechanisms requires the measurement of *time intervals*, the time elapsed between two events; we also need a *global concept of time* shared by all entities that cooperate with one another. For example, a computer chip has an *internal clock* and a predefined set of actions occurs at each clock tick. Each chip has an *interval timer* that helps enhance the system's fault tolerance; when the effects of an action are not sensed after a predefined interval, the action is repeated.

When the entities communicating with each other are networked computers, clock synchronization precision is critical [291]. The event rates are very high, each system goes through state changes at a very fast pace; modern processors run at a 2–4 GHz clock rate. This explains why we need to measure time very accurately; indeed, we have atomic clocks with an accuracy of about 10^{-6} seconds per year.

An isolated system can be characterized by its *history* expressed as a sequence of events, each event corresponding to a change of the state of the system. Local timers provide relative time measurements. A more accurate description adds to the system's history the time of occurrence of each event as measured by the local timer.

Determining the global state of a large-scale distributed system is a very challenging problem. Messages sent by processes may be lost or distorted during transmission. There are no means to ensure a perfect synchronization of local clocks and no obvious methods to ensure a global ordering of events occurring in different processes unless we restrict message delays and the type of errors.

The mechanisms just described are insufficient once we approach the problem of cooperating entities. To coordinate their actions, two entities need a common perception of time. Timers are not enough; clocks provide the only way to measure distributed duration, that is, actions that start in one process and terminate in another.

Global agreement on time is necessary to trigger actions that should occur concurrently. For example, in a real-time control system of a power plant, several circuits must be switched on at the same time. Agreement on the *time when events occur* is necessary for a distributed recording of events, for example, to determine a precedence relationship through a temporal ordering of events. To ensure that a system functions correctly, we need to determine that the event causing a change of state occurred before the state change, e.g., the sensor triggering an alarm has to change its value before the emergency procedure to handle the event was activated. Another example of the need for agreement on the time of occurrence of events is in replicated actions. In this case, several replicas of a process must log the time of an event in a consistent manner.

Time stamps are often used for event ordering using a global time-base constructed on local virtual clocks [338]. The Δ -protocols [119] achieve total temporal order using a global time base. Assume that local virtual clock readings do not differ by more than π , called *precision* of the global time base. Call g the *granularity of physical clocks*. First, observe that the granularity should not be smaller than the precision; given two events a and b occurring in different processes, if $t_b - t_a \leq \pi + g$, we cannot tell which of a or b occurred first [493]. Based on these observations, it follows that the order discrimination of clock-driven protocols cannot be better than twice the clock granularity.

System specification, design, and analysis require a clear understanding of *cause–effect relationships*. During the system specification phase, we view the system as a state machine and define the

actions that cause transitions from one state to another. During the system analysis phase, we need to determine the cause that brought the system to a certain state.

The activity of any process is modeled as a sequence of *events*; hence, the binary relationship cause–effect should be expressed in terms of events and should support our intuition that *the cause must precede the effect*. Again, we need to distinguish between local events and communication events. The latter ones affect more than one process and are essential for constructing a global history of an ensemble of processes. Let h_i denote the local history of process p_i , and let e_i^k denote the k -th event in this history.

The binary cause–effect relationship between two events has the following properties:

1. Causality of local events can be derived from the process history. Given two events e_i^k and e_i^l local to process p_i ,

$$\text{if } e_i^k, e_i^l \in h_i \text{ and } k < l, \text{ then } e_i^k \rightarrow e_i^l. \quad (10.15)$$

2. Causality of communication events. Given two processes p_i and p_j and two events e_i^k and e_j^l ,

$$\text{if } e_i^k = \text{send}(m) \text{ and } e_j^l = \text{receive}(m), \text{ then } e_i^k \rightarrow e_j^l. \quad (10.16)$$

3. Transitivity of the causal relationship. Given three processes, p_i , p_j , and p_m and the events e_i^k , e_j^l , and e_m^n occurring in p_i , p_j , and p_m , respectively,

$$\text{if } e_i^k \rightarrow e_j^l \text{ and } e_j^l \rightarrow e_m^n, \text{ then } e_i^k \rightarrow e_m^n. \quad (10.17)$$

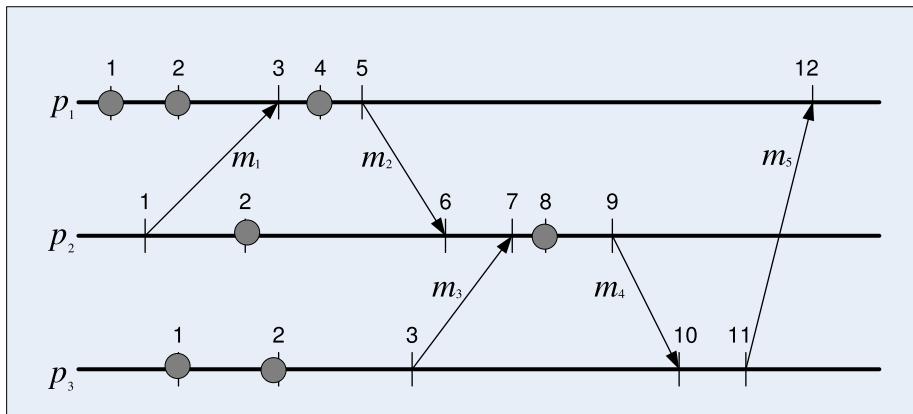
Two events in the global history may be unrelated, neither one being the cause of the other; such events are said to be *concurrent events*.

10.9 Communication, logical clocks, and message delivery rules

We need to bridge the gap between the physical systems and the abstractions used to describe interacting processes. This section addresses the means to bridge this gap. Communicating processes often run on distant systems whose physical clocks cannot be perfectly synchronized due to communication latency. Global ordering of events in communicating processes running on such systems is not feasible, and logical clocks are used instead. Also, messages travel through physical channels with different speeds and follow different paths. As a result, the order in which messages are delivered to processes may be different than the order they were sent.

Logical clocks. A *logical clock (LC)* is an abstraction necessary to ensure the clock condition given by Eqs. (10.24) and (10.25) in the absence of a global clock. Each process p_i maps events to positive integers. Call $LC(e)$ the local variable associated with event e . Each process time stamps each message m sent with the value of the logical clock at the time of sending, $TS(m) = LC(\text{send}(m))$. The rules to update the logical clock are specified by the following relationship:

$$LC(e) := \begin{cases} LC + 1 & \text{if } e \text{ is a local event or a } \text{send}(m) \text{ event} \\ \max(LC, TS(m) + 1) & \text{if } e = \text{receive}(m). \end{cases} \quad (10.18)$$

**FIGURE 10.16**

Three processes and their logical clocks; The usual labeling of events as $e_1^1, e_1^2, e_1^3, \dots$ is omitted to avoid overloading the figure; only the logical clock values for the local and for the communication events are marked. The correspondence between the events and the logical clock values is obvious: $e_1^1, e_2^1, e_3^1 \rightarrow 1, e_1^5 \rightarrow 5, e_2^4 \rightarrow 7, e_3^4 \rightarrow 10, e_1^6 \rightarrow 12$, etc. Global ordering of all events is not possible; there is no way to establish the ordering of events e_1^1, e_2^1 , and e_3^1 .

The concept of logical clocks is illustrated in Fig. 10.16 using a modified *space-time diagram*, where the events are labeled with the logical clock value. Messages exchanged between processes are shown as lines from the sender to the receiver; the communication events corresponding to sending and receiving messages are marked on these diagrams.

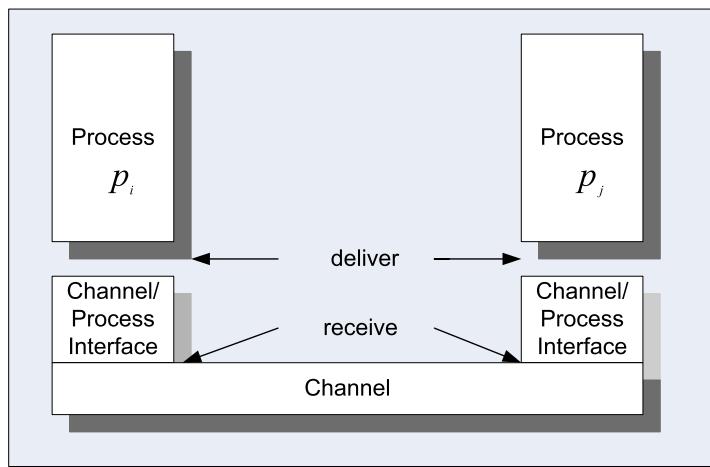
Each process labels local events and send events sequentially until it receives a message marked with a logical clock value larger than the next local logical clock value, as shown in Eq. (10.18). It follows that logical clocks do not allow a global ordering of all events. For example, there is no way to establish the ordering of events e_1^1, e_2^1 , and e_3^1 in Fig. 10.16. Nevertheless, communication events allow different processes to coordinate their logical clocks; for example, process p_2 labels the event e_2^3 as 6 because of message m_2 , which carries the information about the logical clock value as 5 at the time message m_2 was sent. Recall that e_i^j is the j -th event in process p_i .

Logical clocks lack an important property, *gap detection*; given two events e and e' and their logical clock values, $LC(e)$ and $LC(e')$, it is impossible to establish if an event e'' exists such that

$$LC(e) < LC(e'') < LC(e'). \quad (10.19)$$

For example, for process p_1 , there is an event, e_1^4 , between the events e_1^3 and e_1^5 in Fig. 10.16; indeed, $LC(e_1^3) = 3$, $LC(e_1^5) = 5$, $LC(e_1^4) = 4$, and $LC(e_1^3) < LC(e_1^4) < LC(e_1^5)$. However, for process p_3 , the events e_3^3 and e_3^4 are consecutive though $LC(e_3^3) = 3$ and $LC(e_3^4) = 10$.

Message delivery rules. The communication channel abstraction makes no assumptions about the order of messages; a real-life network might reorder messages. This fact has profound implications

**FIGURE 10.17**

Message receiving and message delivery are two distinct operations. The channel-process interface implements the delivery rules, e.g., FIFO delivery.

for a distributed application. Consider, for example, a robot getting instructions to navigate from a monitoring facility with two messages, “turn left” and “turn right”, being delivered out of order.

Message receiving and message delivery are two distinct operations; a *delivery rule* is an additional assumption about the channel-process interface. This rule establishes when a message received is actually delivered to the destination process. The receiving of a message m and its delivery are two distinct events in a causal relationship with one another: A message can only be delivered after being received; see Fig. 10.17:

$$\text{receive}(m) \rightarrow \text{deliver}(m). \quad (10.20)$$

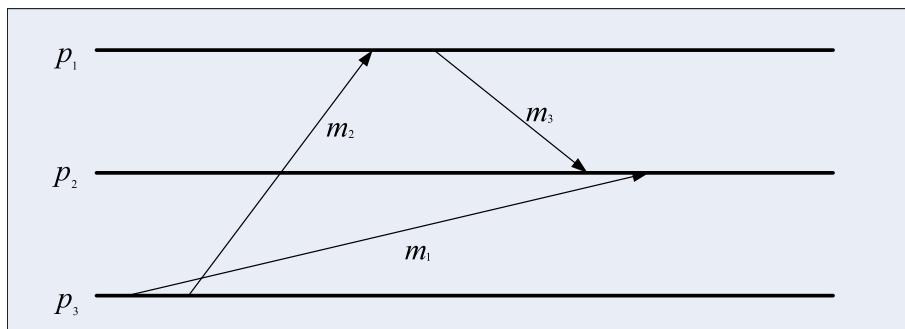
First-In-First-Out (FIFO) delivery implies that messages are delivered in the same order they are sent. For each pair of source-destination processes (p_i, p_j), FIFO delivery requires that the following relationship should be satisfied:

$$\text{send}_i(m) \rightarrow \text{send}_i(m') \Rightarrow \text{deliver}_j(m) \rightarrow \text{deliver}_j(m'). \quad (10.21)$$

Even if the communication channel does not guarantee FIFO delivery, FIFO delivery can be enforced by attaching a sequence number to each message sent. The sequence numbers are also used to reassemble messages out of individual packets.

Causal message delivery. Causal delivery is an extension of the FIFO delivery to the case when a process receives messages from different sources. Assume a group of three processes (p_i, p_j, p_k) and two messages m and m' . Causal delivery requires that

$$\text{send}_i(m) \rightarrow \text{send}_j(m') \Rightarrow \text{deliver}_k(m) \rightarrow \text{deliver}_k(m'). \quad (10.22)$$

**FIGURE 10.18**

Violation of causal delivery when more than two processes are involved; message m_1 is delivered to process p_2 after message m_3 , though message m_1 was sent before m_3 . Indeed, message m_3 was sent by process p_1 after receiving m_2 , which in turn was sent by process p_3 after sending message m_1 .

When more than two processes are involved in a message exchange, the message delivery may be FIFO, but not causal as shown in Fig. 10.18 where we see that

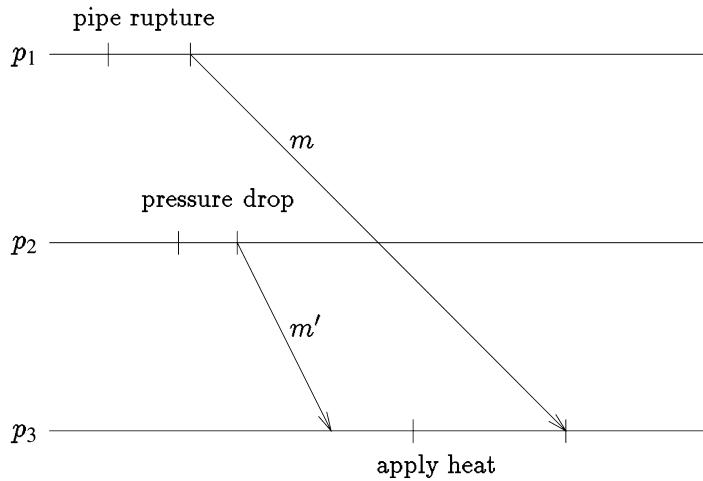
- $\text{deliver}(m_3) \rightarrow \text{deliver}(m_1)$; according to the local history of process p_2 .
- $\text{deliver}(m_2) \rightarrow \text{send}(m_3)$; according to the local history of process p_1 .
- $\text{send}(m_1) \rightarrow \text{send}(m_2)$; according to the local history of process p_3 .
- $\text{send}(m_2) \rightarrow \text{deliver}(m_2)$.
- $\text{send}(m_3) \rightarrow \text{deliver}(m_3)$.

The transitivity property and the causality relationship imply that

$$\text{send}(m_1) \rightarrow \text{deliver}(m_3). \quad (10.23)$$

Call $TS(m)$ the *time stamp* carried by message m . A message received by process p_i is *stable* if no future messages with a time stamp smaller than $TS(m)$ can be received by process p_i . When using logical clocks, a process p_i can construct consistent observations of the system if it implements the following delivery rule: *deliver all stable messages in increasing time stamp order*.

Let us now examine the problem of *consistent message delivery* under several sets of assumptions. First, assume that processes cooperating with each other in a distributed environment have access to a *global real-time clock*, that the message delays are bounded by δ , and that there is no clock drift. Call $RC(e)$ the time of occurrence of event e . A process includes in every message it sends the $RC(e)$, where e is the send message event. The delivery rule in this case is: *at time t deliver all received messages with time stamps up to $(t - \delta)$ in increasing time stamp order*. Indeed, this delivery rule guarantees that, under the bounded delay assumption, the message delivery is consistent. All messages delivered at time t are in order, and no future message with a time stamp lower than any of the messages delivered may arrive.

**FIGURE 10.19**

The causal sequence of events observed by the controller process p_3 , namely, pressure drop \rightarrow apply heat \rightarrow pipe rupture leads to the erroneous conclusion that the cause of the pipe rupture is the increased temperature.

For any two events, e and e' , occurring in different processes, the so-called *clock condition* is satisfied if

$$e \rightarrow e' \Rightarrow RC(e) < RC(e'), \forall e, e'. \quad (10.24)$$

Oftentimes, we are interested in determining the set of events that caused an event knowing the time stamps associated with all events; in other words, we are interested in deducing the causal precedence relation between events from their time stamps. To do so, we need to define the so-called *strong clock condition*. The strong clock condition requires an equivalence between the causal precedence and the ordering of the time stamps

$$\forall e, e', \quad e \rightarrow e' \equiv TS(e) < TS(e'). \quad (10.25)$$

Causal delivery is very important because it allows processes to reason about the entire system using only local information. This is only true in a closed system where all communication channels are known. Sometimes, a system has a *hidden channel* and reasoning based on causal analysis may lead to incorrect conclusions [40]. This is the case of the example in Fig. 10.19 in which the environment acts as a hidden channel.

In this example process p_1 detects the rupture of a steam pipe and sends message m to process p_3 . Process p_2 monitors the pipe pressure, and a few seconds after the rupture of the pipe detects a drop in the pressure and sends message m' to p_3 . Messages m and m' are concurrent from the point of view of explicit communication. p_3 , the controller process, reacts to the pressure drop by applying more heat to increase the temperature. Message m reporting the rupture of the pipe arrives moments later. The causal sequence of events observed by the controller process, pressure drop, apply heat, pipe rupture, leads to the erroneous conclusion that the pipe rupture is due to increased temperature [281]!

10.10 Runs and cuts; causal history

We often need to know the state of several, possibly all, processes of a distributed system. For example, a supervisory process must be able to detect when a subset of processes is deadlocked; a process might migrate from one location to another or be replicated only after an agreement with others. In all these examples, a process needs to evaluate a predicate function of the global state of the system.

We call the process responsible for constructing the global state of the system the *monitor*; a monitor sends messages requesting information about the local state of every process and gathers the replies to construct the global state. Intuitively, the construction of the global state is equivalent to taking snapshots of individual processes and then combining these snapshots into a global view. Yet, combining snapshots is straightforward if and only if all processes have access to a global clock and the snapshots are taken at the same time; hence, the snapshots are contemporaneous with one another.

A *run* is a total ordering R of all the events in the global history of a distributed computation consistent with the local history of each participant process; a run

$$R = (e_1^{j_1}, e_2^{j_2}, \dots, e_n^{j_n}) \quad (10.26)$$

implies a sequence of events, as well as a sequence of global states.

For example, consider the three processes in Fig. 10.20. We can construct a three-dimensional lattice of global states following a procedure similar to the one in Fig. 10.13, starting from the initial state $\Sigma^{(000)}$ and proceeding to any reachable state $\Sigma^{(ijk)}$ with i, j, k being the events in processes p_1, p_2, p_3 , respectively. The run $R_1 = (e_1^1, e_2^1, e_3^1, e_1^2)$ is consistent with both the local history of each process and the global history; this run is valid, and the system has traversed the global states

$$\Sigma^{000}, \Sigma^{100}, \Sigma^{110}, \Sigma^{111}, \Sigma^{211}. \quad (10.27)$$

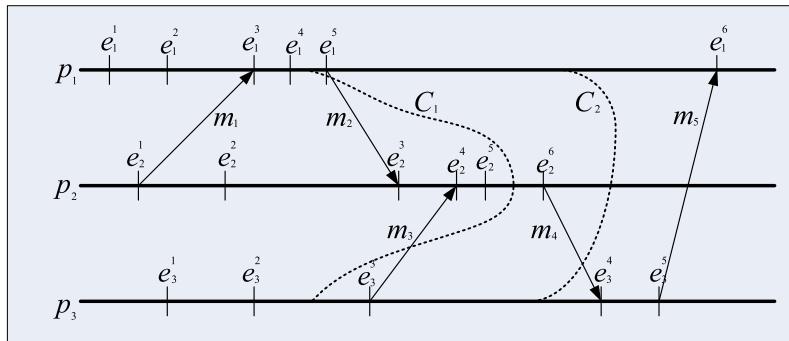
On the other hand, the run $R_2 = (e_1^1, e_1^2, e_3^1, e_3^2)$ is invalid because it is inconsistent with the global history. The system cannot ever reach the state Σ^{301} ; message m_1 must be sent before it is received, so event e_2^1 must occur in any run before event e_1^3 .

A *cut* is a subset of the local history of all processes. If h_i^j denotes the history of process p_i up to and including its j -th event, e_i^j , then a cut C is an n -tuple

$$C = \{h_i^j\} \quad \text{with } i \in \{1, n\} \text{ and } j \in \{1, n_i\}. \quad (10.28)$$

The *frontier* of the cut is an n -tuple consisting of the last event of every process included in the cut. Fig. 10.20 illustrates a *space-time diagram* for a group of three processes, p_1, p_2, p_3 , and it shows two cuts, C_1 and C_2 . C_1 has the frontier $(4, 5, 2)$ frozen after the fourth event of process p_1 , the fifth event of process p_2 and the second event of process p_3 , and C_2 has the frontier $(5, 6, 3)$.

Cuts support the intuition to generate global states based on an exchange of messages between a monitor and a group of processes. A cut represents the instance when requests to report individual state are received by the members of the group. Clearly, not all cuts are meaningful. For example, the cut C_1 with the frontier $(4, 5, 2)$ in Fig. 10.20 violates our intuition regarding causality; it includes e_2^4 , the event triggered by the arrival of message m_3 at process p_2 but does not include e_3^3 , the event triggered

**FIGURE 10.20**

Inconsistent and consistent cuts: The cut $C_1 = (e_1^4, e_2^5, e_3^2)$ is inconsistent because it includes e_2^4 , the event triggered by the arrival of the message m_3 at process p_2 , but does not include e_3^3 , the event triggered by process p_3 sending m_3 , thus the cut C_1 violates causality. On the other hand, $C_2 = (e_1^5, e_2^6, e_3^3)$ is a consistent cut; there is no causal inconsistency; it includes event e_2^6 , the sending of message m_4 , without the effect of it, the event e_3^4 receiving the message by process p_3 .

by process p_3 sending m_3 . In this snapshot, p_3 was frozen after its second event, e_3^2 , before it had the chance to send message m_3 . Causality is violated, and the system cannot ever reach such a state.

Next, we introduce the concepts of consistent and inconsistent cuts and runs. A cut closed under the *causal precedence relationship* is a *consistent cut*. C is a consistent cut if and only if for all events:

$$\forall e, e', (e \in C) \wedge (e' \rightarrow e) \Rightarrow e' \in C. \quad (10.29)$$

A consistent cut establishes an “instance” of a distributed computation; given a consistent cut, we can determine if an event e occurred before the cut. A run R is said to be consistent if the total ordering of events imposed by the run is consistent with the partial order imposed by the causal relation; for all events, $e \rightarrow e'$ implies that e appears before e' in R .

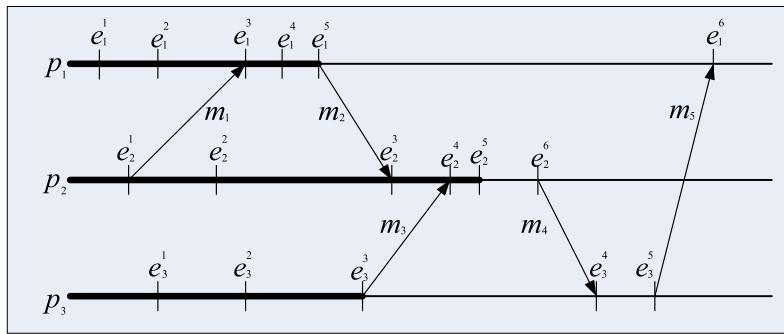
Consider a distributed computation consisting of a group $G = \{p_1, p_2, \dots, p_n\}$ of communicating processes. The *causal history* of event e , $\gamma(e)$, is the smallest consistent cut of G including event e :

$$\gamma(e) = \{e' \in G \mid e' \rightarrow e\} \cup \{e\}. \quad (10.30)$$

The causal history of event e_2^5 in Fig. 10.21 is:

$$\gamma(e_2^5) = \{e_1^1, e_1^2, e_1^3, e_1^4, e_1^5, e_2^1, e_2^2, e_2^3, e_2^4, e_2^5, e_3^1, e_3^2, e_3^3\}. \quad (10.31)$$

This is the smallest consistent cut including e_2^5 . Indeed, if we omit e_3^3 , then the cut $(5, 5, 2)$ would be inconsistent because it would include e_2^4 , the communication event for receiving m_3 , but not e_3^3 , the event caused by sending m_3 . If we omit e_1^5 , the cut $(4, 5, 3)$ would also be inconsistent; it would include e_2^3 but not e_1^5 .

**FIGURE 10.21**

Causal history of event e_2^5 : $\gamma(2^5) = \{e_1^1, e_1^2, e_1^3, e_1^4, e_1^5, e_2^1, e_2^2, e_2^3, e_2^4, e_2^5, e_3^1, e_3^2, e_3^3\}$ is the smallest consistent cut including e_2^5 .

Causal histories can be used as clock values and satisfy the strong clock condition, provided that we equate clock comparison with set inclusion. Indeed,

$$e \rightarrow e' \equiv \gamma(e) \subset \gamma(e'). \quad (10.32)$$

The following algorithm can be used to construct causal histories:

- Each $p_i \in G$ starts with $\theta = \emptyset$.
- Every time p_i receives a message m from p_j it constructs

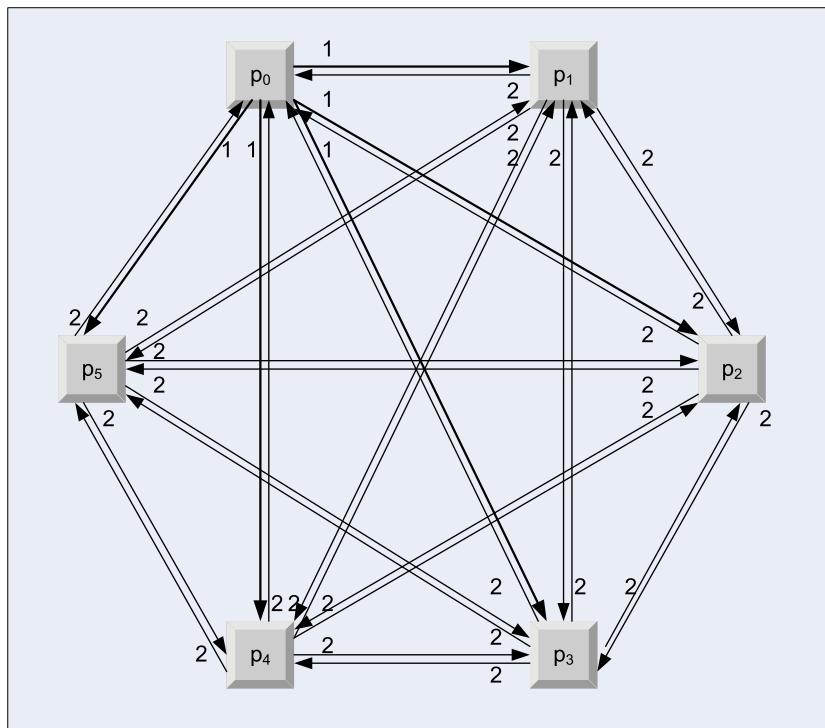
$$\gamma(e_i) = \gamma(e_j) \cup \gamma(e_k), \quad (10.33)$$

with e_i the *receive* event, e_j the previous local event of p_i , and e_k the *send* event of process p_j .

Unfortunately, this concatenation of histories is impractical because the causal histories grow very fast.

Now, we present a protocol to construct consistent global states based on the monitoring concepts discussed in this section. We assume a fully connected network. Recall that given two processes p_i and p_j , the state $\xi_{i,j}$ of the channel from p_i to p_j consists of messages sent by p_i but not yet received by p_j . The snapshot protocol of Chandy and Lamport consists of three steps [93]:

1. Process p_0 sends to itself a “take snapshot” message.
2. Let p_f be the process from which p_i receives the “take snapshot” message for the first time. Upon receiving the message, the process p_i records its local state, σ_i , and relays the “take snapshot” along all its outgoing channels without executing any events on behalf of its underlying computation; channel state $\xi_{f,i}$ is set to empty, and process p_i starts recording messages received over each of its incoming channels.
3. Let p_s be the process from which p_i receives the “take snapshot” message beyond the first time; process p_i stops recording messages along the incoming channel from p_s and declares channel state $\xi_{s,i}$ as those messages that have been recorded.

**FIGURE 10.22**

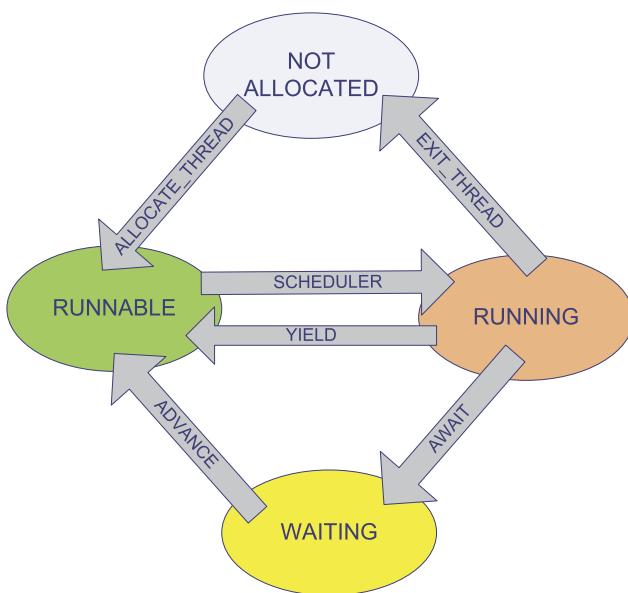
Six processes executing the snapshot protocol.

Each “take snapshot” message crosses each channel exactly once, and every process p_i has made its contribution to the global state; a process records its state the first time it receives a “take snapshot” message and then stops executing the underlying computation for some time. Thus, in a fully connected network with n processes, the protocol requires $n \times (n - 1)$ messages, one on each channel.

For example, consider a set of six processes, each pair of processes being connected by two unidirectional channels, as shown in Fig. 10.22. Assume that all channels are empty, $\xi_{i,j} = 0$, $i \in \{0, 5\}$, $j \in \{0, 5\}$, at the time when process p_0 issues the “take snapshot” message. The actual flow of messages is

- In step 0, p_0 sends to itself the “take snapshot” message.
- In step 1, process p_0 sends five “take snapshot” messages labeled (1) in Fig. 10.22.
- In step 2, each of the five processes, p_1, p_2, p_3, p_4 , and p_5 sends a “take snapshot” message labeled (2) to every other process.

A “take snapshot” message crosses each channel from process p_i to p_j , $i, j \in \{0, 5\}$ exactly once and $6 \times 5 = 30$ messages are exchanged.

**FIGURE 10.23**

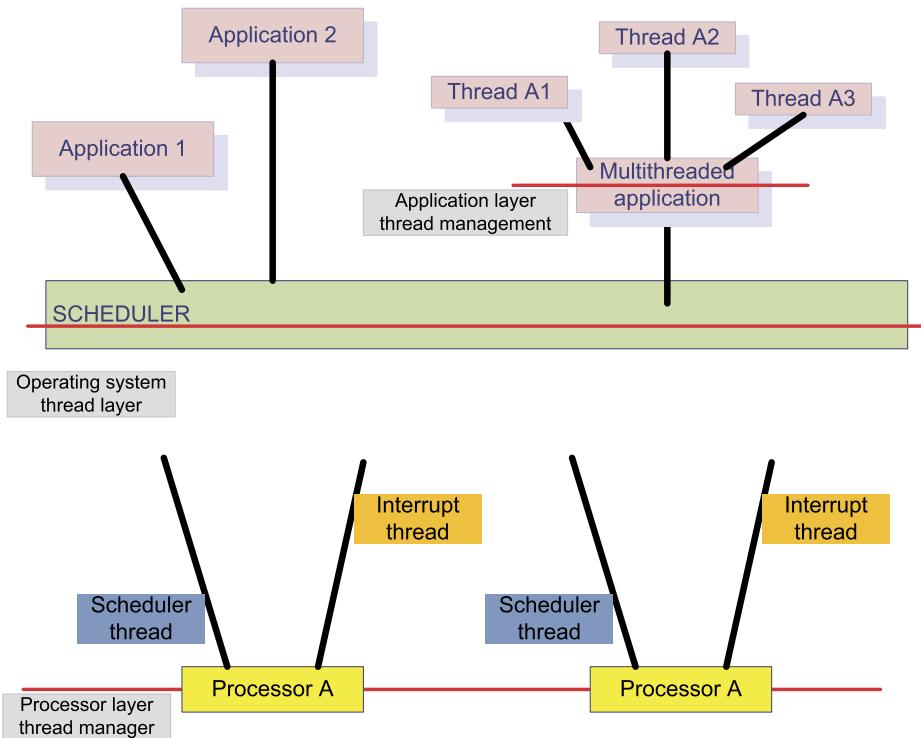
The state of a thread and the actions triggering a change of state.

10.11 Threads and activity coordination

While in the early days of computing, concurrency was analyzed mostly in the context of the system software, nowadays, concurrency is a ubiquitous feature of today's applications. Resources of a single server are insufficient for data-intensive applications and require a careful workload distribution to multiple instances running concurrently on a large number of servers; this is one of the main attractions of cloud computing. Introduction of multi-core processors was a disruptive event; the need to use effectively the cores of a modern processor forced many application developers to implement parallel algorithms and use multithreading.

Many concurrent applications are in the embedded systems area. Embedded systems are a class of reactive systems in which computations are triggered by external events. Such systems populate the Internet of Things (IoT) and are used by the critical infrastructure. A broad spectrum of such applications run multiple threads concurrently to control the ignition of cars, oil processing in a refinery, smart electric meters, heating and cooling systems in homes, or coffee makers. Embedded controllers for reactive real-time applications are implemented as mixed software–hardware systems.

Threads under the microscope. Threads are objects created explicitly to execute streams of instructions by an *allocate thread* action. A thread can be in several states, as shown in Fig. 10.23. Many threads share the core of a processor, and the system scheduler is the authority deciding when a thread gets control of the core and enters a *running* state.

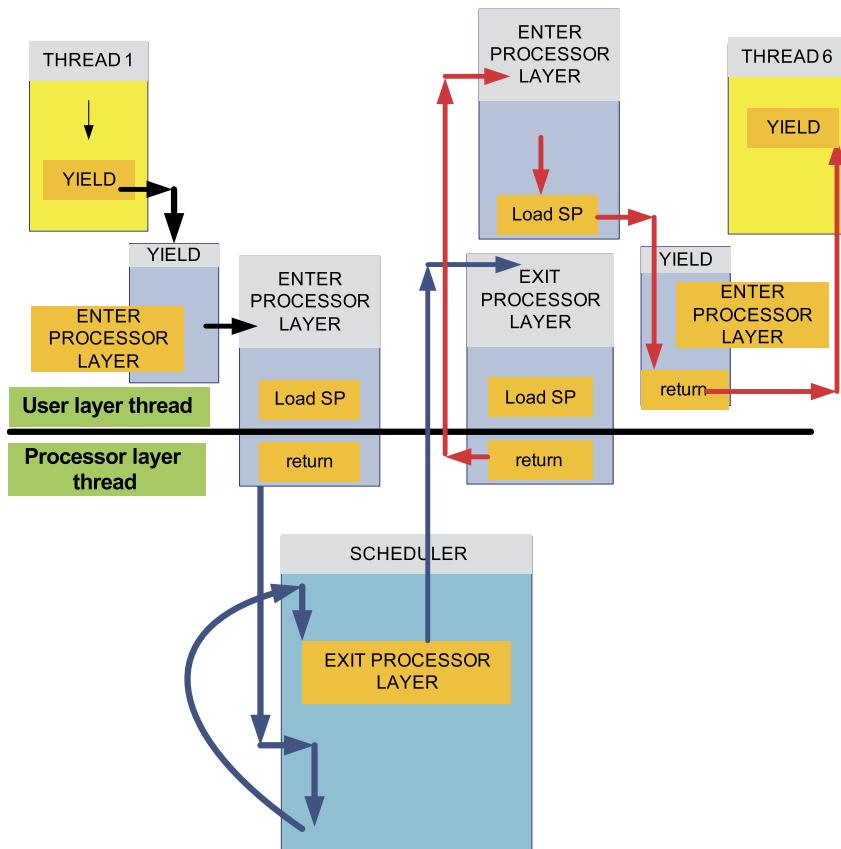
**FIGURE 10.24**

A snapshot of the thread population. Multiple application threads share a core under the control of a scheduler. Multiple operating system-level threads work behind the scene to carry out resource management functions.

The scheduler chooses from the pool of *runnable* threads, the only threads eligible to run. A running thread *yields* the control of the core when it has exhausted the time slot allocated to it or blocks by executing an *await* action, while waiting for the completion of an I/O operation and transition to a *wait* state. The thread could become *runnable* again when the scheduler decides to advance it, e.g., when the I/O operation has finished.

While one may be tempted to think only about application threads, the reality is that the kernel of the operating system operates also a fair number of threads to carry out its functions. Fig. 10.24 provides a snapshot of the thread population including application and operating system threads. Application-level multithreading enables the threads to share the resources allocated to the application, while operating system threads act behind the scene supporting a wide range of resource management functions.

The operation of the scheduler, while totally transparent to application developers, is fairly complex, and multithreading requires several context switches as seen in Fig. 10.25. A context switch of an application thread involves saving the thread state, including registers and the address used when the execution of the suspended thread becomes *runnable* again. Threads are lightweight entities, unlike

**FIGURE 10.25**

Application thread scheduling involves multiple context switches. A context switch saves the current thread state on the system stack. SP is the *stack pointer*.

processes in which information related to the address space, including pointers to the page tables and the process control block, are part of the process state and must be also saved.

Concurrency—the system software side. The kernel of an operating system exploits concurrency for virtualization of system resources such as the processor and the memory. *Virtualization*, covered in depth in Chapter 5, is a system design strategy with a broad range of objectives including:

- Hiding latency and performance enhancement, e.g., schedule a ready-to-run thread when the current thread is waiting for the completion of an I/O operation.
- Avoiding limitations imposed by physical resources, e.g., allow an application to run in a virtual address space of a standard size, rather than be restricted by the physical system memory size.
- Enhancing reliability and performance, as in the case of RAID systems mentioned in Section 2.6.

Sometimes, concurrency is used to describe activities that appear to be executed simultaneously, though only one of them may be active at any given time, as in the case of processor virtualization when multiple threads appear to run concurrently on a single processor. A thread can be suspended due to an external event, and a context switch to a different thread takes place. The state of the suspended thread is saved, the state of another thread ready to proceed is loaded, and then the thread is activated. The suspended thread will be reactivated at a later point in time.

Dealing with some of the effects of concurrency can be very challenging. Context switching could involve multiple components of a OS kernel, including the Virtual Memory Manager (VMM), the Exception Handler (EH), the Scheduler (S), and the Multi-level Memory Manager (MLMM). When a page fault occurs during the fetching of the next instruction, multiple context switches are necessary as shown in Fig. 10.26. The thread experiencing the fault is suspended, and scheduler dispatches another thread ready to run, while the exception handler invokes the multilevel memory manager.

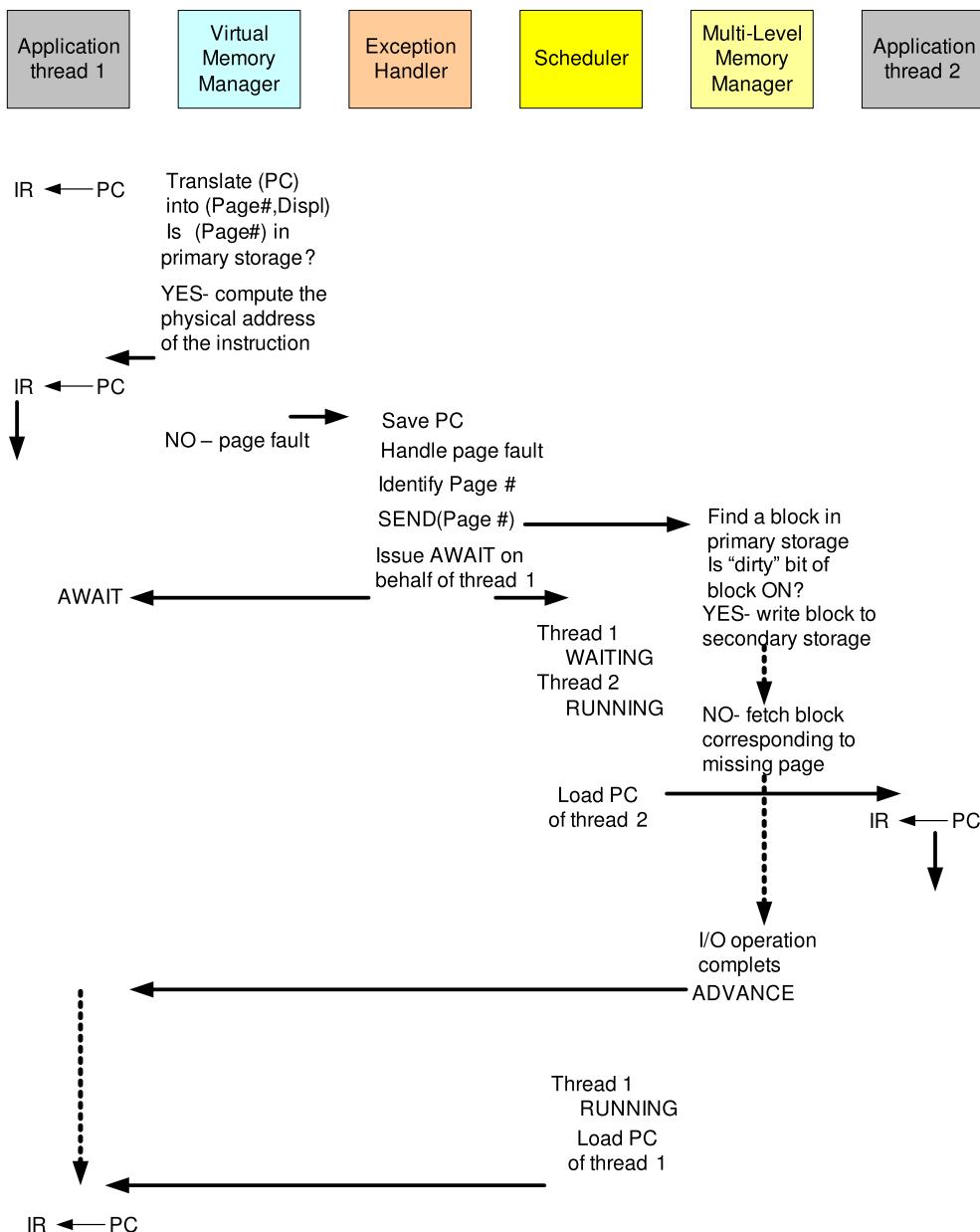
If processor/core sharing seems complicated, operation of a multicore system running multiple virtual machines for applications running under different operating systems is even more complex. Now resource sharing occurs at the operating system level for the threads of one application running under that OS and at the hypervisor level for the threads of different virtual machines, as shown in Fig. 10.27.

Concurrency is often motivated by the desire to enhance the system performance. For example, in a pipelined computer architecture, multiple instructions are in different phases of execution at any given time. Once the pipeline is full, a result is produced at every pipeline cycle; an n -stage pipeline could potentially lead to a speedup by a factor of n . There is always a price to pay for increased performance, and in this example, the price is design complexity and cost. An n -stage pipeline requires n execution units, one for each stage, as well as a coordination unit. It also requires a careful timing analysis to achieve the full speedup.

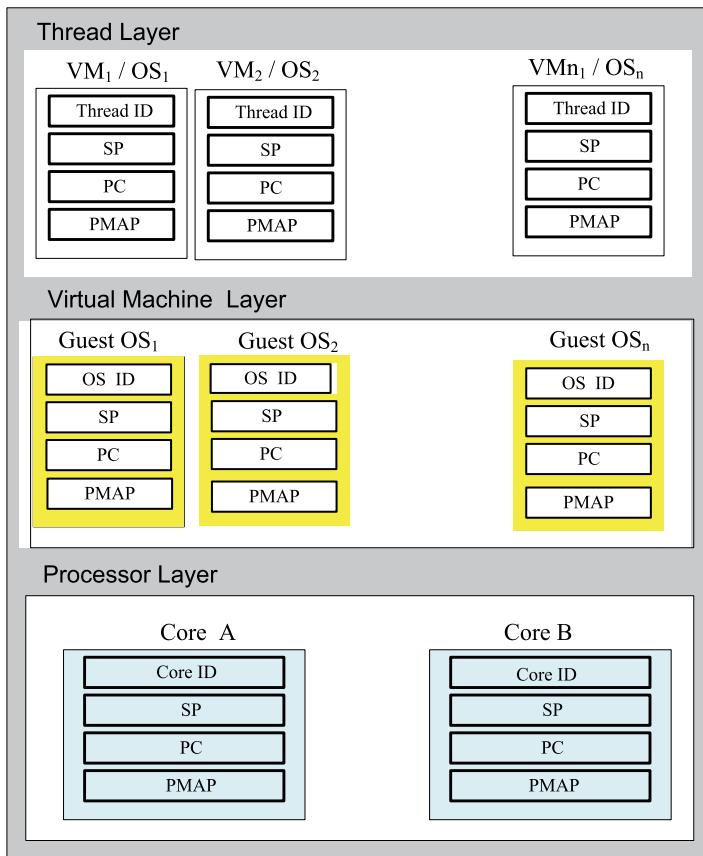
This example shows that the management and the coordination of the concurrent activities increase the complexity of a system. The interaction between pipelining and virtual memory further complicates the functions of the kernel; indeed, one of the instructions in the pipeline could be interrupted due to a page fault and the handling of this case requires special precautions because the state of the processor is difficult to define.

Concurrency—the application software. Concurrency is exploited by application software to speedup a computation and to allow a number of clients to access a service. Parallel applications partition the workload and distribute it to multiple threads running concurrently. Distributed applications, including transaction management systems and applications based on the client-server paradigm discussed in Chapter 3, use extensively concurrency to improve the response time. For example, a web server spawns a new thread when a new request is received, thus multiple server threads run concurrently. A main attraction for hosting web-based applications is cloud elasticity, the ability of a service running on a cloud to acquire resources as needed and to pay for these resources as they are consumed.

Communication channels enable concurrent activities to work in concert and coordinate. Communication protocols enable us to transform noisy and unreliable channels into reliable ones that deliver messages in order. As mentioned earlier, concurrent activities communicate with one another via shared memory or via message passing. Multiple instances of a cloud application and a cloud service and its clients communicate via message passing. The Message Passing Interface (MPI) supports both synchronous and asynchronous communication, and it is often used by parallel and distributed applications.

**FIGURE 10.26**

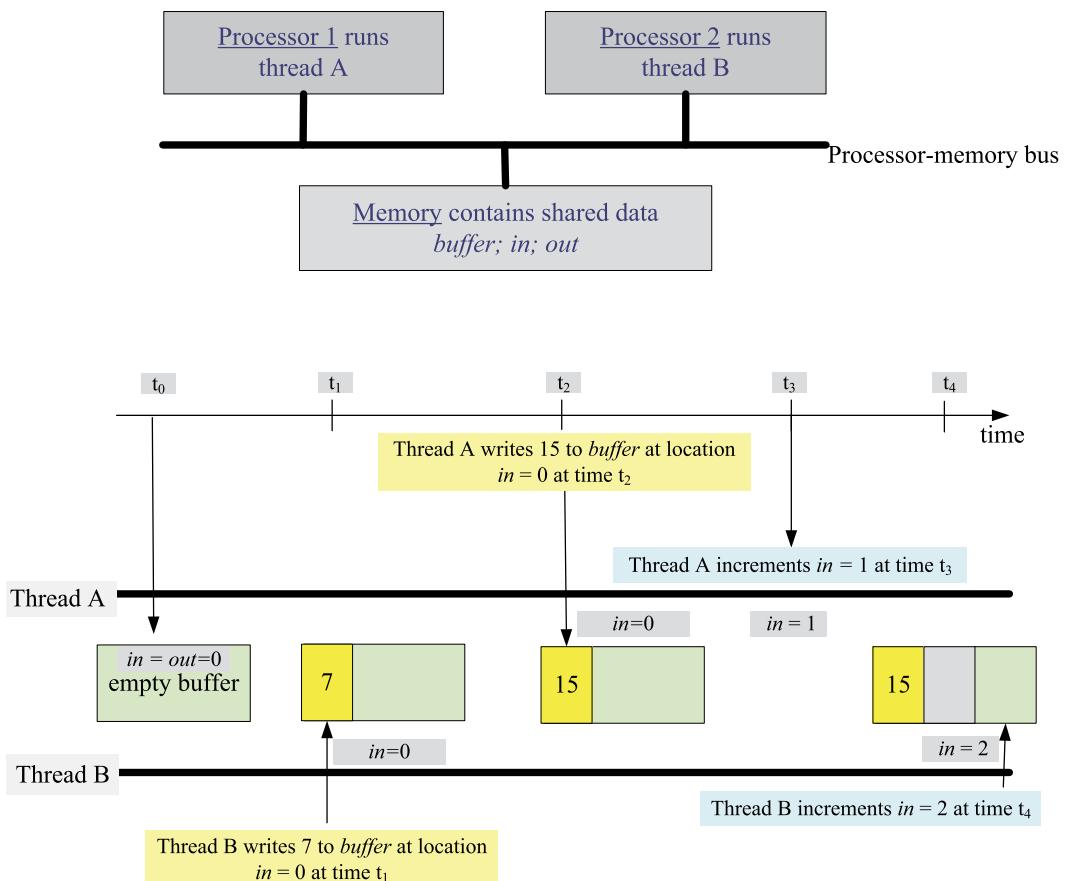
Context switching when a page fault occurs during the instruction fetch phase. IR is the instruction register containing the current instruction, and PC is the program counter pointing to the next instruction to be executed. VMM attempts to translate the virtual address of the next instruction of thread 1 and encounters a page fault. Then, thread 1 is suspended waiting for an event when the page is brought in the physical memory from the disk. The Scheduler dispatches thread 2. To handle the fault, the Exception Handler invokes the MLMM.

**FIGURE 10.27**

Thread multiplexing for a multicore server running multiple virtual machines. Multiple system data structures keep track of the contexts of all threads. The information required for context switching includes the ID, the stack pointer (SP), the program counter (PC), and the page table pointer (PMAP).

Message passing enforces modularity and prevents the communicating activities from *sharing their fate*; a server could fail without affecting the clients not using the service during the period the server was unavailable.

The communication patterns in the case of a parallel application are more structured, while patterns of communication for concurrent activities of a distributed application are more dynamic and unstructured. Barrier synchronization requires the threads running concurrently to wait until all of them have completed the current task before proceeding to the next. Sometimes, one of the activities, a coordinator, mediates communication among concurrent activities, while in other instances, individual threads communicate directly with one another.

**FIGURE 10.28**

Race condition. Initially, at time t_0 , the buffer is empty and $in = 0$. Thread B writes the integer 7 to the buffer at time t_1 . Thread B is slow, incrementing the pointer in takes time and occurs at time t_4 . In the meantime, at time $t_2 < t_4$, a faster thread A writes integer 15 to the buffer, overwrites the content of the first buffer location, and increments the pointer, $in = 1$ at time t_3 . Finally, at time t_4 , thread B increments the pointer $in = 2$.

Coordination of concurrent computations could be quite challenging and involves overhead that ultimately reduces the speedup of parallel computations. Concurrent execution could be very challenging, e.g., it could lead to *race conditions*, an undesirable effect in which the results of concurrent execution depend on the sequence of events. Fig. 10.28 illustrates a race condition in which two threads communicate using a shared data *buffer*. Both threads can write to the *buffer* location pointed at by *in* and can read from *buffer* location pointed at by *out*. When both threads attempt to write at about the same time, the item written by the second thread overwrites the item written by the first thread.

10.12 Critical sections, locks, deadlocks, and atomic actions

Parallel and distributed applications must take special precautions for handling shared resources. For example, consider a financial application when the shared resource is an account record. A thread running on behalf of a transaction first accesses the account to read the current balance, then updates the balance, and, finally, writes back the new balance.

If the thread is interrupted and another thread operating on the same account is allowed to proceed, before the first thread was able to complete the three steps for updating the account, the results of the financial transactions are incorrect. Another challenge is to deal with a transaction involving the transfer from one account to another. A system crash after operation completion on the first account will again lead to an inconsistency; the amount debited from the first account is not credited to the second.

In these cases, as in many other similar situations, a multistep operation should be allowed to proceed to completion without any interruptions, i.e., the operation should be *atomic*. An important observation is that such *atomic actions should not expose the state of the system until the action is completed*. Hiding the internal state of an atomic action reduces the number of states a system can be in, thus it simplifies the design and maintenance of the system. An atomic action is composed of several steps, each of which may fail. Therefore, we have to take additional precautions to avoid exposing the internal state of the system in case of a failure.

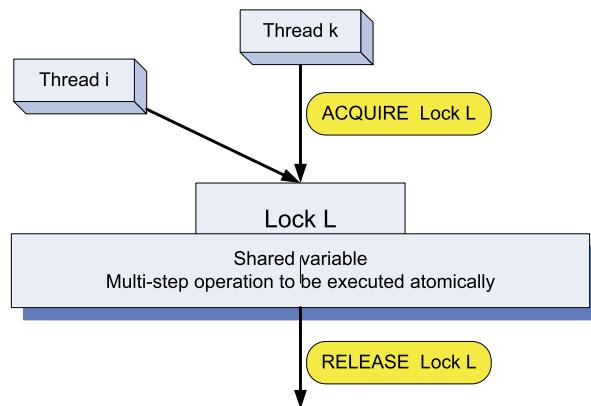
Locks and deadlocks. Concurrency requires a rigorous discipline when threads access shared resources. Concurrent reading of a shared data item is not restricted, while writing a shared data item should be subject to concurrency control. The race conditions discussed in Section 10.11 illustrate the problems created by a hazardous access to a shared resource, the buffer, that both threads attempt to write to. The hazards posed by a lack of concurrency control are ubiquitous. Imagine an embedded system at a power plant in which multiple events occur concurrently and the event signaling a dangerous malfunction of one subsystem is lost.

In all these cases, only one thread should be allowed to modify shared data at any given time, and other threads should only be allowed to read or write this data item only after the first one has finished. This process called *serialization* applies to segments of code called *critical sections* that need to be protected by control mechanisms called *locks*, permitting access to one and only one thread at a time.

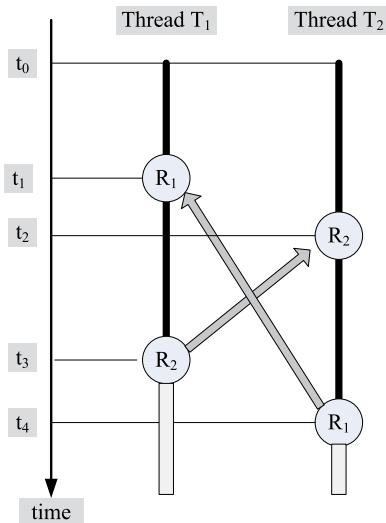
A lock is an object that grants access to a critical section. To enter a critical section, a thread must acquire the lock of that section and after finishing must release the lock, as depicted in Fig. 10.29. Only one thread should be successful when multiple threads attempt to acquire the lock at the same time; the other threads must wait until the lock is released.

One may argue that serialization by locking a data structure is against the very nature of concurrency, allowing multiple computations to run at the same, but, without some form of concurrency control, it is not possible to guarantee the correctness of results of any computation. Lock-free programming [233] is rather challenging and will not be discussed in this chapter.

A lock should be seen as an antidote to uncontrolled concurrency and should be used sparingly and only to protect a critical section. Like any medication, locking has side effects; it does not only increase the execution time, but could lead to deadlocks. Indeed, another potential problem for concurrent execution of multiple processes/threads is the presence of deadlocks. A *deadlock* occurs when processes/threads competing with one another for resources are forced to wait for additional resources held by other processes/threads and none of the processes/threads can finish, as depicted in Fig. 10.30.

**FIGURE 10.29**

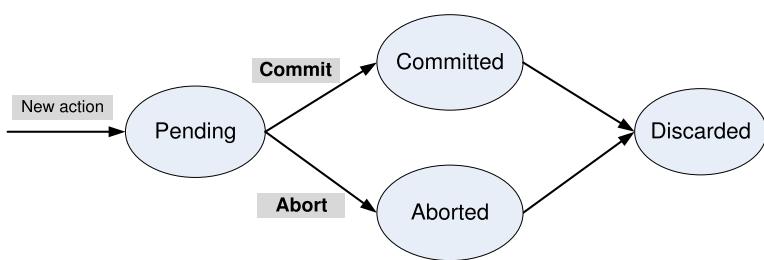
A lock protects a critical section consisting of multiple operations that have to be executed atomically.

**FIGURE 10.30**

Thread deadlock. Threads T_1 and T_2 start concurrent execution at time t_0 . Both need resources R_1 and R_2 to complete execution. T_1 acquires R_1 at time t_1 and T_2 acquires R_2 at time t_2 . At time $t_3 > t_2$, thread T_1 attempts to acquire resource R_2 held by thread T_2 and blocks waiting for it to be released. At time $t_4 > t_3$, thread T_2 attempts to acquire resource R_1 held by thread T_1 and blocks waiting for it to be released. Neither thread can make any progress.

The four Coffman conditions [112] must hold simultaneously for a deadlock to occur:

- Mutual exclusion:* at least one resource must be nonsharable, and only one process or one thread may use the resource at any given time.

**FIGURE 10.31**

The states of an *all-or-nothing* action.

2. *Hold and wait*: at least one process or one thread must hold one or more resources and wait for others.
3. *Nonpreemption*: the scheduler or a monitor should not be able to force a process or a thread holding a resource to relinquish it.
4. *Circular wait*: given the set of n processes or threads $\{P_1, P_2, P_3, \dots, P_n\}$, P_1 should wait for a resource held by P_2 , P_2 should wait for a resource held by P_3 , and so on, P_n should wait for a resource held by P_1 .

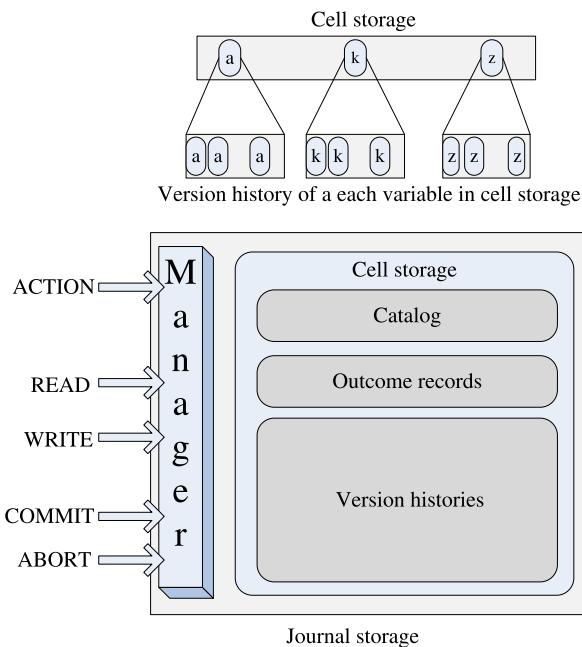
There are other potential problems related to concurrency. When two or more processes/threads continually change their state in response to changes in the other processes, we have a *livelock* condition; the result is that none of the processes can complete its execution. Very often, processes/threads running concurrently are assigned priorities and scheduled based on these priorities. *Priority inversion* occurs when a higher priority process/task is indirectly preempted by a lower priority one.

Atomicity. The discussion of the transaction system suggests that an analysis of atomicity should pay special attention to the basic operation of updating the value of an object in storage. Even to modify the contents of a memory location, several machine instructions must be executed: load the current value in a register, modify the contents of the register, and store back the result.

Atomicity cannot be implemented without some hardware support; indeed, the instruction set of most processors support the *Test-and-Set* instruction that writes to a memory location and returns the old content of that memory cell as non-interruptible operations. Other architectures support *Compare-and-Swap*, an atomic instruction that compares the contents of a memory location to a given value and, only if the two values are the same, modifies the contents of that memory location to a given new value.

Two flavors of atomicity can be distinguished: *all-or-nothing* and *before-or-after* atomicity. *All-or-nothing* means that either the entire atomic action is carried out, or the system is left in the same state it was before the atomic action was attempted; in our examples a transaction is either carried out successfully, or the record targeted by the transaction is returned to its original state. The states of an *all-or-nothing* action are shown in Fig. 10.31.

To guarantee the all-or-nothing property of an action, we have to distinguish preparatory actions that can be undone, from irreversible ones, such as the alteration of the only copy of an object. Such preparatory actions are: allocation of a resource, fetching a page from secondary storage, allocation of memory on the stack, and so on. One of the golden rules of data management is never to change the

**FIGURE 10.32**

Storage models. Cell storage does not support all-or-nothing actions. When we maintain the version histories, it is possible to restore the original content, but we need to encapsulate the data access and provide mechanisms to implement the two phases of an atomic all-or-nothing action. The journal storage does precisely that.

only copy; maintaining the history of changes and a log of all activities allow us to deal with system failures and to ensure consistency.

An all-or-nothing action consists of a *pre-commit* and a *post-commit* phase; during the former it should be possible to backup from it without leaving any trace, while the latter phase should be able to run to completion. The transition from the first to the second phase is called a *commit point*. During the *pre-commit* phase, all steps necessary to prepare the post-commit phase, e.g., check permissions, swap in main memory all pages that may be needed, mount removable media, and allocate stack space, must be carried out; during this phase, no results should be exposed, and no actions that are irreversible should be carried out. Shared resources allocated during the pre-commit cannot be released until after the commit point. The commit step should be the last step of an all-or-nothing action.

A discussion of storage models illustrates the effort required to support all-or-nothing atomicity; see Fig. 10.32. The common storage model implemented by hardware is the so-called *cell storage*, a collection of cells each capable to hold an object, e.g., the primary memory of a computer in which each cell is addressable. Cell storage does not support all-or-nothing actions; once the contents of a cell is changed by an action, there is no way to abort the action and restore the original content of the cell.

To be able to restore a previous value, we have to maintain a *version history* for each variable in the cell storage. The storage model that supports all-or-nothing actions is called *journal storage*. Now,

the cell storage is no longer accessible to the action, but the access is mitigated by a *storage manager*. In addition to the basic primitives to *Read* an existing value and to *Write* a new value in cell storage, the storage manager uniquely identifies an action that changes the value in cell storage and, when the action is aborted, is able to retrieve the version of the variable before the action and restore it. When the action is committed, then the new value should be written to the cell.

Fig. 10.32 shows that for a journal storage, in addition to the version histories of all variables affected by the action, we have to implement a catalog of variables and also to maintain a record to identify each new action. A new action first invokes the *Action* primitive; at that time, an outcome record uniquely identifying the action is created. Then, every time the action accesses a variable, the version history is modified and, finally, the action either invokes a *Commit* or an *Abort* primitive. In the journal storage model, the action is atomic and follows the state transition diagram in Fig. 10.31.

Before-or-after atomicity means that, from the point of view of an external observer, the effect of multiple actions is the same as if these actions have occurred one after another, in some order; a stronger condition is to impose a sequential order among transitions. In our example, the transaction acting on two accounts should either debit the first account and then credit the second one, or leave both accounts unchanged. The order is important because the first account cannot be left with a negative balance.

Atomicity is a critical concept for our efforts to build reliable systems from unreliable components and, at the same time, to support as much parallelism as possible for better performance. Atomicity allows us to deal with unforeseen events and to support coordination of concurrent activities. The unforeseen event could be a system crash, a request to share a control structure, the need to suspend an activity, and so on; in all these cases, we have to save the state of the process or of the entire system to be able to restart it at a later time.

As atomicity is required in many contexts, it is desirable to have a systematic approach rather than an ad hoc one. A systematic approach to atomicity must address several delicate questions:

- How to guarantee that only one atomic action has access to a shared resource at any given time?
- How to return to the original state of the system when an atomic action fails to complete?
- How to ensure that the order of several atomic actions leads to consistent results?

Answers to these questions increase the complexity of the system and often generate additional problems. For example, access to shared resources can be protected by locks, but, when there are multiple shared resources protected by locks, concurrent activities may deadlock. A *lock* is a construct that enforces sequential access to a shared resource; such actions are packaged in the *critical sections* of the code. If the lock is not set, a thread first locks the access, then enters the critical section, and finally unlocks it; a thread wishing to enter the critical section finds the lock set and waits for the lock to be reset. A lock can be implemented using the hardware instructions supporting atomicity.

Semaphores and monitors are more elaborate structures ensuring serial access. Semaphores force processes to queue up when the lock is set and are released from this queue and allowed to enter the critical section one by one. Monitors provide special procedures to access the shared data; see Fig. 10.33. *The mechanisms for the process coordination we described require the cooperation of all activities*, the same way traffic lights prevent accidents only as long as the drivers follow the rules.

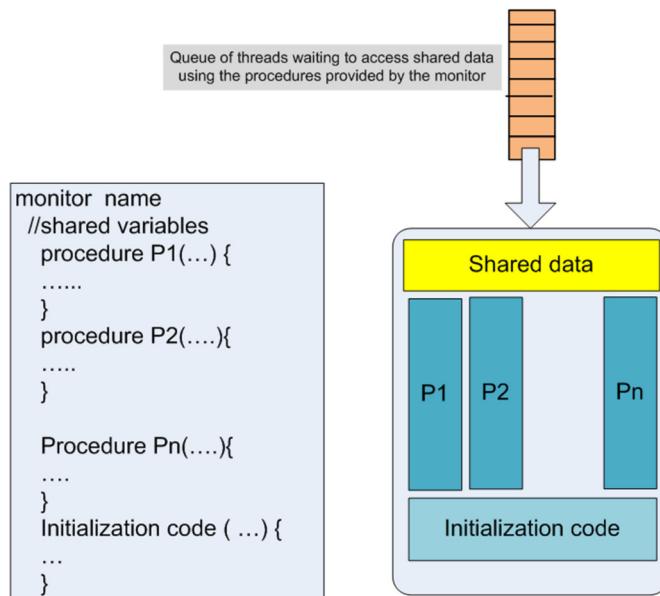


FIGURE 10.33

A monitor provides special procedures to access the data in a critical section.

10.13 Consensus protocols

Consensus is a pervasive problem in many areas of human endeavor; consensus is the process of agreeing to one of several alternates proposed by a number of agents. We restrict our discussion to a distributed system when a set of processes must reach consensus on a single proposed value.

No fault-tolerant consensus protocol can guarantee progress [175], but protocols that guarantee freedom from inconsistencies (safety) have been developed. A family of protocols to reach consensus based on a finite state machine approach is called *Paxos*.⁵

A fair number of contributions to the family of Paxos protocols are discussed in the literature. Leslie Lamport proposed several versions of the protocol including Disk Paxos, Cheap Paxos, Fast Paxos, Vertical Paxos, Stoppable Paxos, Byzantizing Paxos by Refinement, Generalized Consensus and Paxos, and Leaderless Byzantine Paxos. Lamport also published a paper on the fictional part-time parliament in Paxos [292] and a layman's dissection of the protocol [293].

The *consensus service* consists of a set of n processes. *Clients* send requests to processes and propose a value and wait for a response; the goal is to get the set of processes to reach consensus on a

⁵ Paxos is a small Greek island in the Ionian Sea; a fictional consensus procedure is attributed to an ancient Paxos legislative body. The island had a part-time parliament because its inhabitants were more interested in other activities than in civic work; “the problem of governing with a part-time parliament bears a remarkable correspondence to the problem faced by today’s fault-tolerant distributed systems, where legislators correspond to processes and leaving the Chamber corresponds to failing” according to Leslie Lamport [292] (for additional papers, see <http://research.microsoft.com/en-us/um/people/lamport/pubs/pubs.html>).

single proposed value. A *quorum* is a subset of all acceptors. A proposal has a proposal number pn and contains a value v . Several types of requests such as *prepare* and *accept* flow through the system.

The *basic Paxos* considers several types of entities: (a) *client*, an agent that issues a request and waits for a response; (b) *proposer*, an agent with the mission to advocate a request from a client, convince the acceptors to agree on the value proposed by a client, and act as a coordinator to move the protocol forward in case of conflicts; (c) *acceptor*, an agent acting as the fault-tolerant “memory” of the protocol; (d) *learner*, an agent acting as the replication factor of the protocol and taking action once a request has been agreed upon; and finally, (e) the *leader*, a distinguished proposer.

In a typical deployment of the algorithm, an entity plays three roles, as proposer, acceptor, and learner. Then, the flow of messages can be described as follows [293]: “clients send messages to a leader; during normal operations the leader receives the client’s command, assigns it a new command number i , and then begins the i -th instance of the consensus algorithm by sending messages to a set of acceptor processes.” By merging the roles, the protocol “collapses” into an efficient client–master–replica style protocol.

The *basic Paxos* protocol is based on several assumptions about the processors and the network:

- Processes run on processors and communicate through a network; the processors and the network may experience failures, but not Byzantine failures. A *Byzantine failure* is a fault presenting different symptoms to different observers. In a distributed system, a Byzantine failure could be an *omission failure*, e.g., a crash failure, failure to receive a request or to send a response; it could also be a *commission failure*, e.g., process a request incorrectly, corrupt the local state, and/or send an incorrect or inconsistent response to a request.
- Processors: (i) operate at arbitrary speeds; (ii) have stable storage and may rejoin the protocol after a failure; and (iii) can send messages to any other processor.
- The network: (i) may lose, reorder, or duplicate messages; and (ii) messages are sent asynchronously and may take arbitrarily long time to reach the destination.

A proposal consists of a pair of numbers, a unique proposal number and a proposed value, (pn, v) ; multiple proposals may propose the same value v . A value is chosen if a simple majority of acceptors have accepted it. We need to guarantee that at most one value can be chosen; otherwise, there is no consensus. The two phases of the algorithm are:

Phase I.

1. *Proposal preparation*: a proposer (the leader) sends a proposal ($pn = k, v$). The proposer chooses a proposal number $pn = k$ and sends a *prepare message* to a majority of acceptors requesting:
 - that a proposal with $pn < k$ should not be accepted;
 - the $pn < k$ of the highest number proposal already accepted by each acceptor.
2. *Proposal promise*: An acceptor must remember the highest proposal number it has ever accepted and the highest proposal number it has ever responded to. The acceptor can accept a proposal with $pn = k$ if and only if it has not responded to a prepare request with $pn > k$; if it has already replied to a prepare request for a proposal with $pn > k$, then it should not reply. Lost messages are treated as an acceptor that chooses not to respond.

Phase II.

1. *Accept request*: if the majority of acceptors respond, then the proposer chooses the value v of the proposal as follows:

- the value v of the highest proposal number selected from all the responses;
- an arbitrary value if no proposal was issued by any of the proposers.

The proposer sends an *accept request* message to a quorum of acceptors including $(pn = k, v)$

2. *Accept*: If an acceptor receives an *accept message* for a proposal with the proposal number $pn = k$, it must accept it if and only if it has not already promised to consider proposals with a $pn > k$. If it accepts the proposal, it should register the value v and send an *accept* message to the proposer and to every learner; if it does not accept the proposal, it should ignore the request.

Several algorithm properties are important to show its correctness: (1) a proposal number is unique; (2) any two sets of acceptors have at least one acceptor in common; and (3) the value sent out in Phase 2 of the algorithm is the value of the highest numbered proposal of all the responses in Phase 1.

Fig. 10.34 illustrates the flow of messages for the consensus protocol. A detailed analysis of the message flows for various failure scenarios and of the properties of the protocol can be found in [293]. We only mention that the protocol defines three safety properties: (1) nontriviality—the only values that can be learned are proposed values; (2) consistency—at most one value can be learned; and (3) liveness—if a value v has been proposed, eventually, every learner will learn some value, provided that sufficient processors remain non-faulty. Fig. 10.35 shows the message exchange when there are three actors involved.

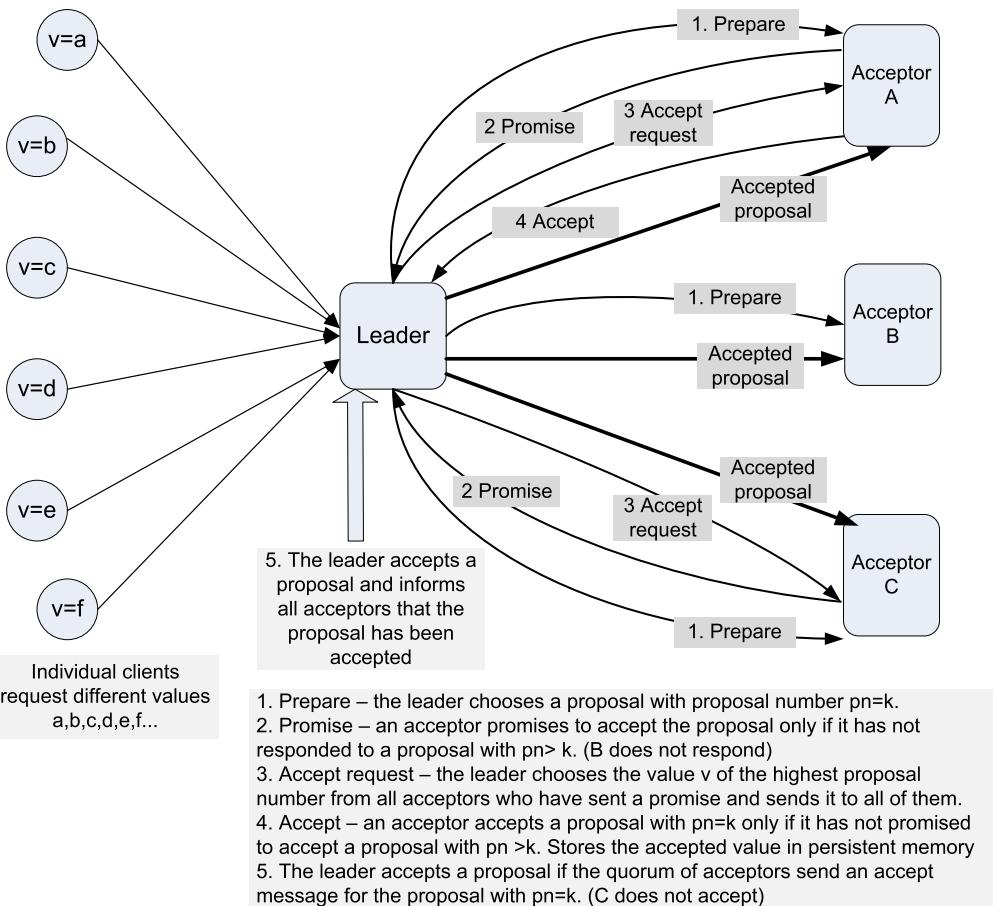
Chubby, a locking service based on the Paxos algorithm, is discussed in Section 7.7. A distributed coordination system discussed in Section 11.4, the Zookeeper, borrows several ideas from the Paxos algorithm:

- i. A leader proposes values to the followers.
- ii. Leaders wait for acknowledgments from a quorum of followers before considering a proposal committed (learned);
- iii. Proposals include epoch numbers, which are similar to ballot numbers in Paxos.

10.14 Load balancing

A general formulation of the load balancing problem is to evenly distribute \mathcal{N} objects to \mathcal{P} places. Another formulation of the load balancing is in the context of placing m balls into $n < m$ bins, chosen independently and uniformly at random. The question in this case is to find the maximum number of balls in any bin. Load balancing can also be formulated in the context of hashing as the problem of placing m items sequentially in $n < m$ bins; the task is determining the longest time for finding an item.

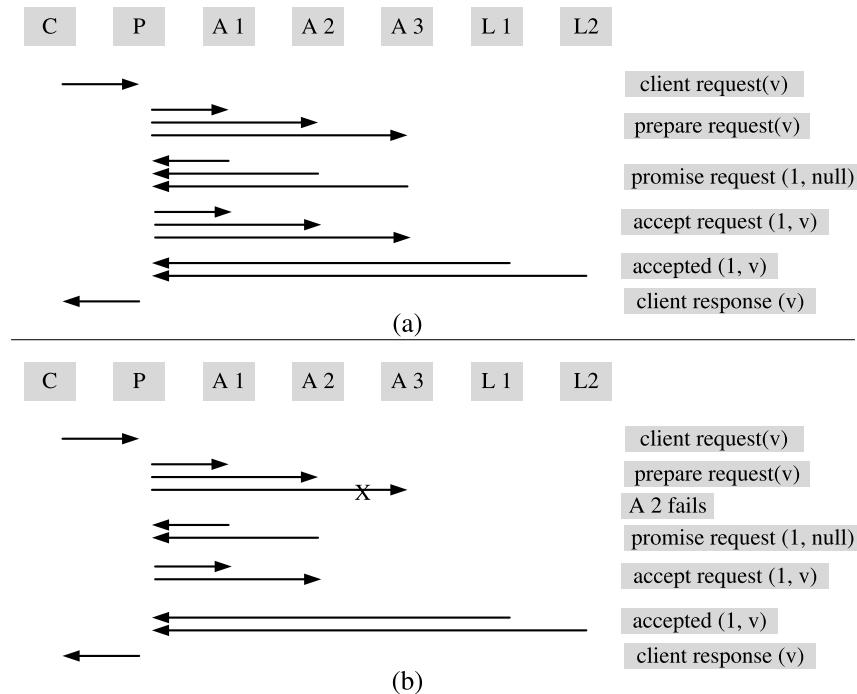
The load balancing problem in a distributed system is formulated as follows: given a set \mathcal{T} of tasks, distribute them to a set of \mathcal{P} processors, which compute at the same rate, such that only one task can run at any given time on one processor; there is no preemption and each task runs to completion. Knowing the execution time of each task, the question is how to distribute them to minimize the completion time. Unfortunately, load balancing, as well as scheduling problems, are NP-complete [188].

**FIGURE 10.34**

The flow of messages for the Paxos consensus algorithm. Individual clients propose different values to the leader who initiates the algorithm. Acceptor A accepts the value in message with proposal number $pn = k$; acceptor B does not respond with a promise, while acceptor C responds with a promise, but ultimately does not accept the proposal.

The importance of load balancing is undeniable, and practical solutions to overcome the algorithmic complexity are widely used. For example, randomization suggests to distribute the tasks to processors chosen independently and uniformly at random. This random distribution strategy should lead to an almost equal load of processors, provided that there are enough tasks and that the distribution of the task execution times is rather narrow. Several other heuristics are used in practice.

The balls-and-bins model. The load balancing problem is often discussed using the balls-and-bins model. In this model, we define the *load* of a bin as the number of balls in the bin. The question asked is: What is $\max(\mathcal{L}_i)$, $1 \leq i \leq n$, the maximum load in any bin, once all the n balls have chosen a

**FIGURE 10.35**

The basic Paxos with three actors: proposer (P), three acceptors (A1, A2, A3), and two learners (L1, L2). The client (C) sends a request to one of the actors playing the role of a proposer. The entities involved are: (a) successful first round when there are no failures; (b) successful first round of Paxos when an acceptor fails.

bin independently and uniformly at random? The answer is that with *high probability*, namely, with a probability $p \geq 1 - \mathcal{O}(1/n)$ [202],

$$\max(\mathcal{L}_i) \approx \frac{\log n}{\log \log n}. \quad (10.34)$$

This result applies also to task scheduling in a distributed system. Interestingly enough, this solution does not involve communication among the tasks, the processors, or among the tasks and the processors.

A rather surprising result proven in [39] is that better load balance is achieved when the balls are placed sequentially, and for each ball, we choose two bins independently and uniformly at random, then place the ball into the less full bin. In this case, the maximal load, $\max(\mathcal{L}_i^2)$, $1 \leq i \leq n$ is

$$\max(\mathcal{L}_i^2) \leq \frac{\log \log n}{\log 2} + \mathcal{O}(1) \text{ with high probability.} \quad (10.35)$$

This result, discussed in [353,355] and called *the power of two choices*, shows that having two or more choices leads to an exponential improvement of the load balance. The following discussion is based on the *layered induction approach* in [39], when the number of bins that contain at least j balls conditioned on the number of bins that contain at least $(j - 1)$ balls is inductively bound.

If the choice for each ball is extended from 2 to d bins, then the result is further improved. The GREEDY algorithm in [39] considers an (m, n, d) problem: n initially empty bins, m balls to be placed sequentially in the bins, and d choices made independently and uniformly at random with replacement. Each ball is placed in the least loaded of the d bins, and ties are broken arbitrarily. Then, the maximum load, the maximum number of balls in a bin, has an upper bound of

$$\max(\mathcal{L}_i^d) \leq \frac{\log \log n}{\log d} + \mathcal{O}(1) \text{ with high probability.} \quad (10.36)$$

An intuitive justification of this results is discussed next. Call β_k the number of bins with *at least* k balls stacked on top of one another in the order they have been placed in the bin. The *height* of the top ball in a bin with k balls is k .

We wish to determine β_{k+1} , a high probability upper bound of the cardinality of bins loaded with at least $k + 1$ balls. A bin will contain at least $k + 1$ balls if in the previous round, it had at least k balls. Recall that there are at least β_k such bins; thus the probability of choosing a bin with k or more balls from the set of n bins is β_k/n . But there are $d > 2$ choices, therefore the probability that a ball lands in a bin already containing k or more balls drops at each step at least quadratically and is

$$p_{n,k,d} = \left(\frac{\beta_k}{n} \right)^d. \quad (10.37)$$

The number of balls with height at least $(k + 1)$ is dominated by a Bernoulli random variable with the probability of success equals to $p_{n,k,d}$. This implies that

$$\beta_{k+1} \leq c \times \left[n \times \left(\frac{\beta_k}{n} \right)^d \right] \quad (10.38)$$

with c a constant. It follows that, after $j = \mathcal{O}(\log \log n)$ steps, the fraction β_k/n drops below $1/n$; thus, $\beta_j < 1$.

The sequential ball placement required in the algorithms discussed in this section and the decision to choose one of two bins, or one of d bins, deserves further scrutiny. It implies either a centralized system where an agent makes the choice, or some form of communication between balls to reach an agreement about the placement strategy. Communication is expensive, e.g., during the time of a short message exchange, a modern processor could execute several billion floating-point operations.

A tradeoff between load balance and communication is inevitable; we can reduce the maximum load only through coordination, thus the price to pay is increased communication. In a distributed system, a server is not aware of the tasks it has not communicated with, while tasks are unaware of the actions of other tasks and may only know the load of the servers. Global coordination among tasks is prohibitively expensive.

Parallelization of the randomized load balance. One of the questions examined in [353,355] is how to parallelize the randomized load balance. Once again, an extended balls-and-bins model is used and

the goal is to minimize the maximum load and the communication complexity. Each one of the m balls begins by choosing d out of the n bins as prospective destination.

The choices are made independently and uniformly at random with replacement of balls, and the final decision of the destination bin requires r rounds of communication, with two stages per round. At each stage communication is done in parallel using short messages including an ID or an index. In the first stage each ball sends messages to all prospective bins, and in the second one each bin sends messages to all balls the bin has received a message from. During the final round, balls commit to a bin.

The strategies analyzed are symmetric, all balls and bins use the same algorithm, and all possible destinations are chosen independently and uniformly at random. An algorithm is *asynchronous* if an entity, a ball or a bin, does not have to wait for a round to complete; it only has to wait for messages addressed to it, rather than waiting for messages addressed to another entity. A round is *synchronous* if a barrier synchronization is required between some pairs of rounds.

A load lower bound for a broad class of algorithms like the one in [39] with r -rounds of communication derived in [355] is

$$\Omega\left(\sqrt[r]{\frac{\log n}{\log \log n}}\right), \quad (10.39)$$

with at least constant probability. Therefore, no algorithm can achieve a maximum load $\mathcal{O}(\log \log n)$ with high probability in a constant number of communication rounds.

A random graph $G(v, e)$ is used to represent the model and to derive this result. Each bin is a vertex v in this graph, and each ball is an undirected edge, e . When $d = 2$, the vertices of the two edges of a ball correspond to the two prospective bins. There are no self-loops in this graph where S denotes the set of edges. Multiple edges correspond to two balls that have chosen the same pair of bins. Selection of a bin by a ball transforms the undirected edge representing the ball into a directed edge oriented toward the vertex, or bin, the ball chooses as its destination. The goal is to minimize the maximum in-degree over the set of all graph vertices to avoid conflicts.

$\mathcal{N}(e)$, the *neighborhood* of an edge $e \in S$, is the set of all edges incident to an endpoint of e , and

$$\mathcal{N}(S) = \cup_{e \in S} \mathcal{N}(e). \quad (10.40)$$

Similarly, $\mathcal{N}(v)$ is the *neighborhood* of a vertex v . $\mathcal{N}_l(e)$, the l -*neighborhood* of an edge $e \in S$ is defined inductively as

$$\mathcal{N}_1(e) = \mathcal{N}(e), \quad \mathcal{N}_l(e) = \mathcal{N}(\mathcal{N}_{l-1}(e)). \quad (10.41)$$

$\mathcal{N}_{l,x}(e)$, the (l, x) -*neighborhood* of an edge $e(x, y)$, is defined inductively as

$$\mathcal{N}_{1,x}(e) = \mathcal{N}(x) - \{e\}, \quad \mathcal{N}_{l,x}(e) = \mathcal{N}(\mathcal{N}_{l-1,x}(e)) - \{e\}. \quad (10.42)$$

In an r round protocol, the balls make their choice in the final round; therefore each ball knows everything only about the balls in its $(r - 1)$ neighborhood.

A ball $e = (x; y)$ learns from each bin about its l -*neighborhood* consisting of two subgraphs corresponding to $\mathcal{N}_{l,x}(e)$ and $\mathcal{N}_{l,y}(e)$. When the two subgraphs of the ball's l -*neighborhood* are isomorphic rooted trees, with the roots x and y , we say that the ball has a symmetric l -*neighborhood* or, that the ball is *confused*, and then the ball chooses the destination bin using a fair coin flip.

A tree of depth r , in which the root has degree T and each internal node has $T - 1$ children, is called a (T, r) -rooted, *balanced tree*. A (T, r) tree in graph G is *isolated* if it is a connected component of G with no edges of multiplicity greater than one. A random graph with n vertices and n edges contains an isolated $(T, 2)$ tree with

$$T = \left(\sqrt{2} - \mathcal{O}(1) \right) \sqrt{\frac{\log n}{\log \log n}}, \quad (10.43)$$

with constant probability as shown in [355]. A corollary of this statement is that any nonadaptive, symmetric load distribution strategy for the balls-and-bins problem, with n balls and n bins and with $d = 2$ and $r = 2$, has a final load of at least

$$\left(\frac{\sqrt{2}}{2} - \mathcal{O}(1) \right) \sqrt{\frac{\log n}{\log \log n}}, \quad (10.44)$$

with at least constant probability. Indeed, half of the confused balls (edges) in an isolated $(T, 2)$ tree adjacent to the root will orient themselves towards the root because we have assumed that the balls flip an unbiased coin to choose the bin. This result can be extended for a range of r and d .

The balls-and-bins model has applications to hashing. The hashing implementation discussed in [277] uses a single hash function to map keys to entries in a table, and in case of a collision, i.e., when two or more keys map to the same table entry, all the colliding keys are stored in a linked list called a *chain*. The table entries are heads of chains and the longest search time occurs for the longest chain. The length of the longest chain is $\mathcal{O}(\frac{\log n}{\log \log n})$ with a high probability, when n keys are inserted into a table with n entries and each key is mapped to an entry of the table independently and uniformly, a process known as *perfect random hashing*.

The search time can be substantially reduced by using two hash functions and placing an item in the shorter of the two chains [266]. To search for an element, we have to search through the chains linked to the two entries given by both hash functions. If the n keys are sequentially inserted into the table with n entries, the length of the longest chain, thus, the maximum time to find an item, is $\mathcal{O}(\log \log n)$ with high probability.

The two-choice paradigm can be applied effectively to routing virtual circuits in interconnection networks with low congestion. The paradigm is also used to optimize the emulation of shared-memory multiprocessor (SMM) systems on a distributed-memory multiprocessor systems (DMM). The emulation algorithm should minimize the time needed by the DMM to emulate one step of the SMM.

The layered induction approach is also used in dynamic scenarios, e.g., when a new ball is inserted in the system [354]. Another technique to analyze load balancing based on the balls-and-bins model is the *witness tree*. To compute a bound for the probability of a “heavily-loaded” system event, we have to identify a *witness tree* of events and then estimate the probability that the witness tree occurs. This probability can be bounded by enumerating all possible witness trees and summing their individual probabilities of occurrence.

10.15 Multithreading in Java; FlumeJava; Apache Crunch

Java is a general-purpose computer programming language designed with portability in mind at Sun Microsystems.⁶ Java applications are typically compiled to bytecode and can run on a Java Virtual Machine (JVM), regardless of the computer architecture. Java is a class-based, object-oriented language with support for concurrency. It is one of the most popular programming languages, and it is widely used for a wide range of applications running on mobile devices and computer clouds.

Java Threads. Java supports processes and threads. Recall that a process has a self-contained execution environment and has its own private address space and run-time resources. A thread is a lightweight entity within a process. A Java application starts with one thread, the *main thread*, which can create additional threads.

Memory consistency errors occur when different threads have inconsistent views of the same data. Synchronized methods and synchronized statements are the two idioms for synchronization. Serialization of *critical sections* is protected by specifying the *synchronized* attribute in the definition of a class or method. This guarantees that only one thread can execute the critical section and each thread entering the section sees the modification done. Synchronized statements must specify the object that provides the intrinsic lock.

The current versions of Java support atomic operations of several datatypes with methods, such as *getAndDecrement()*, *getAndIncrement()* and *getAndSet()*. An effective way to control data sharing among threads is to share only immutable data among threads. A class is made *immutable* by marking all its fields as *final* and declaring the class as *final*.

A *Thread* in the *java.lang.Thread* class executes an object of type *java.lang.Runnable*. The *java.util.concurrent* package provides better support for concurrency than the *Thread* class. This package reduces the overhead for thread creation and prevents too many threads overloading the CPU and depleting the available storage. A *thread pool* is a collection of *Runnable* objects and contains a queue of tasks waiting to get executed.

Threads can communicate with one another via *interrupts*. A thread sends an interrupt by invoking an *interrupt* on the *Thread* object to the thread to be interrupted. The thread to be interrupted is expected to support its own interruption. *Thread.sleep* causes the current thread to suspend execution for a specified period.

The executor framework works with *Runnable* objects which cannot return results to the caller. The alternative is to use *java.util.concurrent.Callable*. A *Callable* object returns an object of type *java.util.concurrent.Future*. The *Future* object can be used to check the status of a *Callable* object and to retrieve the result from it. Yet, the *Future* interface has limitations for asynchronous execution, and the *CompletableFuture* extends the functionality of the *Future* interface for asynchronous execution.

Nonblocking algorithms based on low-level atomic hardware primitives such as compare-and-swap (CAS) are supported by Java 5.0 and later versions. The fork-join framework introduced in Java 7 supports the distribution of work to several workers and then waiting for their completion. The *join* method allows one thread to wait for completion of another.

⁶ The design of Java was initiated in 1991 by James Gosling, Mike Sheridan, and Patrick Naughton with a C/C++-style syntax. Five principles guided its design: (1) simple, object-oriented, and familiar; (2) architecture-neutral and portable; (3) robust and secure; (4) interpreted, threaded, and dynamic; and (5) high performance. Java 1.0 was released in 1995. Java 8 is the only version currently supported for free by Oracle, a company that acquired Sun Microsystems in 2010.

FlumeJava. A Java library used to develop, test, and run efficiently data parallel pipelines is described in [90]. FlumeJava is used to develop data parallel applications such as MapReduce, discussed in Section 11.5.

At the heart of the system is the concept of *parallel collection*, which abstracts the details of data representation. Data in a parallel collection can be an in-memory data structure, one or more files, BigTable, discussed in Section 7.11, or a MySQL database. Data-parallel computations are implemented by composition of several operations for parallel collections.

In turn, parallel operations are implemented using *deferred evaluation*. The invocation of a parallel operation records the operation and its arguments in an internal graph structure representing the execution plan. Once completed, the execution plan is optimized.

The most important classes of the FlumeJava library are the *Pcollection* $< T >$ used to specify an immutable bag of elements of type T and the *PTable* $< K, V >$ representing an immutable multimap with keys of type K and values of type V . The internal state of a *PCollection* object is either *deferred* or *materialized*, i.e., not yet computed or computed, respectively. The *PObject* $< T >$ class is a container for a single Java object of type T and can be either deferred or materialized.

parallelDo() supports element-wise computation over an input *PCollection* $< T >$ to produce a new output *PCollection* $< S >$. This primitive takes as the main argument a *DoFn* $< T, S >$, a function-like object defining how to map each value in the input into zero or more values in the output. In the following example from [90], *collectionOf(strings())* specifies that the *parallelDo()* operation should produce an unordered *PCollection* whose *String* elements should be encoded using UTF-8.⁷

```
Pcollection<String> words =
    lines.parallelDo(new DoFn<String, String> () {
        void process (String line, EmitFn<String> emitFn {
            for (String word : splitIntoWords(line) ) {
                emitFn.emit(word);
            }
        }
    }, collectionOf(strings( )) );
```

Other primitive operations are: *groupByKey()*, *combineValues()* and *flatten()*.

- *groupByKey()* converts a multimap of type *PTable* $< K, V >$. Multiple key/value pairs may share the same key into a unimap of type *PTable* $< K, Collection < V >$, where each key maps to an unordered, plain Java Collection of all the values with that key.
- *combineValues()* takes an input *PTable* $< K, Collection < V >$ and an associative combining function on V s, and returns a *PTable* $< K, V >$ where each input collection of values has been combined into a single output value.
- *flatten()* takes a list of *PCollection* $< T >$ s and returns a single *PCollection* $< T >$ that contains all the elements of the input *PCollections*.

⁷ UTF-8 is a character encoding standard defined by Unicode capable of encoding all possible characters. The encoding is variable-length and uses 8-bit code units.

Pipelined operations are implemented by concatenation of functions. For example, if the output of function f is applied as input of function g in a *Parallel Do* operation, then two *Parallel Do* compute f and $f \otimes g$. The optimizer is only concerned with the structure of the execution plan and not with the optimization of user-defined functions.

FlumeJava traverses the operations in the plan of a batch application in forward topological order and executes each operation in turn. Independent operations are executed simultaneously. FlumeJava exploits not only the task parallelism but also the data parallelism within operations.

Apache Crunch, modeled after FlumeJava, simplifies creation of data pipelines on top of Apache Hadoop. Crunch was designed to use MapReduce effectively and is often used in conjunction with Hive and Pig, discussed in Section 11.8.

A Crunch pipeline creates a collection of Pig scripts and Hive queries. Crunch is fast, but slightly slower than a hand-tuned pipeline developed with the MapReduce APIs, according to <https://attic.apache.org/projects/crunch.html>.

10.16 History notes and further readings

In 1965, Edsger Dijkstra posed the problem of synchronizing N processes, each with a *critical section*, subject to two conditions: *mutual exclusion*—no two critical sections are executed concurrently; and *livelock freedom*—if several processes wait to enter critical sections, one of them will eventually succeed, [146]. Lamport comments “Dijkstra was aware from the beginning of how subtle concurrent algorithms are and how easy it is to get them wrong. He wrote a careful proof of his algorithm. The computational model implicit in his reasoning is that an execution is represented as a sequence of states, where a state consists of an assignment of values to the algorithm’s variables plus other necessary information such as the control state of each process (what code it will execute next).”

Producer-consumer synchronization was the second fundamental concurrent programming problem identified by Dijkstra. An equivalent formulation of the problem is: Given a bounded FIFO (first-in-first-out), the producer stores data into an N -element buffer and the consumer retrieves the data. The algorithm uses three variables: N —the buffer size, in —the infinite sequence of unread input values, and out —the sequence of values output so far. In his discussion of the producer–consumer synchronization algorithm, Lamport notes that “The most important class of properties one proves about an algorithm are invariance properties. A state predicate is an invariant if it is true in every state of every execution.”

Lamport notes that ‘Petri nets are a model of concurrent computation especially well-suited for expressing the need for arbitration. Although simple and elegant, Petri nets are not expressive enough to formally describe most interesting concurrent algorithms.’ He also mentioned that the first scientific examination of fault tolerance was Dijkstra’s 1974 seminal paper on self-stabilization [147], a work ahead of its time. Arguably, the most influential study of concurrency models was Milner’s Calculus of Communicating Systems (CCS) [348,349]. A number of formalisms based on the standard model were introduced for describing and reasoning about concurrent algorithms, including Amir Pnueli’s temporal logic introduced in 1977 [400].

Further readings. A fair number of textbooks discuss theoretical, as well as practical, aspects of concurrency. For example, [509] is dedicated to concurrency in transactional processing systems, and [397]

analyzes concurrency and consistency. The text [298] covers concurrent programming in Java, while [513] presents multithreading in C++.

The von Neumann architecture was introduced in [79]. The BSP and Multi-BSP models were introduced by Valiant in [485] and [486], respectively. Models of computations are discussed in [439]. PNs were introduced by Carl Adam Petri in [398]. An in-depth discussion of concurrency theory and system modeling with PNs can be found in [399]. The discussion of distributed systems leads to the observation that the analysis of communicating processes requires a more formal framework. Tony Hoare realized that a language based on execution traces is insufficient to abstract the behavior of communicating processes and developed *communicating sequential processes* (CSP) [241]. Milner initiated an axiomatic theory called the Calculus of Communicating System (CCS), [349]. Process algebra is the study of concurrent communicating processes within an algebraic framework. The process behavior is modeled as a set of equational axioms and a set of operators. This approach has its own limitations, the real-time behavior of the processes, the true concurrency, still escapes this axiomatization.

Seminal papers in distributed systems are authored by Manu Chandy and Leslie Lamport [93], by Leslie Lamport [291], [292], [293], Tony Hoare [241], and Robin Milner [349]. The collection of contributions with the title “Distributed Systems,” edited by Sape Mullender includes some of these papers. A survey of techniques and results related to the power of two random choices is presented in [354]. Seminal results on this subject are due to Azar [39], Karp [202,266], Mitzenmacher [353,355], and others.

10.17 Exercises and problems

Problem 1. Nonlinear algorithms do not obey the rules of scaled speed-up. For example, it was shown that, when the concurrency of an $\mathcal{O}(N^3)$ algorithm doubles, the problem size increases only by slightly more than 25%. Read [451] and explain this result.

Problem 2. Given a system of four concurrent threads t_1 , t_2 , t_3 , and t_4 we take a snapshot of the consistent state of the system after 3, 2, 4, and 3 events in each thread, respectively; all but the second event in each thread are local events. The only communication event in thread t_1 is to send a message to t_4 and the only communication event in thread t_3 is to send a message to t_2 . Draw a space–time diagram showing the consistent cut; mark individual events on thread t_i as e_i^j .

How many messages are exchanged to obtain the snapshot in this case? The snapshot protocol allows the application developers to create a checkpoint. An examination of the checkpoint data shows that an error has occurred, and it is decided to trace the execution. How many potential execution paths must be examined to debug the system?

Problem 3. The run-time of a data-intensive application could be days, or possibly weeks, even on a powerful supercomputer. Checkpoints are taken for a long-running computation periodically, and when a crash occurs, the computation is restarted from the latest checkpoint. This strategy is also useful for program and model debugging; when one observes wrong partial results, the computation can be restarted from a checkpoint where the partial results seem to be right. Express η , the *slowdown due to checkpointing*, for a computation in which checkpoints are taken after a run lasting τ units of time, and each checkpoint requires κ units of time. Discuss optimal choices for τ and κ . The checkpoint data can

be stored locally, on the secondary storage of each processor, or on a dedicated storage server accessible via a high-speed network. Which solution is optimal and why?

Problem 4. What is in your opinion the critical step in the development of a systematic approach to all-or-nothing atomicity? What does a systematic approach mean? What are the advantages of a systematic versus an ad hoc approach to atomicity? The support for atomicity affects the complexity of a system. Explain how the support for atomicity requires new functions/mechanisms and how these new functions increase the system complexity. At the same time, atomicity could simplify the description of a system; discuss how it accomplishes this. Support for atomicity is critical for system features that lead to increased performance and functionality, such as: virtual memory, processor virtualization, system calls, and user-provided exception handlers. Analyze how atomicity is used in each case.

Problem 5. The PN in Fig. 10.10(d) models a group of n concurrent processes in a shared-memory environment. At any given time, only one process may write, but any subset of the n processes may read at the same time, provided that no process writes. Identify the firing sequences, the markings of the net, the postsets of all transition, and the presets of all places. Can you construct a state machine to model the same process?

Problem 6*. Consider a computation consisting of n stages with a barrier synchronization among the N threads at the end of each stage. Assuming that you know the distribution of the random execution time of each thread for each stage, show how one could use order statistics [128] to estimate the completion time of the computation.

Problem 7. Consider the data flow graph of a computation C in Fig. 10.6. Call t_1, t_2, t_3 , and t_4 the time when the data inputs *data1*, *data2*, *data3*, and *data4* become available. Call T_i , $1 \leq i \leq 13$ the time required by computations C_i , $1 \leq i \leq 13$, respectively, to complete and express T the total time required by C to complete function of t_i and T_i .

Problem 8. Discuss the factors affecting parallel slackness including characteristics of the parallel computation, such as fine versus coarse grain, and the characteristics of the workload and of the computing substrate.

Problem 9*. In Section 10.14, we discussed the *power of two choices* for the balls and bins problem with n bins. Instead of placing each ball in one random bin, we choose two random bins for each ball, place it in the one that currently has fewest balls, and proceed in this manner sequentially for each ball. Prove Eq. (10.35).

Hint—the idea of the proof: Call B_i the number of bins with more than i balls at the end. We wish to find an upper bound, β_i for B_i . The probability that a ball is placed in bin q with at least $i+1$ balls in it is $Pr(N_q \geq i+1) \leq \left(\frac{\beta_1}{n}\right)^2$. Indeed, both choices of placing this ball must be in bins with at least i balls. The distribution of bins B_{i+1} is dominated by the binomial distribution $Bin\left(n, \left(\frac{\beta_1}{n}\right)^2\right)$. The mean of this distribution is $\left(\frac{\beta_1}{n}\right)^2$. According to Chernoff bound $\beta_{i+1} = c \left(\frac{\beta_1}{n}\right)^2$ with some constant c . Therefore, $\frac{\beta_1}{n}$ decreases quadratically, and the following holds $i \approx \frac{\ln \ln n}{\ln 2} \Rightarrow \beta_1 < 1$. It follows that the maximum number of balls in a bin is $\frac{\ln \ln n}{\ln 2}$ with high probability.

Problem 10. What is the difference between wait for graph and resource allocation graph?

Cloud applications

11

Over the years the users of large-scale computing systems in data centers discovered how difficult it was to develop efficient data-intensive and computationally intensive applications. It was equally difficult to locate systems best-suited to run an application, to determine when an application was able to run on these systems, and to estimate the time when the results could be expected. Porting an application from one system to another frequently constituted a challenging endeavor. Often, an application optimized for one system performed poorly on other systems.

There were also formidable challenges for the providers of service because system resources could not be effectively managed, and it was not feasible to provide QoS guarantees. Accommodating a dynamic load and supporting security and rapid recovery after a system-wide failure were daunting tasks due to the scale of the system. Any economic advantage offered by resource concentration was offset by the relatively low utilization of costly resources.

Cloud computing has changed the views and perceptions of users and providers of service on how to compute more efficiently and at a lower cost. Most of the challenges faced by application developers and service providers have either disappeared, or are significantly diminished. Cloud application developers are now able to work in a familiar environment and enjoy the advantages of a *just-in-time infrastructure*; they are free to design an application without being concerned where the application will run.

Cloud elasticity allows applications to seamlessly absorb additional workload. Cloud users also benefit from the speedup due to parallelization. When the workload can be partitioned in n segments, the application can spawn n instances of itself and run them concurrently, resulting in dramatic speedups. This is particularly useful for computer-aided design and for modeling complex systems when multiple design alternatives and multiple models of a physical system can be evaluated at the same time.

Cloud computing is focused on enterprise computing which clearly differentiates it from the grid computing effort largely focused on scientific and engineering applications. An important advantage of cloud computing over grid computing is that the resources offered by a cloud service provider are in one administrative domain.

Cloud computing is beneficial for the providers of computing cycles because it typically leads to more efficient resource utilization. It soon became obvious that a significant percentage of the typical workloads are dominated by frameworks such as MapReduce and that multiple frameworks must share the large computer clusters populating the cloud infrastructure.

The future success of cloud computing rests on the ability of companies promoting utility computing to convince an increasingly larger segment of the user population of the advantages of network-centric computing and network-centric content. This translates into the ability to provide satisfactory solutions to critical aspects of security, scalability, reliability, quality of service, and the requirements agreed upon in SLAs.

Sections 11.1, 11.2, and 11.3 cover cloud application developments and provide insights into workflow management. Coordination based on a state machine model is presented in Section 11.4, followed by the in-depth discussion of the MapReduce programming model and one of its applications in Sections 11.5 and 11.6. Hadoop, Yarn, and Tez are covered in Section 11.7, while systems such as Pig, Hive, and Impala are discussed in Section 11.8.

An overview of current cloud applications and new application opportunities is presented in Section 11.9, followed by a discussion of cloud applications in science and engineering and in biology research in Sections 11.10 and 11.11, respectively. Social computing and software fault isolations are the subjects of Sections 11.12 and 11.13.

11.1 Cloud application development and architectural styles

Web services, database services, and transaction-based services are ideal applications for cloud computing. The cost-performance profiles of such applications benefit from an elastic environment where resources are available when needed and where a cloud user pays only for the resources consumed by her application.

Not all applications are suitable for cloud computing. Applications where the workload cannot be arbitrarily partitioned, or require intensive communication among concurrent instances are unlikely to perform well on a cloud. An application with a complex workflow and multiple dependencies, as is often the case in high-performance computing, could require longer execution times and higher costs on a cloud. Benchmarks for high-performance computing, discussed in Section 11.10, show that communication and memory-intensive applications may not exhibit the performance levels shown when running on supercomputers with low-latency and high-bandwidth interconnects.

Cloud application development challenges. The development of efficient cloud applications faces challenges posed by the inherent imbalance between computing, I/O, and the communication bandwidth of processors. These challenges are greatly amplified due to the scale of the cloud infrastructure, its distributed nature, and by the very nature of data-intensive applications. Though cloud computing infrastructures attempt to automatically distribute and balance the workload, application developers are still left with the responsibility to identify optimal storage for the data, exploit spatial and temporal data and code locality, and minimize communication among running threads and instances.

A main attraction of cloud computing is the ability to use as many servers as necessary to optimally respond to the cost and the timing constraints of an application. This is only possible if the workload can be partitioned in segments of arbitrary size and can be processed in parallel by the servers available on the cloud. The *arbitrarily divisible load sharing model* describes workloads that can be partitioned into an arbitrarily large number of units and can be processed concurrently by multiple cloud instances. Applications most suitable for cloud computing enjoy this model.

The shared infrastructure, a defining quality of cloud computing, has side effects. Performance isolation, discussed in Section 5.2, is nearly impossible under real conditions, especially when a system is heavily loaded. The performance of virtual machines fluctuates based on the workload and the environment. Security isolation is challenging on multitenant systems.

Reliability of the cloud infrastructure is also a major concern. The frequent failures of servers built with off-the-shelf components is a consequence of the sheer number of servers and communication

systems. Choosing an optimal instance from those offered by the cloud infrastructure is another critical factor to be considered. Instances differ in terms of performance isolation, reliability, and security. Cost considerations also play a role in the choice of the instance type.

Many applications consist of multiple stages. In turn, each stage may involve multiple instances running in parallel on cloud servers and communicating among them. Thus, efficiency, consistency, and communication scalability are major concerns for an application developer. A cloud infrastructure exhibits latency and bandwidth fluctuations affecting the performance of all applications and, in particular, of data-intensive ones.

Data storage plays a critical role in the performance of any data-intensive application. The organization and the location of data storage, as well as the storage bandwidth, must be carefully analyzed to ensure optimal application performance. Clouds support many storage options to set up a file system similar to the Hadoop file system discussed in Section 11.7. Among them are off-instance cloud storage, e.g. S3, and mountable off-instance block storage, e.g., EBS, as well as storage persistent for the lifetime of the instance.

Many data-intensive applications use metadata associated with individual data records. For example, the metadata for an MPEG audio file may include the title of the song, the name of the singer, recording information, etc. Metadata should be stored for easy access; the storage should be scalable, and reliable.

Another important consideration for the application developer is logging. Performance considerations limit the amount of data logging, while the ability to identify the source of unexpected results and errors is helped by frequent logging. Logging is typically done using instance storage preserved only for the lifetime of the instance; thus measures to preserve the logs for a post-mortem analysis must be taken. Another challenge waiting resolution is related to software licensing, discussed in Section 2.11.

Cloud application architectural styles. Cloud computing is based on the client–server paradigm discussed in Section 3.12. The vast majority of cloud applications take advantage of request–response communication between clients and stateless servers. A *stateless server* does not require a client to first establish a connection to the server; instead, it views a client request as an independent transaction and responds to it.

The advantages of stateless servers are obvious. Recovering from a server failure requires a considerable overhead for a server that maintains the state of all its connections, but, in case of a stateless server, a client is not affected while a server goes down and then comes back up between two consecutive requests. A stateless system is simpler, more robust, and scalable; a client does not have to be concerned with the state of the server; if the client receives a response to a request, this means that the server is up and running; if not, it should resend the request later. A connection-based service must reserve space to maintain the state of each connection with a client; therefore, such a system is not scalable, and the number of clients a server could interact with at any given time is limited by the storage space available to the server.

The Hypertext Transfer Protocol used by a browser to communicate with a web server is a request–response application protocol. HTTP uses TCP, a connection-oriented and reliable transport protocol. While TCP ensures reliable delivery of large objects, it exposes web servers to denial of service attacks. In such attacks, malicious clients fake attempts to establish a TCP connection and force the server to allocate space for the connection. A basic web server is stateless; it responds to an HTTP request without maintaining a history of past interactions with the client. The client, a browser, is also stateless since it sends requests and waits for responses.

A critical aspect of the development of networked applications is how processes and threads running on systems with different architectures, possibly compiled from different programming languages, can *communicate structured information with one another*. First, the internal representation of the two structures at the two sites may be different; one system may use *Big-Endian* and the other *Little-Endian* representation. The character representations also may be different. A communication channel transmits a sequence of bits/bytes and the data structure must be serialized at the sending site and reconstructed at the receiving site.

Several considerations including neutrality, extensibility, and independence, must be analyzed before choosing the architectural style of an application. *Neutrality* refers to application-level protocol ability to use different transport protocols such as TCP and UDP and to run on top of a different protocol stack. For example, SOAP can use not only TCP but also UDP, SMTP, and Java Message Service as transport vehicles. *Extensibility* refers to the ability to incorporate additional functions such as security. *Independence* refers to the ability to accommodate different programming styles.

Very often application clients and the servers running on the cloud communicate using RPCs, discussed in Section 3.12, yet other communication styles are possible. RPC-based applications use *stubs* to convert the parameters involved in an RPC call. A stub performs two functions, marshaling the data structures and serialization.

A more general concept is that of an Object Request Broker (ORB), middleware facilitating communication of networked applications. The sending site ORB transforms data structures used internally and transmits a byte sequence over the network. The receiving site ORB maps the byte sequence to data structures used internally by the receiving process.

Common Object Request Broker Architecture (CORBA) was developed in early 1990s to enable networked applications developed using different programming languages and running on systems with different architecture and system software to work with one another. At the heart of the system is the Interface Definition Language (IDL) used to specify the interface of an object. An IDL representation is then mapped to the set of programming languages including: C, C++, Java, Smalltalk, Ruby, Lisp, and Python. Networked applications pass CORBA by reference and pass data by value.

Simple Object Access Protocol (SOAP) was developed in 1998 for web applications. SOAP message format is based on the Extensible Markup Language (XML). SOAP uses TCP and more recently, UDP transport protocols, but it can also be stacked above other application layer protocols such as HTTP, SMTP, or JMS. The processing model of SOAP is based on a network consisting of senders, receivers, intermediaries, message originators, ultimate receivers, and message paths. SOAP is an underlying layer of Web Services.

Web Services Description Language (WSDL) was introduced in 2001 as an XML-based grammar to describe communication between end points of a networked application. The abstract definitions of the elements involved include: *services*, collection of endpoints of communication; *types*, containers for data type definitions; *operations*, description of actions supported by a service; *port types*, operations supported by endpoints; *bindings*, protocols and data format supported by a particular port type; and *port*, an endpoint as a combination of a binding and a network address. These abstractions are mapped to concrete message formats and network protocols to define endpoints and services.

Representational State Transfer (REST) is a software architecture for distributed hypermedia systems supporting client communication with stateless servers. It is platform- and language-independent, supports data caching, and can be used in the presence of firewalls. REST almost always uses HTTP

for all four CRUD (*Create/Read/Update/Delete*) operations; it uses *GET*, *PUT*, and *DELETE* to read, write, and delete the data, respectively.

REST is an easier to use alternative to RPC, CORBA, or web services such as SOAP or WSDL. For example, to retrieve the address of an individual from a database, a REST system sends an URL specifying the network address of the database, the name of the individual, and the specific attribute in the record that the client application wants to retrieve, in this case the address. The corresponding SOAP version of such a request consists of ten lines or more of XML. The REST server responds with the address of the individual. This justifies the statement that REST is a lightweight protocol. As far as usability is concerned, REST is easier to build from scratch and debug, but SOAP is supported by tools that use self-documentation, e.g., WSDL to generate the code to connect.

11.2 Coordination of multiple activities

Many applications require the completion of multiple interdependent tasks [529]. The description of a complex activity involving such an ensemble of tasks is known as a *workflow*. In this section, we discuss workflow models, the lifecycle of a workflow, and the desirable properties of a workflow description. Workflow patterns, reachability of the goal state of a workflow, dynamic workflows, and a parallel between traditional transaction systems and cloud workflows are covered in Section 11.3.

Basic concepts. Workflow models are abstractions revealing the most important properties of the entities participating in a workflow management system. *Task* is the central concept in workflow modeling. A task is a unit of work to be performed on the cloud, and it is characterized by several attributes, such as: (i) name—a string of characters uniquely identifying the task; (ii) description—a natural language description of the task; (iii) actions—an action is a modification of the environment caused by the execution of the task; (iv) preconditions—Boolean expressions that must be true before the action(s) of the task can take place; (v) postconditions—Boolean expressions that must be true after the action(s) of the task do take place; (vi) attributes—provide indications of the type and quantity of resources necessary for the execution of the task, the actors in charge of the tasks, the security requirements, whether the task is reversible or not, and other task characteristics; and (vii) exceptions—provide information on how to handle abnormal events. The exceptions supported by a task consist of a list of `<event, action>` pairs. The exceptions included in the task exception list are called *anticipated exceptions*, as opposed to unanticipated exceptions. Events not included in the exception list trigger replanning. *Replanning* means restructuring of a process, redefinition of the relationship among various tasks.

A *composite task* is a structure describing a subset of tasks and the order of their execution. A *primitive task* is one that cannot be decomposed into simpler tasks. A composite task inherits some properties from workflows; it consists of tasks, has one start symbol, and possibly several end symbols. At the same time, a composite task inherits some properties from tasks; it has a name, preconditions, and postconditions.

A *routing task* is a special-purpose task connecting two tasks in a workflow description. The task that has just completed execution is called the *predecessor task*, the one to be initiated next is called the *successor task*. A routing task could trigger the sequential, concurrent, or iterative execution. Two types of routing tasks exist:

- A *fork routing task* triggers execution of several successor tasks. Several semantics for this construct are possible: (i) all successor tasks are enabled; (ii) each successor task is associated with a condition, the conditions for all tasks are evaluated and only the tasks with a `true` condition are enabled; (iii) each successor task is associated with a condition, the conditions for all tasks are evaluated, but the conditions are mutually exclusive, and only one condition may be `true`, thus only one task is enabled; and (iv) nondeterministic— k out of $n > k$ successors are selected at random to be enabled.
- A *join routing task* waits for completion of its predecessor tasks. There are several semantics for the join routing task: (i) the successor is enabled after all predecessors end; (ii) the successor is enabled after k out of $n > k$ predecessors end; and (iii) iterative—the tasks between the fork and the join are executed repeatedly.

Process descriptions and cases. A *process description*, also called a workflow schema, is a structure describing the *tasks* or *activities* to be executed and the order of their execution; a process description contains one start symbol and one end symbol. A process description can be provided in a Workflow Definition Language (WFDL), supporting constructs for choice, concurrent execution, the classical *fork*, *join* constructs, and iterative execution. Workflow description resembles a *flowchart*, a concept we are familiar with from programming.

The phases in the life-cycle of a workflow are creation, definition, verification, and enactment. There is a striking similarity between the life-cycle of a workflow and that of a traditional computer program, namely, creation, compilation, and execution; see Fig. 11.1. The workflow specification by means of a workflow description language is analogous to writing a program. Planning is equivalent to automatic program generation. Workflow verification corresponds to syntactic verification of a program, and workflow enactment mirrors the execution of a compiled program.

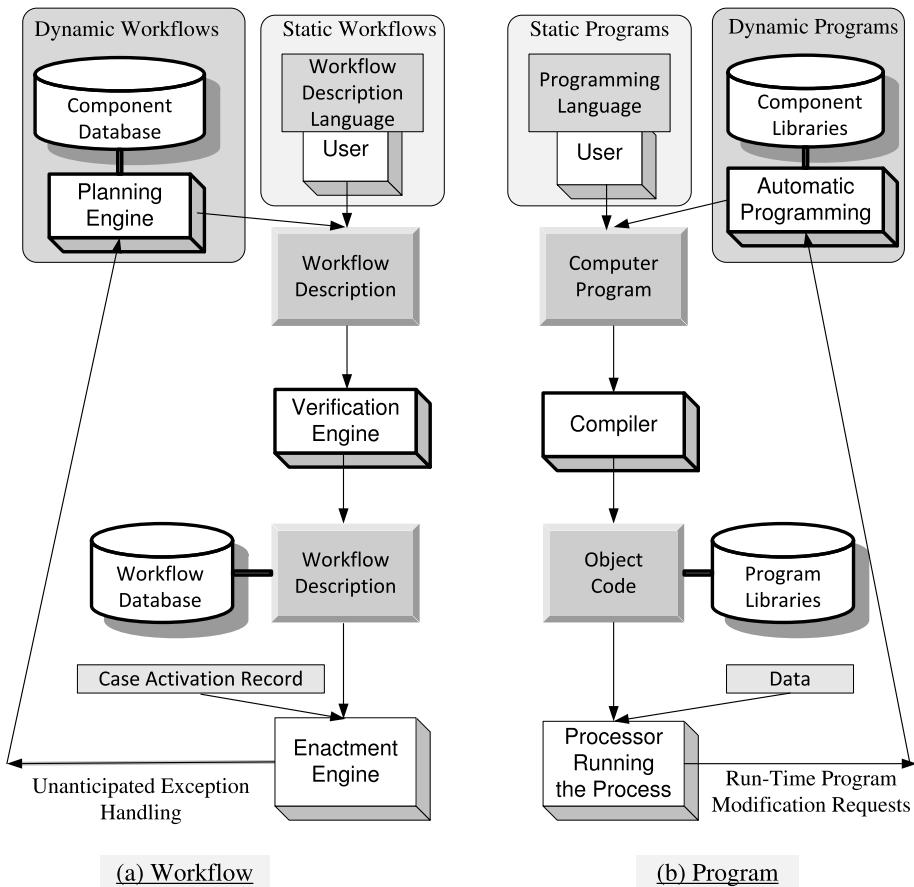
A *case* is an instance of a process description. The start and stop symbols in the workflow description enable the creation and the termination of a case. An *enactment model* describes the steps taken to process a case. When all tasks required by a workflow are executed by a computer, the enactment can be performed by a program called an *enactment engine*.

The *state of a case* at time t is defined in terms of tasks already completed at that time. Events cause transitions between states. Identifying the states of a case consisting of concurrent activities is considerably more difficult than identifying the states of a strictly sequential process. Indeed, when several activities could proceed concurrently, the state has to reflect the progress made on each independent activity.

An alternative description of a workflow can be provided by a transition system describing the possible paths from the current state to a goal state. Sometimes, instead of providing a process description, we may specify only the goal state and expect the system to generate a workflow description that could lead to that state through a set of actions. In this case, the new workflow description is generated automatically, knowing a set of tasks and the preconditions and postconditions for each one of them. In AI, this activity is known as *planning*.

The state space of a process includes one initial state and one goal state. A transition system identifies all possible paths from the initial to the goal state. A case corresponds to a particular path in the transition system. The state of a case tracks the progress made during the enactment of that case.

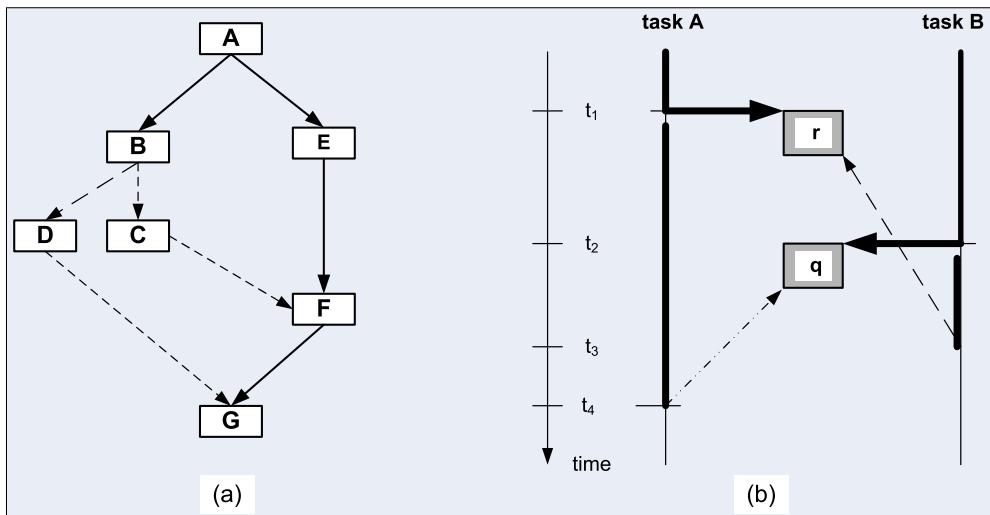
Safety and liveness are the most desirable properties of a process description. Informally, *safety* means that nothing “bad” ever happens, while *liveness* means that something “good” will eventually take place, should a case based on the process be enacted. Not all processes are safe and live. For example, the process description in Fig. 11.2(a) violates the liveness requirement. As long as task C is

**FIGURE 11.1**

Workflows and programs. (a) The life-cycle of a workflow. (b) The life-cycle of a computer program. The workflow definition is analogous to writing a program. Planning is analogous to automatic program generation. Verification corresponds to syntactic verification of a program. Workflow enactment mirrors the execution of a program. A static workflow corresponds to a static program and a dynamic workflow to a dynamic program.

chosen after completion of *B*, the process will terminate. However, if *D* is chosen, then *F* will never be instantiated because it requires the completion of both *C* and *E*. The process will never terminate because *G* requires completion of both *D* and *F*.

A process description language should be unambiguous and should allow a verification of the process description before the enactment of a case. It is entirely possible that a process description may be enacted correctly in some cases but could fail for others. Such enactment failures may be very costly and should be prevented by a thorough verification at the process definition time. To avoid enactment

**FIGURE 11.2**

- (a) A process description that violates the liveness requirement; if task *C* is chosen after completion of *B*, the process will terminate after executing task *G*; if *D* is chosen, then *F* will never be instantiated because it requires the completion of both *C* and *E*. The process will never terminate because *G* requires completion of both *D* and *F*.
- (b) Tasks *A* and *B* need exclusive access to two resources *r* and *q*, and a deadlock may occur if the following sequence of events occur: at time t_1 , task *A* acquires *r*, at time t_2 , task *B* acquires *q* and continues to run; then, at time t_3 , task *B* attempts to acquire *r*, and it blocks because *r* is under the control of *A*; task *A* continues to run at time t_4 , attempts to acquire *q*, and blocks because *q* is under the control of *B*.

errors, we need to verify process description and check for desirable properties such as safety and liveness. Some process description methods are more suitable for verification than others.

A note of caution: Although the original description of a process could be live, the actual enactment of a case may be affected by deadlocks due to resource allocation. To illustrate this situation, consider two tasks, *A* and *B*, running concurrently; each of them needs exclusive access to resources *r* and *q* for a period of time. Two scenarios are possible:

- (1) either *A* or *B* acquires both resources and then releases them, and allows the other task to do the same;
- (2) we face the undesirable situation in Fig. 11.2(b) when at time t_1 task *A* acquires *r* and continues its execution; then, at time t_2 , task *B* acquires *q* and continues to run. Then, at time t_3 , task *B* attempts to acquire *r*, and it blocks because *r* is under the control of *A*. Task *A* continues to run, and at time t_4 attempts to acquire *q* and it blocks because *q* is under the control of *B*.

The deadlock illustrated in Fig. 11.2(b) can be avoided by requesting each task to acquire all resources at the same time; the price to pay is underutilization of resources; indeed, the idle time of each resource increases under this scheme.

11.3 Workflow patterns

The term *workflow pattern* refers to the temporal relationships among the tasks of a process. The workflow description languages and the mechanisms to control the enactment of a case must have provisions to support these temporal relationships. Workflow patterns are analyzed in [1], [325], and [531]. These patterns are classified in several categories: basic, advanced branching and synchronization, structural, state-based, cancelation, and patterns involving multiple instances. The basic workflow patterns in Fig. 11.3 are:

Sequence—occurs when several tasks have to be scheduled one after the completion of the other, Fig. 11.3(a).

AND split—requires several tasks to be executed concurrently. Both tasks B and C are activated when task A terminates; see Fig. 11.3(b). In case of an *explicit AND split*, the activity graph has a routing node, and all activities connected to the routing node are activated as soon as the flow of control reaches the routing node. In the case of an *implicit AND split*, activities are connected directly and conditions can be associated with branches linking an activity with the next ones. The tasks are activated only when the conditions associated with a branch are true.

Synchronization—an activity can start only after several concurrent activities finish execution; task C can start only after both tasks A and B terminate, Fig. 11.3(c).

XOR split—requires a decision; after the completion of task A, either B or C can be activated, Fig. 11.3(d).

XOR join—several alternatives are merged into one; task C is enabled when either A or B terminates; see Fig. 11.3(e).

OR split—choose multiple alternatives out of a set; after completion of task A, one could activate either B or C, or both; see Fig. 11.3(f).

Multiple merge—allows multiple activations of a task and does not require synchronization after the execution of concurrent tasks; once A terminates, tasks B and C execute concurrently; see Fig. 11.3(g). When the first of them, say B, terminates, then task D is activated; then, when C terminates, D is activated again.

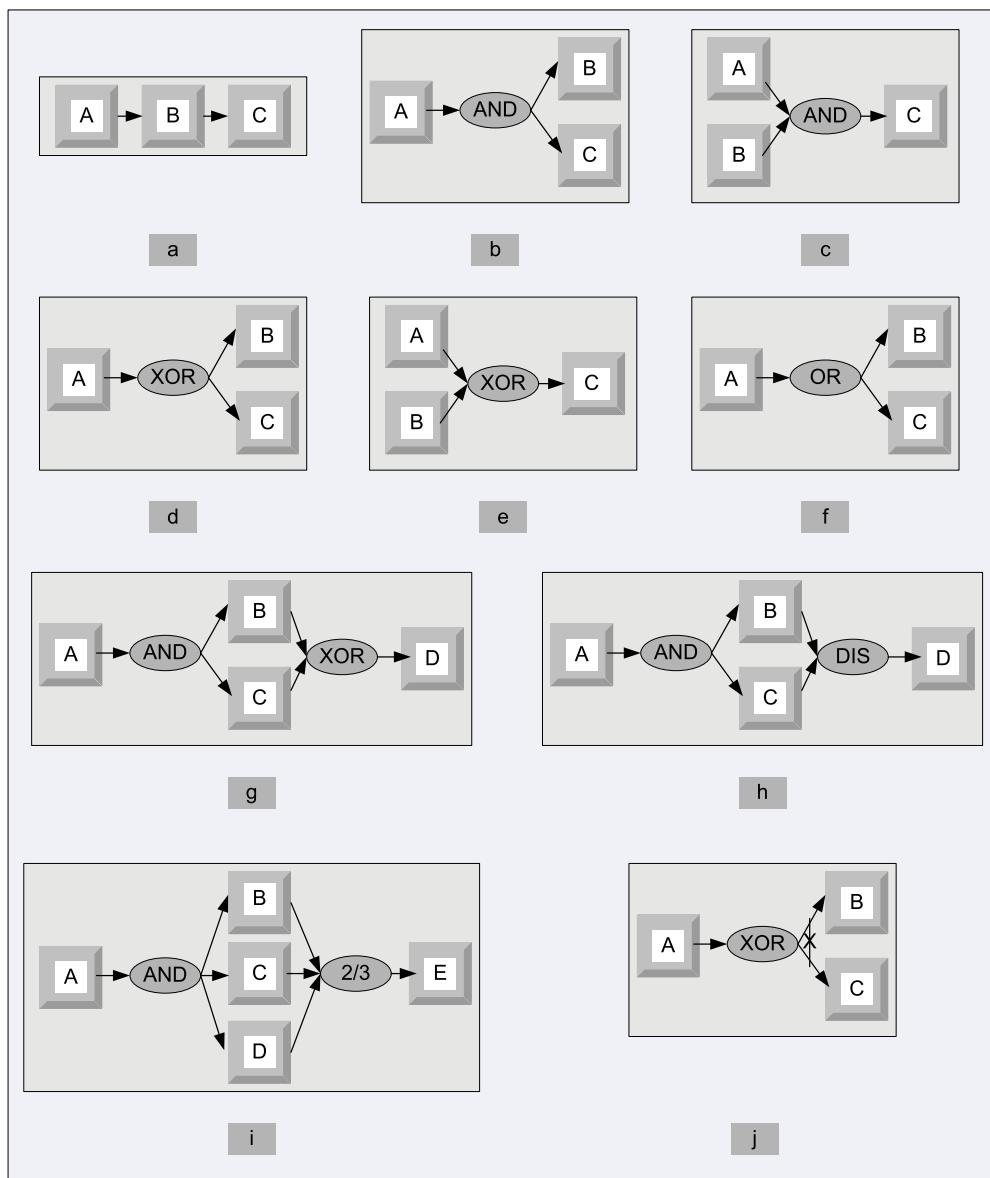
Discriminator—waits for a number of incoming branches to complete before activating the subsequent activity; see Fig. 11.3(h); then, it waits for the remaining branches to finish without taking any action until all of them have terminated; next, it resets itself.

N out of M join—provides a barrier synchronization; assuming $M > N$ tasks run concurrently, N of them have to reach the barrier before the next task is enabled; in our example, any two out of the three tasks A, B, and C have to finish before E is enabled; see Fig. 11.3(i).

Deferred choice—similar to XOR split but the choice is not made explicitly and the run-time environment decides what branch to take; see Fig. 11.3(j).

Next, we discuss the reachability of the goal state, and we consider the following elements:

- A system Σ , an initial state of the system, $\sigma_{initial}$, and a goal state, σ_{goal} .
- A process group, $\mathcal{P} = \{p_1, p_2, \dots, p_n\}$; each process p_i in the process group is characterized by a set of preconditions, $pre(p_i)$, postconditions, $post(p_i)$, and attributes, $atr(p_i)$.
- A workflow described by a directed activity graph \mathcal{A} or by a procedure Π capable to construct \mathcal{A} given the tuple $\langle \mathcal{P}, \sigma_{initial}, \sigma_{goal} \rangle$. The nodes of \mathcal{A} are processes in \mathcal{P} , and the edges define precedence relationships among processes. $P_i \rightarrow P_j$ implies that $pre(p_j) \subset post(p_i)$.

**FIGURE 11.3**

Basic workflow patterns. (a) Sequence; (b) AND split; (c) Synchronization; (d) XOR split; (e) XOR merge; (f) OR split; (g) Multiple Merge; (h) Discriminator; (i) N out of M join; (j) Deferred Choice.

- A set of constraints, $\mathcal{C} = \{C_1, C_2, \dots, C_m\}$.

The coordination problem for system Σ in state $\sigma_{initial}$ is to reach state σ_{goal} , as a result of postconditions of some process $P_{final} \in \mathcal{P}$ subject to constraints $C_i \in \mathcal{C}$. Here, $\sigma_{initial}$ enables the preconditions of some process $P_{initial} \in \mathcal{P}$. Informally, this means that a chain of processes exists such that the postconditions of one process are preconditions of the next process in the chain.

The *preconditions* of a process are either the conditions and/or the events that trigger the execution of the process or the data the process expects as input. The *postconditions* are the results produced by that process. The attributes of a process describe special requirements or properties of the process.

Some workflows are static, and the activity graph does not change during the enactment of a case. *Dynamic workflows* are those that allow the activity graph to be modified during the enactment of a case. Some of the more difficult questions encountered in dynamic workflow management refer to: (i) how to integrate workflow and resource management and guarantee optimality or near optimality of cost functions for individual cases; (ii) how to guarantee consistency after a change in a workflow; and (iii) how to create a dynamic workflow. Static workflows can be described in WFDL (the workflow definition language), but dynamic workflows need a more flexible approach.

We distinguish two basic models for the mechanics of workflow enactment: (i) *strong coordination models* in which the process group \mathcal{P} executes under the supervision of a *coordinator* process and a coordinator process acts as an enactment engine and ensures a seamless transition from one process to another in the activity graph; and (ii) *Weak coordination models* in which there is no supervisory process.

In the first case, we may deploy a *hierarchical coordination scheme* with several levels of coordinators. A supervisor at level i in a hierarchical scheme with $i + 1$ levels coordinates a subset of processes in the process group. A supervisor at level $i - 1$ coordinates a number of supervisors at level i , and the root provides global coordination. Such a hierarchical coordination scheme may be used to reduce the communication overhead; a coordinator and the processes it supervises may be colocated.

An important feature of this coordination model is its ability of supporting dynamic workflows. The coordinator or the global coordinator may respond to a request to modify the workflow by first stopping all the threads of control in a consistent state, then investigating the feasibility of the requested changes, and finally implementing feasible changes.

Weak coordination models are based on peer-to-peer communication between processes in the process group by means of a societal service, such as a *tuple space*. Once a process $p_i \in \mathcal{P}$ finishes, it deposits a token including possibly a subset of its postconditions, $post(p_i)$, in a tuple space. The consumer process p_j is expected to visit at some point in time the tuple space, examine the tokens left by its ancestors in the activity graph, and, if its preconditions $pre(p_j)$ are satisfied, commence the execution. This approach requires individual processes to either have a copy of the activity graph or some timetable to visit the tuple space. An alternative approach is using an *active space*, a tuple space augmented with the ability to generate an event awakening the consumer of a token.

There are similarities and some differences between workflows of traditional transaction-oriented systems and cloud workflows; the similarities are mostly at the modeling level, whereas the differences affect the mechanisms used to implement workflow management systems. Some of the more subtle differences between them are:

- The emphasis in a transactional model is placed on the contractual aspect of a transaction; in a workflow the enactment of a case is sometimes based on a “best-effort” model where the agents involved will do their best to attain the goal state, but there is no guarantee of success.

- A critical aspect of the transactional model in database applications is maintaining a consistent state of the database; a cloud is an open system; thus its state is considerably more difficult to define.
- Database transactions are typically short-lived; the tasks of a cloud workflow could be long lasting.
- A database transaction consists of a set of well-defined actions that are unlikely to be altered during the execution of the transaction. However, the process description of a cloud workflow may change during the lifetime of a case.
- The individual tasks of a cloud workflow may not exhibit the traditional properties of database transactions. For example, consider durability, which means that, at any instance of time before reaching the goal state, a workflow may roll back to some previously encountered state and continue from there on an entirely different path. A task of a workflow could be either reversible or irreversible. Sometimes, paying a penalty for reversing an action is more profitable in the long run than continuing on a wrong path.
- Resource allocation is a critical aspect of the workflow enactment on a cloud without an immediate correspondent for database transactions.

The relatively simple coordination model discussed next is often used in cloud computing.

11.4 Coordination based on a state machine model—zookeeper

Cloud computing elasticity requires the ability to distribute computations and data across multiple systems. In a distributed computing environment coordination among these systems is a critical function. The coordination model depends on the specific task, e.g., coordination of data storage, orchestration of multiple activities, blocking an activity until an event occurs, reaching consensus on the next action, or recovery after an error.

The entities to be coordinated could be processes running on a set of cloud servers or even running on multiple clouds. Servers running critical task are often replicated; when one primary server fails, a backup automatically continues the execution. This is only possible if the backup is in a *hot standby* mode, in other words, if the primary server shares its state with the backup at all times.

For example, in the distributed data store model discussed in Section 2.6, the access to data is mitigated by a proxy. An architecture with multiple proxies is desirable, as a proxy is a single point of failure. All proxies should be in the same state, so, whenever one of them fails, the client could seamlessly continue to access the data from another proxy.

Consider now an advertising service involving a large number of cloud servers. The advertising service runs on a number of servers specialized for tasks such as: database access, monitoring, accounting, event logging, installers, customer dashboards,¹ advertising campaign planners, scenario testing, etc.

These activities can be coordinated using a configuration file shared by all systems. When the service starts, or after a system failure, all servers use the configuration file to coordinate their actions. This solution is static; any change requires an update and re-distribution of the configuration file. Moreover, in case of a system failure the configuration file does not allow recovery from the state each server was in prior to the system crash.

¹ A customer dashboard provides access to key customer information, such as contact name and account number, in an area of the screen that remains persistent as the user navigates through multiple web pages.

A solution for the coordination problem is to implement a proxy as a deterministic finite-state machine transitioning from one state to the next in response to client commands. When P proxies are involved, all must be synchronized and execute the same sequence of state transitions upon receiving client commands. This scenario can be ensured when all proxies implement a version of the Paxos consensus algorithm described in Section 10.13.

Zookeeper is a distributed coordination service based on this model. The high throughput and low latency service is used for coordination in large-scale distributed systems. The open-source software is written in Java and has bindings for Java and C. The information about the project is available at <http://zookeeper.apache.org/>.

Zookeeper software must first be downloaded and installed on several servers. Then clients can connect to any server and access the coordination service. The service is available as long as the majority of servers in the pack are available.

The servers in the pack communicate with one another and elect a *leader*. A database is replicated on each of them, and the consistency of the replicas is maintained. Fig. 11.4(a) shows that the service provides a single system image; a client can connect to any server of the pack. A client uses TCP to connect to one server, then sends requests, receives responses, and watches events. A client synchronizes its clock with the server it is connected to. When a server fails, the TCP connections of all clients connected to it time-out, and the clients detect the failure and connect to other servers.

Figs. 11.4(b) and (c) show that a READ operation directed to any server in the pack returns the same result, while processing of a WRITE operation is more involved. The servers elect a *leader*, and any *follower* server forwards to the leader requests from the clients connected to it. The leader uses atomic broadcast to reach consensus. When the leader fails, the servers elect a new leader.

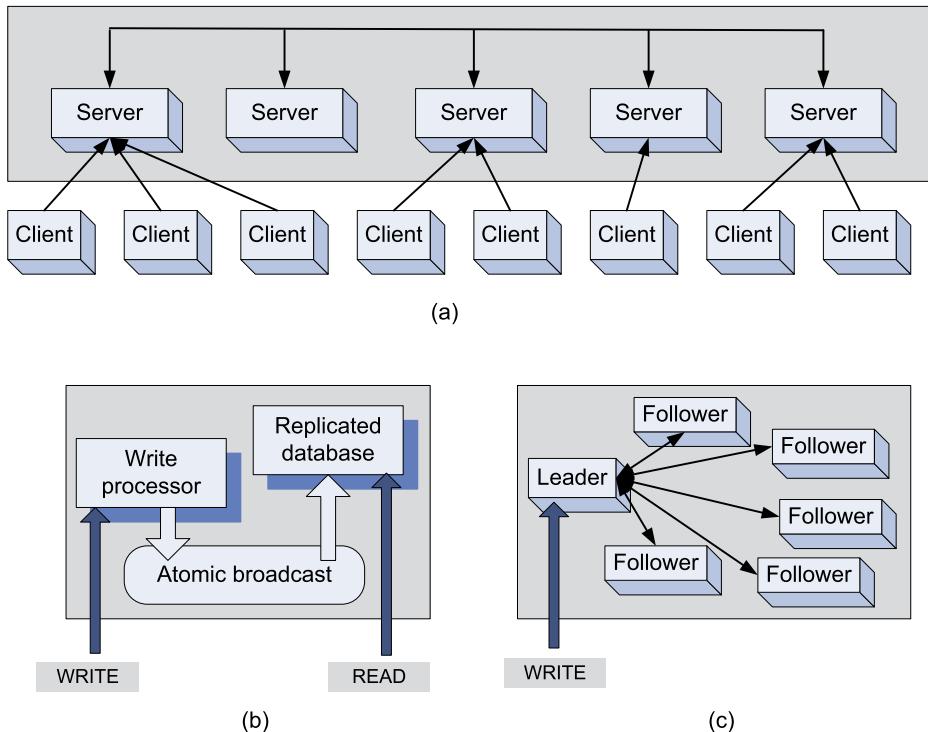
The system is organized as a shared hierarchical namespace similar to the organization of a file system. A name is a sequence of path elements separated by a backslash. Every name in Zookeeper's name space is identified by a unique path; see Fig. 11.5.

A *znode* of Zookeeper, the equivalent of an *inode* of UFS, may have data associated with it. The system is designed to store state information; the data in each node includes version numbers for the data, changes of ACLs,² and time stamps. A client can set a watch on a *znode* and receive a notification when the *znode* changes. This organization allows coordinated updates. The data retrieved by a client contains also a version number. Each update is stamped with a number that reflects the order of the transition.

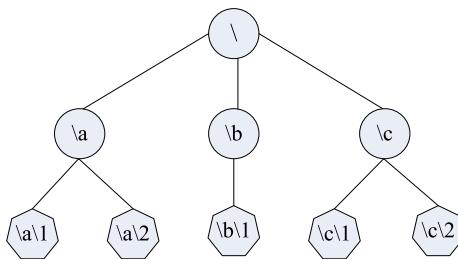
Data stored in each node is read and written atomically, and a READ returns all data stored in a *znode*, while a WRITE replaces all data in the *znode*. Unlike a file system, the Zookeeper data, the image of the state, is stored in the server memory. Updates are logged to disk for recoverability, and WRITEs are serialized to disk before they are applied to the in-memory database containing the entire tree. Zookeeper service guarantees:

1. Atomicity—a transaction either completes or fails.
2. Updates sequential consistency—updates are applied strictly in the order in which they are received.
3. Single system image for the clients—a client receives the same response regardless of the server it connects to.

² An Access Control List (ACL) is a list of pairs (subject, value) that define the list of access rights to an object; for example, read, write, execute permissions for a file.

**FIGURE 11.4**

Zookeeper coordination service. (a) The service provides a single system image; clients can connect to any server in the pack. (b) The functional model of Zookeeper service; the replicated database is accessed directly by READ commands. (c) Processing a WRITE command: (i) a server receiving a command from a client, forwards it to the leader; (ii) the leader uses atomic broadcast to reach consensus among all followers.

**FIGURE 11.5**

The Zookeeper is organized as a shared hierarchical namespace; a name is a sequence of path elements separated by a backslash.

4. Update persistence—once applied, an update persists until it is overwritten by a client.
5. Reliability—the system is guaranteed to function correctly as long as the majority of servers function correctly.

READ requests are serviced from the local replica of the server connected to the client to reduce the response time. When the leader receives a WRITE request, it determines the state of the system where the WRITE will be applied, and then it transforms the state into a transaction capturing the new state.

The messaging layer is responsible for the election of a new leader when the current leader fails. The messaging protocol uses: *packets*—sequence of bytes sent through a FIFO channel, *proposals*—units of agreement, and *messages*—sequence of bytes atomically broadcast to all servers. A message is included into a proposal, and it is agreed upon before it is delivered. Proposals are agreed upon by exchanging packets with a quorum of servers as required by the Paxos algorithm.

An atomic messaging system keeps all of the servers in pack in sync. The messaging system guarantees: (a) Reliable delivery: if a message m is delivered to one server, it will be eventually delivered to all servers; (b) Total order: if message m is delivered before message n to one server, it will be delivered before n to all servers; and (c) Causal order: if message n is sent after m has been delivered by the sender of n , then m must be ordered before n .

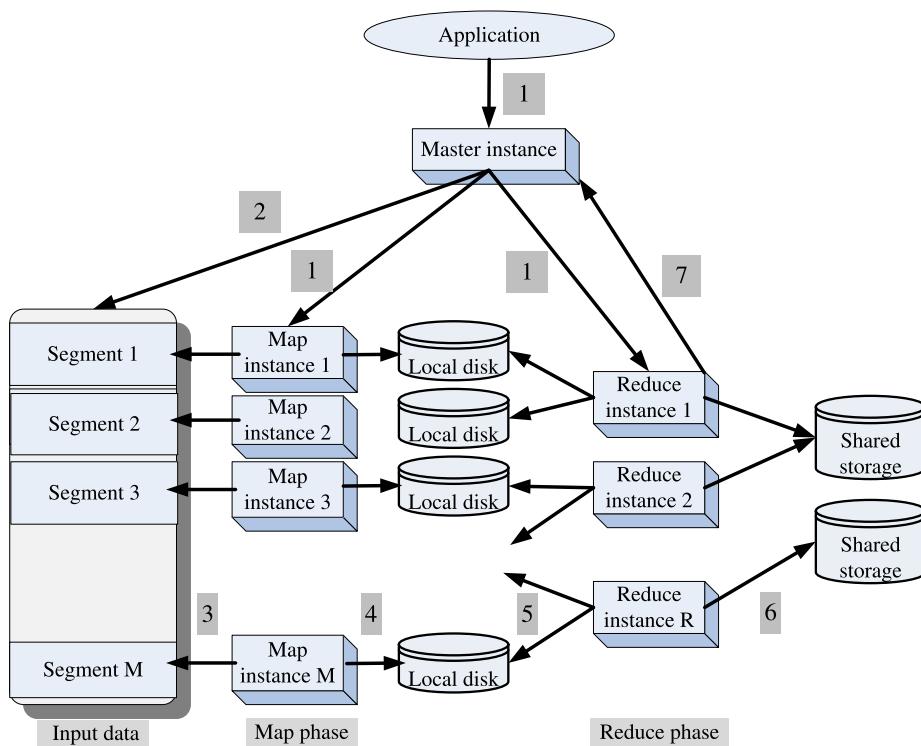
Zookeeper Application Programming Interface (API) consists of seven operations: *create*—add a node at a given location on the tree; *delete*—delete a node; *get data*—read data from a node; *set data*—write data to a node; *get children*—retrieve a list of the children of the node; and *sync*—wait for the data to propagate. The system also supports the creation of *ephemeral* nodes, nodes created when a session starts and deleted when the session ends.

This brief description shows that the Zookeeper service supports the finite-state machine model of coordination where a *znode* stores the state. The Zookeeper service can be used to implement higher-level operations such as group membership, synchronization, and so on. Yahoo’s Message Broker and many other applications use the Zookeeper service.

11.5 MapReduce programming model

A main advantage of cloud computing is elasticity, the ability to use as many servers as necessary to optimally respond to the cost and the timing constraints of an application. Typically, a front end of a transaction processing system distributes incoming transactions to a number of back-end systems and attempts to balance the workload. As the workload increases, new back-end systems are added to the pool. Many realistic applications in physics, biology, and other areas of computational science and engineering obey the arbitrarily divisible load sharing model; the workload can be partitioned into an arbitrarily large number of smaller workloads of equal, or nearly so, size. Yet, partitioning the workload of data-intensive applications is not always trivial.

MapReduce. MapReduce is based on a very simple idea for parallel processing of data-intensive applications supporting arbitrarily divisible load sharing; see Fig. 11.6. First, split the data into blocks, then assign each block to an instance/process and run these instances in parallel. Once all the instances have finished the computations assigned to them, start the second phase and merge the partial results produced by individual instances. The Same Program Multiple Data (SPMD) paradigm, used since the early days of parallel computing, is based on the same idea, but it assumes that a *Master* instance partitions the data and gathers the partial results.

**FIGURE 11.6**

MapReduce philosophy. 1. An application starts a Master instance and M worker instances for the Map phase and later R worker instances for the Reduce phase. 2. The Master partitions the input data in M segments. 3. Each Map instance reads its input data segment and processes the data. 4. The results of the processing are stored on the local disks of the servers where the Map instances run. 5. When all Map instances have finished processing, their data is read by R Reduce instances. 6. The final results are written by Reduce instances to a shared storage server. 7. The Master instance monitors the Reduce instances, and when all of them report task completion, the application is terminated.

MapReduce is a programming model inspired by the *map* and the *reduce* primitives of the Lisp programming language. It was conceived for processing and generating large data sets on computing clusters [129]. As a result of the computation, a set of input $\langle key, value \rangle$ pairs is transformed into a set of output $\langle key, value \rangle$ pairs.

Numerous applications can be easily implemented using this model. For example, one can process logs of web page requests and count the URL access frequency; the Map and Reduce functions produce the pairs $\langle URL, 1 \rangle$ and $\langle URL, totalcount \rangle$, respectively. Another trivial example is *distributed sort* when the map function extracts the key from each record and produces a $\langle key, record \rangle$ pair, and the Reduce function outputs these pairs unchanged. The following example [129] shows the two

user-defined functions for an application that counts the number of occurrences of each word in a set of documents.

```
map(String key, String value):
    // key: document name;  value: document contents
    for each word w in value:
        EmitIntermediate(w, "1");

reduce(String key, Iterator values):
    // key: a word;  values: a list of counts
    int result = 0;
    for each v in values:
        result += ParseInt(v);
    Emit(AsString(result));
```

Call M and R the number of Map and Reduce tasks, respectively, and N the number of systems used by the MapReduce. When a user program invokes the MapReduce function, the following sequence of actions take place:

- The run-time library splits the input files into M *splits* of 16 to 64 MB each, identifies a number N of systems to run, and starts multiple copies of the program, one of the systems being a Master and the others Workers. The Master assigns to each idle system either a *map* or a *reduce* task. The Master makes $\mathcal{O}(M + R)$ scheduling decisions and keeps $\mathcal{O}(M \times R)$ worker state vectors in memory. These considerations limit the size of M and R ; at the same time, efficiency considerations require that $M, R \gg N$.
- A Worker being assigned a Map task reads the corresponding input split, parses $\langle key, value \rangle$ pairs, and passes each pair to a user-defined Map function. The intermediate $\langle key, value \rangle$ pairs produced by the Map function are buffered in memory before being written to a local disk, partitioned into R regions by the partitioning function.
- The locations of these buffered pairs on the local disk are passed back to the Master, which is responsible for forwarding these locations to the Reduce Workers. A Reduce Worker uses remote procedure calls to read the buffered data from the local disks of the Map Workers; after reading all the intermediate data, it sorts it by the intermediate keys. For each unique intermediate key, the key and the corresponding set of intermediate values are passed to a user-defined Reduce function. The output of the Reduce function is appended to a final output file.
- The Master wakes up the user program when all Map and Reduce tasks finish.

The system is fault-tolerant; for each Map and Reduce task, the Master stores the state (idle, in-progress, or completed) and the identity of the worker machine. The Master pings every worker periodically and marks the worker as failed if it does not respond; a task in progress on a failed worker is reset to idle and becomes eligible for rescheduling. The Master writes periodic checkpoints of its control data structures, and, if the task fails, it can be restarted from the last checkpoint. The data is stored using GFS, the Google File System discussed in Section 7.6.

An environment for experimenting with MapReduce is described in [129]: the computers are typically dual-processor x86 processors running Linux, with 2–4 GB of memory per machine and commodity networking hardware typically 100–1 000 Mbps. A cluster consists of hundreds or thousands

of machines. Data is stored on IDE³ disks attached directly to individual machines. The file system uses replication to provide availability and reliability with unreliable hardware. To minimize network bandwidth, the input data is stored on the local disks of each system.

MapReduce with FlumeJava. The Java library discussed in Section 10.15 supports a new operation called MapShuffleCombineReduce. This operation combines *ParallelDo*, *GroupByKey*, *CombineValues*, and *Flatten* operations into a single MapReduce [90]. This generalization of MapReduce supports multiple reducers and combiners and enables each reducer to produce multiple outputs, rather than enforcing the requirement that the reducer must produce outputs with the same key as its input. This solution enables the FlumeJava optimizer to produce better results.

M input channels, each performing a Map operation, feed into R output channels, each optionally performing a shuffle, an optional combine, and a Reduce operation. The *executor* of FlumeJava will run an operation locally if the input is relatively small. It will run parallel MapReduce remotely, though the overhead of launching a remote execution is larger, but the advantages for data larger inputs are significant. Temporary files for the outputs of all operations are created automatically and deleted when no longer necessary for later operations of the pipeline.

The system supports a *cached* execution mode when it first attempts to reuse the result of an operation from the previous run if it was saved in a (internal or user-visible) file and if it was not changed since. A result is unchanged if the inputs and the operation's code and saved state have not changed. This execution mode is useful for debugging an extended pipeline. Reference [90] reports that the largest pipeline had 820 unoptimized stages and 149 optimized stages.

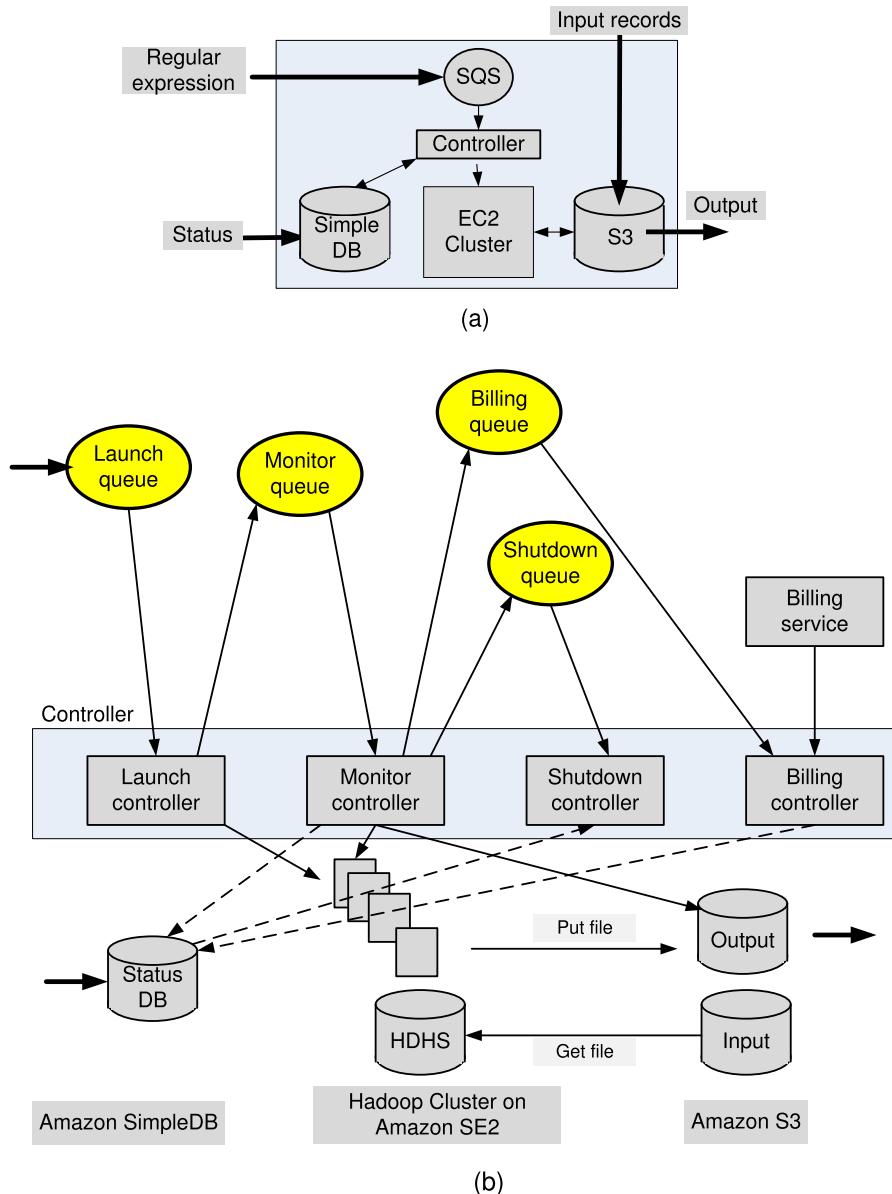
11.6 Case study: the GrepTheWeb application

Many applications process massive amounts of data using the MapReduce programming model. An application, GrepTheWeb [487], now in production at Amazon, illustrates the power and the appeal of cloud computing. The application allows a user to define a regular expression and search the web for records that match it. GrepTheWeb is analogous to the *grep* Unix command used to search a file for a given regular expression.

This application performs a search of a very large set of records attempting to identify records that satisfy a regular expression. The source of this search is a collection of document URLs produced by Alexa Web Search, a software system that crawls the web every night. The inputs to the applications are: (i) the large data set produced by the web-crawling software, and (ii) a regular expression. The output is the set of records that satisfy the regular expression. The user is able to interact with the application and get the current status; see Fig. 11.7(a).

The application uses message passing to trigger the activities of multiple controller threads that launch the application, initiate processing, shutdown the system, and create billing records. GrepTheWeb uses Hadoop MapReduce, an open-source software package that splits a large data set into chunks, distributes them across multiple systems, launches the processing, and, when the processing is complete, aggregates the outputs from various systems into a final result. Apache Hadoop is a

³ IDE (Integrated Drive Electronics) is an interface for connecting disk drives; the drive controller is integrated into the drive, as opposed to a separate controller on, or connected to, the motherboard.

**FIGURE 11.7**

The organization of the GrepTheWeb application. The application uses the Hadoop MapReduce software and four Amazon services: EC2, Simple DB, S3, and SQS. (a) The simplified workflow showing the two inputs, the regular expression and the input records generated by the web crawler; a third type of input are the user commands to report the current status and to terminate the processing. (b) The detailed workflow; the system is based on message passing between several queues; four controller threads periodically poll their associated input queues, retrieve messages, and carry out the required actions.

software library for distributed processing of large data sets across clusters of computers using a simple programming model.

GrepTheWeb workflow, illustrated in Fig. 11.7(b), consists of the following steps [487]:

1. *The start-up phase:* create several queues to launch, monitor, billing, and shutdown queues; start the corresponding controller threads. Each thread polls periodically its input queue and when a message is available, retrieves the message, parses it, and takes the required actions.
2. *The processing phase:* it is triggered by a *StartGrep* user request; then a launch message is enqueued in the launch queue. The launch controller thread picks up the message and executes the launch task; then, it updates the status and time stamps in the Amazon *Simple DB* domain. Lastly, it enqueues a message in the monitor queue and deletes the message from the launch queue. The processing phase consists of the following steps:
 - a. The launch task starts Amazon EC2 instances: it uses a Java Runtime Environment preinstalled Amazon Machine Image (AMI), deploys required Hadoop libraries, and starts a Hadoop job (run Map/Reduce tasks).
 - b. Hadoop runs Map tasks on EC2 slave nodes in parallel: a Map task takes files from S3, runs a regular expression, and writes locally the match results along with a description of up to five matches; then, the Combine/Reduce task combines and sorts the results and consolidates the output.
 - c. Final results are stored on Amazon S3 in the output bucket.
3. *The monitoring phase:* the monitor controller thread retrieves the message left at the beginning of the processing phase, validates the status/error in Simple DB, and executes the monitor task; it updates the status in the Simple DB domain and enqueues messages in the shutdown and the billing queues. The monitor task checks for the Hadoop status periodically and updates the Simple DB items with status/error and the S3 output file. Finally, it deletes the message from the monitor queue when the processing is completed.
4. *The shutdown phase:* the shutdown controller thread retrieves the message from the shutdown queue and executes the shutdown task that updates the status and time stamps in the Simple DB domain; finally, it deletes the message from the shutdown queue after processing. The shutdown phase consists of the following steps:
 - a. The shutdown task kills the Hadoop processes, terminates the EC2 instances after getting EC2 topology information from Simple DB, and disposes of the infrastructure.
 - b. The billing task gets the EC2 topology information, Simple DB usage, S3 file and query input, calculates the charges, and passes the information to the billing service.
5. *The cleanup phase:* archives the Simple DB data with user info.
6. *User interactions with the system:* get the status and output results. The *GetStatus* is applied to the service endpoint to obtain the status of the overall system (all controllers and Hadoop) and download the filtered results from S3 after completion.

Multiple S3 files are bundled up and stored as S3 objects to optimize the end-to-end transfer rates in the S3 storage system. Another performance optimization is to run a script and sort the keys, the URL pointers, and upload them in sorted order in S3. Multiple fetch threads are started to fetch the objects. This application illustrates the means to create an on-demand infrastructure and run it on a massively distributed system in a manner that enables it to run in parallel and scale up and down, based on the number of users and the problem size.

11.7 Hadoop, Yarn, and Tez

A wide range of data-intensive applications, such as marketing analytics, image processing, machine learning, and web crawling, use the Apache Hadoop, an open-source, Java-based software system.⁴ Hadoop supports distributed applications handling extremely large volumes of data. Many members of the community contributed to the development and optimization of Hadoop and of several related Apache projects such as Hive and HBase.

Hadoop is used by many organization from industry, government, and research. The long list of Hadoop users includes major IT companies, e.g., Apple, IBM, HP, Microsoft, Yahoo, and Amazon, media companies, e.g., *The New York Times* and Fox, social networks including, Twitter, Facebook, and LinkedIn, and government agencies such as the Federal Reserve. In 2012, the Facebook Hadoop cluster had a capacity of 100 petabytes and was growing at a rate of 0.5 petabytes a day. In 2013, more than half of Fortune 500 companies were using Hadoop. Azure HDInsight service deploys Hadoop on Microsoft Azure.

Hadoop. Apache Hadoop is an open-source software framework for distributed storage and distributed processing based on the MapReduce programming model. Recall that in MapReduce the Map stage processes the raw input data, one data item at a time, and produces a stream of data items annotated with keys. Next, a local sort stage orders the data produced during the Map stage by key. The locally ordered data is then passed to an (optional) combiner stage for partial aggregation by key. The shuffle stage then redistributes data among machines to achieve a global organization of data by key.

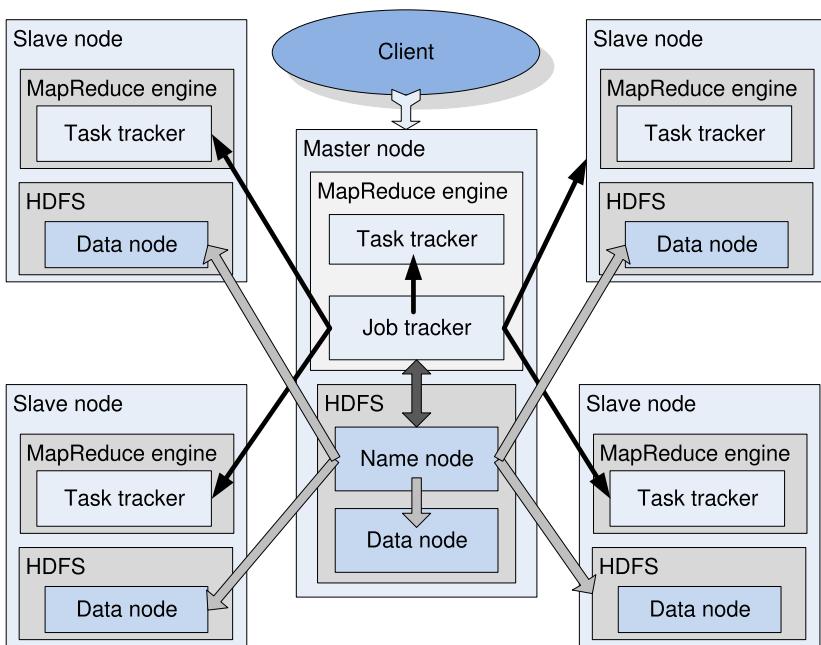
A Hadoop system has two components, a MapReduce engine and a database; see Fig. 11.8. The database could be the Hadoop File System (HDFS), Amazon S3, or CloudStore, an implementation of GFS, discussed in Section 7.6. HDFS is a highly performant distributed file system written in Java. HDFS is portable but cannot be directly mounted on an existing operating system; it is not fully POSIX compliant.

The Hadoop engine on the master of a multinode cluster consists of a *job tracker* and a *task tracker*, while the engine on a slave has only a *task tracker*. The *job tracker* receives a MapReduce job from a client and dispatches the work to the *task trackers* running on the nodes of a cluster. To increase efficiency, the *job tracker* attempts to dispatch the tasks to the available slaves closest to the place where the task data was stored. The *task tracker* supervises the execution of the work allocated to the node; several scheduling algorithms have been implemented in Hadoop engines, including Facebook's fair scheduler and Yahoo's capacity schedulers.

HDFS replicates data on multiple nodes; the default is three replicas, and a large dataset is distributed over many nodes. The *name node* running on the master manages the data distribution and data replication and communicates with *data nodes* running on all cluster nodes; it shares with the *job tracker* information about the data placement to minimize communication between the nodes where data is located and the ones where it is needed. Although HDFS can be used for applications other than those based on the MapReduce model, its performance for such applications is not on par with the ones it was originally designed for.

Hadoop brings computations to the data on clusters built with off-the-shelf components. This strategy is pushed further by Spark, which stores data in processor's memory instead of the disk. Data

⁴ Hadoop requires JRE (Java Runtime Environment) 1.6 or higher.

**FIGURE 11.8**

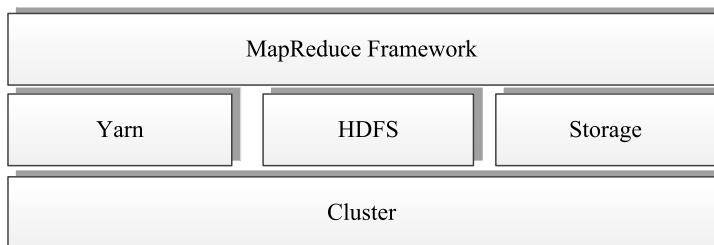
A Hadoop cluster using HDFS; the cluster includes a master and four slave nodes. Each node runs a MapReduce engine and a database engine, often HDFS. The *job tracker* of the master's engine communicates with the *task trackers* on all the nodes and with the *name node* of HDFS. The *name node* of the HDFS shares information about the data placement with the *job tracker* to minimize communication between the nodes where data is located and the ones where it is needed.

locality allows Hadoop and Spark to compete with traditional High Performance Computing (HPC) running on supercomputers with high-bandwidth storage and faster interconnection networks.

The Apache Hadoop framework has the following modules:

1. Common—contains libraries and utilities needed by all Hadoop modules.
2. Distributed File System (HDFS)—a distributed file-system that stores data on commodity machines, providing very high aggregate bandwidth across the cluster.
3. Yarn—a resource-management platform responsible for managing computing resources in clusters and using them for scheduling of users' applications.
4. MapReduce—an implementation of the MapReduce programming model.

Additional software packages such as Apache Pig, Apache Hive, Apache HBase, Apache Phoenix, Apache Spark, Apache ZooKeeper, Cloudera Impala, Apache Flume, Apache Sqoop, Apache Oozie, and Apache Storm are also available.

**FIGURE 11.9**

Resources needed by the MapReduce framework are provided by the cluster and are managed by Yarn; HDFS provides permanent, reliable, and distributed storage; multiple organizations of the storage system are supported, e.g., AWS implementation of Hadoop offers S3.

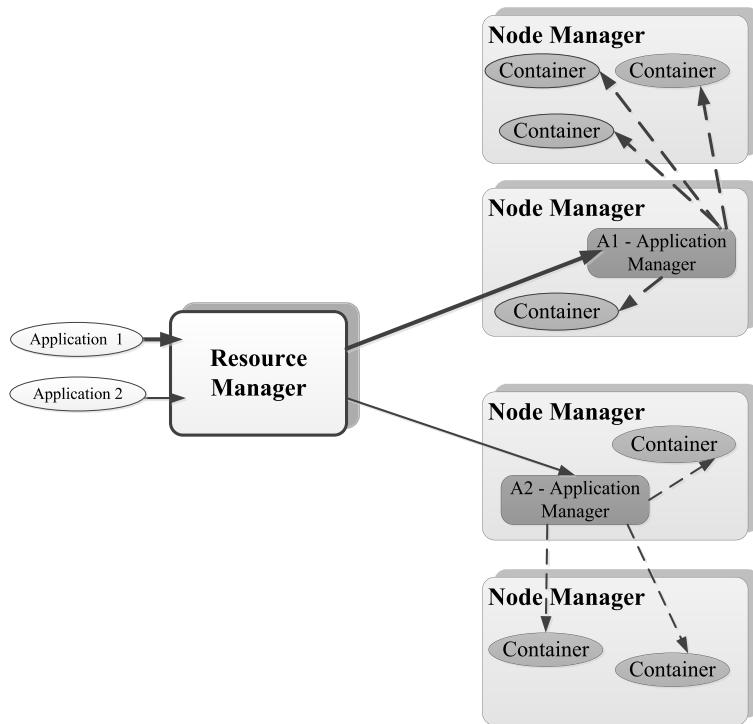
Several Hadoop frameworks are used to manage and run deep analytics. SQL processing is particularly important to gain insights from large collections of data and the number of SQL-on-Hadoop systems has increased. There are many SQL engines running on Hadoop, including BigSQL from IBM, Impala from Cloudera, and HAWQ from Pivotal. All these engines implement a standard language specification and compete on performance and extended services. Users are insulated from difficulties because the applications that talk to those engines are portable.

Systems such as Pig, Hive, and Impala, discussed in Section 11.8, are native Hadoop-based systems, or database-Hadoop hybrids. HadoopOthers, such as *Hadapt* [5], exploits Hadoop scheduling and fault-tolerance, but uses a relational database, PostgreSQL, to execute query fragments.

Yarn. Yarn is a resource management system supplying CPU cycles, memory, and other resources needed by a single job or to a DAG of MapReduce applications. Resource allocation and job scheduling in Hadoop versions prior to 2.0 were done by the MapReduce framework. In the newer versions of Hadoop, Yarn carries out these functions. This new system organization allows frameworks such as Spark to share cluster resources. The organization of Hadoop including Yarn in Fig. 11.9 shows the three system elements: (1) the MapReduce framework; (2) Yarn, HDFS, and the storage substrate; and (3) the cluster where the application is running.

Fig. 11.10 presents the organization of Yarn and shows the Resource Manager and Node Managers running in each node. A node manager is responsible for containers, monitors their resource usage (CPU, memory, disk, network), and reports to the resource manager tasked to arbitrate resources sharing among all applications. Each application has an Application Manager that negotiates with the resource manager the access to resources needed by the application. Once resources are allocated, the application manager interacts with the node managers of each node allocated to the application to start the tasks and then monitors their execution.

The Scheduler component of the resource manager uses the resource container abstraction that incorporates memory, CPU, disk, network, etc. and bases the resource allocation decisions on applications. The scheduler performs no monitoring or tracking application status and offers no guarantees about restarting failed tasks. A pluggable policy is responsible for sharing cluster resources among applications. The Capacity Scheduler and the Fair Scheduler are examples of scheduler plug-ins.

**FIGURE 11.10**

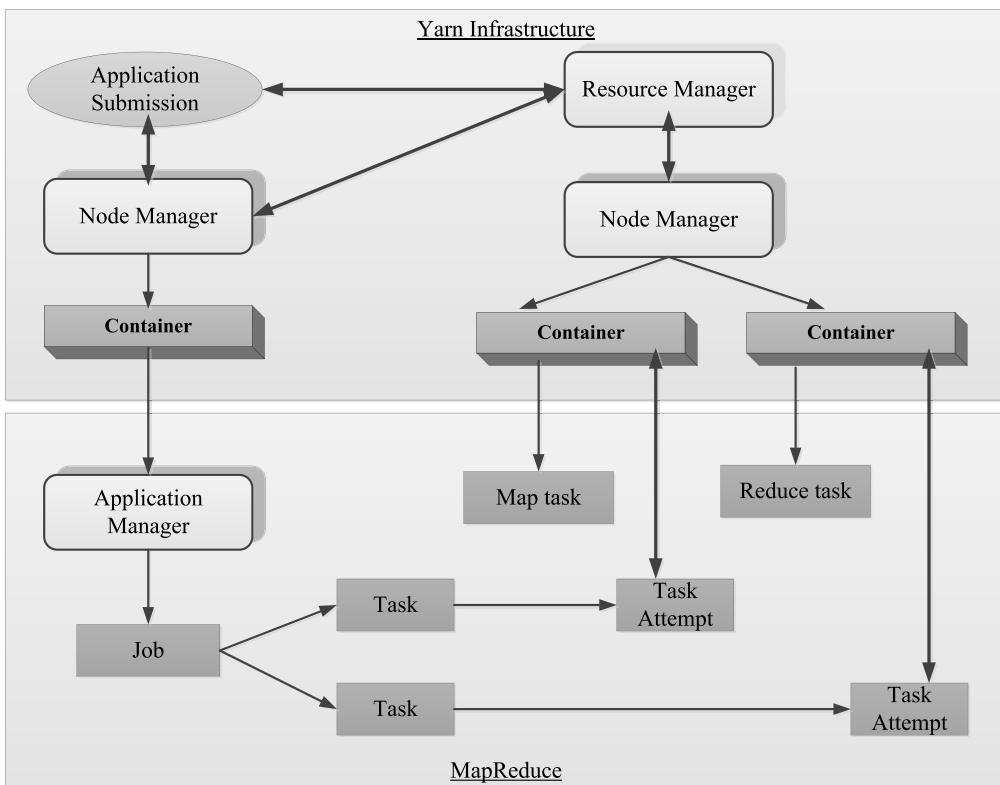
Yarn organization and application processing. Applications are submitted to the *Resource Manager*. The Resource Manager communicates with Node Managers running on every node of the cluster. The tasks of every application are managed by an Application Manager. Each task is packaged in a container.

MapReduce in Hadoop-2.x maintains API compatibility with previous stable release, thus MapReduce jobs should still run unchanged on top of Yarn after a recompile.

The process of starting an application involves several steps, illustrated in Fig. 11.11:

1. The user submits an application to the Resource Manager.
2. The Resource Manager invokes the Scheduler and allocates a container for the Application Manager.
3. The Resource Manager contacts the Node Manager where the container will be launched.
4. The Node Manager launches the container.
5. The container executes the Application Manager.
6. The Resource Manager contacts the Node Manager(s) where the tasks of the application will run.
7. Containers for the tasks of the application are created.
8. The Application Manager monitors the execution of the tasks until termination.

Tez. Tez is an extensible framework for building high-performance batch and interactive processing for Yarn-based applications in Hadoop. A job is decomposed into individual tasks, and each task of the job

**FIGURE 11.11**

The interaction among the components of Yarn and MapReduce. Once an application is submitted, Yarn's Resource Manager contacts a Node Manager to create a container for the Application Manager. Then, the application manager is activated. The resource manager with the assistance of the scheduler selects the node manager(s) where the containers for the application tasks are created. Then, the containers start the execution of these tasks.

runs as an Yarn process. Tez models data processing as a DAG; the graph vertices represent application logic, and edges represent movement of data. A Java API is used to express the DAG representation of the workflow. The execution engine uses Yarn to acquire resources and reuses every component in the pipeline to avoid operation duplication. Apache Hive and Apache Pig use Apache Tez to improve the speed of MapReduce applications; see <http://hortonworks.com/apache/tez/>.

11.8 SQL on Hadoop: Pig, Hive, and Impala

The landscape of parallel SQL database vendors prior to 2009 included IBM, Oracle, Microsoft, ParAccel, Greenplum, Teradata, Netezza, and Vertica, but none of them were supporting SQL queries along

with MapReduce jobs. From their early years in business, Yahoo and Facebook used the MapReduce platform extensively to store, process, and analyze huge amounts of data. Both companies wanted a faster and easier-to-work-with platform supporting SQL queries.

MapReduce is a heavyweight, high-latency execution framework and does not support workflows, join operations for combined processing of several datasets, filtering, aggregation, top-k thresholding, and high-level operations. To address these challenges, Yahoo created a dataflow system called Pig in 2009. Facebook built Hive on MapReduce because it was the shortest path to SQL on Hadoop. Pig and Hive have their own language, Pig Latin and HiveQL, respectively. In both systems, a user types a query, then a parser reads the query, figures out what the user wants, and runs a series of MapReduce jobs. That was a sensible decision given the requirements of the time.

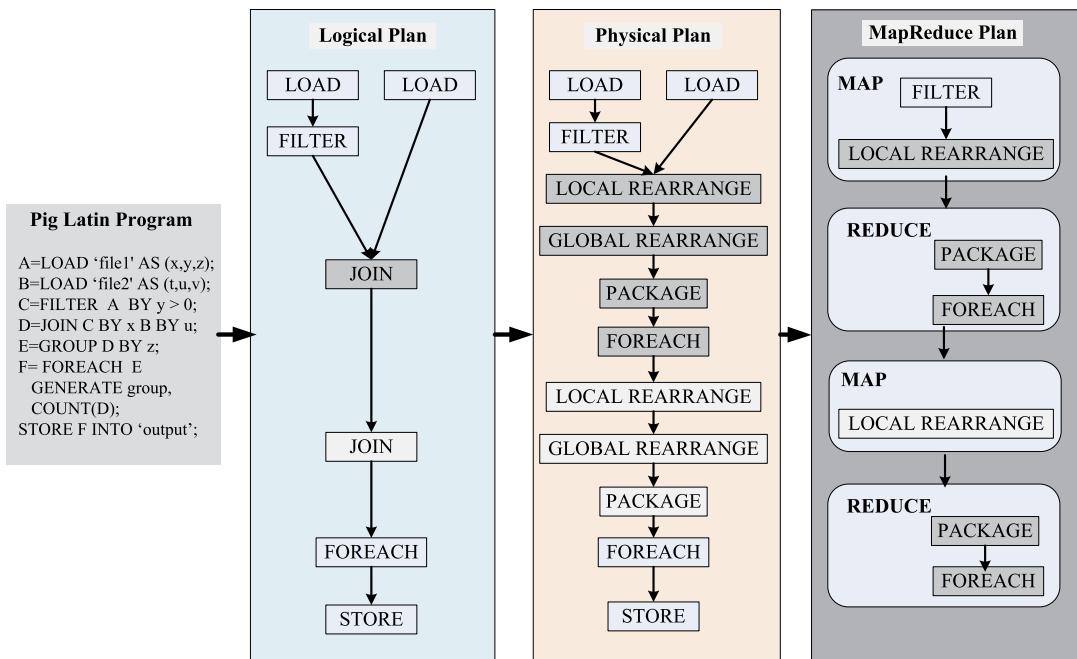
In 2012, Cloudera developed Impala and published it under an Apache license. Later, Facebook developed Presto, its next-generation query processing engine for real-time access to data via SQL. It was built, like Hive, from the ground up, as a distributed query processing engine.

Pig. The system presented in [190] supports workflows, joint operations for combined processing of several data sets, filtering, aggregation, and high-level operations. Pig compiles programs written in a language called Pig Latin into a set of Hadoop jobs and coordinates their execution. It also supports several user interaction modes: (1) Interactive—using a shell for Pig commands; (2) Batch—a user submits a script containing a series of Pig commands; and (3) Embedded—commands are submitted via method invocation from a Java program. The job processing stages are:

1. Parsing—the parser performs syntactic verification, type checking, and schema inference and produces a DAG called a Logical Plan. In this plan, an operator is annotated with the schema of its output data, with braces indicating a bag of tuples.
2. Logical optimization of the DAG and creation of a Physical Plan describing data distribution.
3. Compilation of the optimized Physical Plan into a set of MapReduce jobs, followed by optimization phase, e.g., partial aggregation, resulting in an optimized DAG. Distributive and algebraic aggregation functions, e.g., AVERAGE, are broken into series of three steps: initial (e.g., generate [sum, count] pairs), intermediate (e.g., combine n [sum, count] pairs into a single pair), final (e.g., combine n [sum, count] pairs and take the quotient). These steps are assigned to the Map, Combine, and Reduce stages, respectively.
4. The DAG is topologically sorted and jobs are submitted to Hadoop for execution. The flow control is implemented using a pull model with a single threaded implementation and a simple API for user-defined functions for moving tuples through the execution pipeline. An operator can respond in one of three ways when asked to produce a tuple: return the tuple, declare that it has finished, or return a pause signal indicating that either it is not finished or unable to produce an output tuple.

Compilation and execution is triggered by the STORE command. If a Pig Latin program contains more than one STORE command, the generated physical plan contains a SPLIT physical operator. Fig. 11.12 illustrates the transformation of a Pig Latin program into a Logical Plan, followed by the transformation of the Logical Plan into a Physical Plan, and finally the transformation of the Physical Plan into a MapReduce Plan.

Memory management is challenging because Pig is implemented in Java, and program controlled memory allocation and deallocation is not feasible. A handler can be registered using the *MemoryPoolMXBean* Java class. The handler is notified whenever a configurable memory threshold is reached; then the system releases registered bags until enough memory is freed.

**FIGURE 11.12**

The transformation of a Pig Latin program into a Logical Plan, followed by the transformation of the Logical Plan into a Physical Plan, and finally the transformation of the Physical Plan into a MapReduce Plan. Each one of the two JOINs in the Logical Plan generate LOCAL REARRANGE, GLOBAL REARRANGE, PACKAGE, and FOREACH statements. A pair of MAP, REDUCE (is then generated in the MapReduce Plan.

Pig Mix is a benchmark used at Yahoo and elsewhere to evaluate system's performance and to exercise a wide range of the system functionality. About 60% of the ad hoc Hadoop jobs at Yahoo use Pig. The use of the system outside Yahoo has been increasing.

Hive. Hive is an open-source system for data-warehousing supporting queries expressed in an SQL-like declarative language called HiveQL [475]. These queries are then compiled into MapReduce jobs executed on Hadoop. The system includes the Metastore, a catalog for schemas, and the statistics used for query optimization.

Hive supports data organized as tables, partitions, and buckets. *Tables* are inherited from relational databases and can be stored internally or externally in HDFS, NFS, or a local directory. A table is serialized and stored in the files of an HDFS directory. The format of each serialized table is stored in the system catalog, and it is accessed during query compilation and execution.

Partitions are components of tables described by the subdirectories of table directory. In turn, *buckets* consists of partition data stored as a file in the partition's directory and selected based on the hash of one of the columns of the table. The query language supports data-definition statements to create ta-

bles with specific serialization formats and partitioning and bucketing columns, as well as user-defined column transformation and aggregation functions implemented in Java.

HiveQL accepts as input DDL, DML and allows user-defined MapReduce scripts written in any language using a simple row-based streaming interface. Data Description Language (DDL) has a syntax for defining data structures similar to a computer programming language, and it is widely used for database schemas. Data Manipulation Language (DML) is used to retrieve, store, modify, delete, insert, and update data in a database; SELECT, UPDATE, INSERT statements or query statements are examples of DML.

The system has several components:

- *External Interface*—includes a command line (CLI), a web user-interface (UI), and language APIs such as JDBC and ODBC.⁵
- *Thrift Server*—a client API for execution of HiveQL statements.⁶ Java and C++ clients can be used to build JDBC or ODBC common drivers, respectively.
- *Megastore*—the system catalog. It contains several objects: (1) Database—a namespace for tables; (2) Table—contains the list of columns and their types, owner, storage, the location of the table data, data formats, and bucketing information; (3) Partition—each partition can have its own columns and storage information.
- *Driver*—manages the compilation, optimization, and execution of HiveQL statements.
- *Database*—the namespace for tables.
- *Table*—the metadata for tables containing the list of columns and their types, owner, storage, and a wealth of other data including the location of the table data, data formats, and bucketing information. It also includes SerDe metadata regarding the implementation class of serializer and deserializer methods and information required by their implementation.
- *Partition*—information about the columns in a partition including SerDe and storage information.

The HiveQL compiler has several components: the *Parser* transforms an input string into a parse tree; then, the *Semantic Analyzer* transforms this tree into an internal representation, the *Logical Plan Generator* converts this internal representation into a *Logical Plan*, and finally the *Optimizer* rewrites the logical plan.

Facebook's Hive warehouse contains over 700 TB of data and supports more than 5 000 daily queries. Hive is an Apache project, with an active user and developer community.

Impala is a query engine exploiting a shared-nothing parallel database architecture written in C++ and Java and designed to use standard Hadoop components (HDFS, HBase, Metastore, Yarn, Sentry) [282]. The system, developed by Cloudera and published under an Apache license in 2012, is designed from the ground up for SQL query execution on Hadoop, rather than a general-purpose distributed processing system. It delivers better performance than Hive because it does not translate an SQL query into another processing framework as Hive does. It supports most of the SQL-92 SELECT statement syntax and SQL-2003 analytic functions. It does not support UPDATE or DELETE, but supports bulk insertions *INSERT INTO ... SELECT*

⁵ Java Database Connectivity (JDBC) is an API for Java, defining how a client may access a database. Open Database Connectivity (ODBC) is an open standard application API for accessing a database.

⁶ Thrift is a framework for cross-language services; see Apache Thrift, <http://incubator.apache.org/thrift>.

Every node in a Cloudera cluster has Impala code installed and waiting for SQL queries to execute. The Impala code lives on every node alongside MapReduce, Apache HBase, and third-party engines such as SAS and Apache Spark that a customer may choose to deploy. All these engines have access to the same data, and users can choose one of them depending on the application. Impala runs queries using long-running daemons on every HDFS DataNode, and pipelines intermediate results between computation stages.

The I/O layer of Impala spawns one I/O thread per disk on each node to read data stored in HDFS and achieves high utilization of both CPU and disks by decoupling asynchronous read requests from the synchronous actual reading of data. The system exploits Intel's SSE instructions discussed in Section 3.4 to parse and process textual data efficiently. Impala requires the working set of a query to fit in the aggregate physical memory of the cluster.

The system operates basic services offered by three daemons:

1. *Impalad*—accepts, plans, and coordinates query execution. An Impalad daemon is deployed on every server and operates a query planner, a query coordinator, and a query executor. The front end compiles SQL text into query plans executable by the back ends. In this stage, a parse tree is translated into a single-node plan tree including: HDFS/HBase scan, hash join, cross join, union, hash aggregation, sort, top-n, and analytic evaluation nodes. A second phase transforms the single-node plan into a distributed execution plan, and in this process, the goal is to minimize data movement and maximize scan locality.
2. *Statestored*—provides a metadata publish-subscribe service and disseminates cluster-wide metadata to all processes. It maintains a set of topics (key, value, version), triplets defined by an application. Processes wishing to receive updates to any topic express their interest by registering at start-up and providing a list of topics.
3. *Catalogd*—is the catalog repository and metadata access gateway. It pulls information from third-party metadata stores and aggregates it. Only a skeleton entry for each table it discovers is loaded at startup; then, table metadata is loaded in the background from third-party stores.

A recent paper [176] compares the performance of Impala and Hive. Both systems input their data from columnar storage in a Parquet and the Optimized Row Columnar, the Apache columnar storage formats shared by the software in the Hadoop ecosystem, regardless of the choice of data processing framework, data model, or programming language. Columnar formats improve the performance of queries in the context of relational databases.

Parquet stores data grouped together in logical horizontal partitions called *row groups*. Every row group contains a *column chunk* for each column in the table. A column chunk consists of multiple *pages* and is guaranteed to be stored contiguously on disk. Compression and encoding schemes work at a page level. Metadata is stored at all the levels in the hierarchy, i.e., file, column chunk, and page. An ORC file stores multiple groups of row data as *stripes* and has a *footer* containing the list of the stripes in the file, the number of rows stored in each stripe, the data type of each column, and column-level aggregates, such as: count, sum, min, and max.

The experiments in [176] used Hive version 0.12 (Hive-MR) and Impala version 1.2.2 on top of Hadoop 2.0.0-cdh4.5.0, Hive version 0.13 (Hive-Tez) on top of Tez 0.3.0, and Apache Hadoop 2.3.0.1. Hadoop was configured to run 12 containers per node, 1 per core. The HDFS replication factor is three, and the maximum JVM heap size is set to 7.5 GB per task. Impala uses MySQL as the metastore. One Impalad process runs on each compute node and has access to 90 GB of memory.

One of the nodes of a 21-node cluster used for the measurements hosts the HDFS NameNode, and there are 20 compute nodes. Each node runs 64-bit Ubuntu Linux 12.04 and has: one Intel Xeon CPUs @ 2.20GHz with 6 cores, eleven SATA disks (2TB, 7k RPM), one 10-Gigabit Ethernet card, and 96 GB of RAM. One out of the eleven SATA disks in each node hosts the OS, and the rest are used for HDFS.

The experiments reported in [176] run the 22 TPC-H⁷ queries for Hive and Impala and report the execution time for each query. The file cache is flushed before each run in all compute nodes. The results show that Impala outperforms both *Hive-MapReduce* and Hive-Tez for all file formats, with or without compression; the improvements are in the range 1.5–13.5 times. Several factors contribute to the drastic performance boost:

- Impala has a more efficient I/O subsystem than Hive-MR and Hive-Tez.
- Long-running daemon processes handle queries in each node thus, the overhead of job initialization and scheduling due to MapReduce in Hive-MR is eliminated.
- The query execution is pipelined, while in Hive-MR data is written out at the end of each step and read in by subsequent step(s).

Impala uses code generation to eliminate the overhead of virtual function calls, inefficient instruction branching due to large switch statements, and of other sources of inefficiency. Query execution time improves by about 1.3 times for all 21 TPC-H queries combined when this feature is enabled at runtime. TPC-H Query 1 shows the largest improvement, 5.5 times, while the remaining queries improve up to 1.75 times.

A second set of experiments involve TPC-DS benchmark.⁸ Results show that Impala is on average 8.2 times faster than Hive-MR and 4.3 times faster than Hive-Tez on TPC-DS and 10 times faster than Hive-MR and 4.4 times faster than Hive-Tez on a second workload. The first workload consists of 20 queries that access a single fact table and six dimension tables, while the second uses the same workload but removes the explicit partitioning predicate and uses the correct predicate values.

11.9 Current cloud applications and new applications opportunities

Existing cloud applications can be divided in several broad categories: (i) processing pipelines; (ii) batch processing systems; and (iii) web applications [487].

Processing pipelines are data-intensive and sometimes compute-intensive applications that represent a fairly large segment of applications currently running on a cloud. Several types of data processing applications can be identified:

- Indexing; processing pipeline supports indexing of large datasets created by web-crawler engines.
- Data mining; processing pipeline supports searching large sets of records to locate items of interests.

⁷ The TPC BenchmarkH (TPC-H) is a decision-support benchmark consisting of a suite of business-oriented ad hoc queries and concurrent data modifications chosen to have broad industry-wide relevance. This benchmark is relevant for applications examining a large volume of data and executing queries with a high degree of complexity.

⁸ TPC-DS is a de-facto industry standard benchmark for assessing the performance of decision support systems.

- Image processing; a number of companies allow users to store their images on the cloud, e.g., Flickr (flickr.com) and Picasa (<http://picasa.google.com/>). Image-processing pipelines support image conversion, e.g., enlarge an image or create thumbnails, and compress or encrypt images.
- Video transcoding; processing pipeline translates one video format to another, e.g., AVI to MPEG.
- Document processing; processing pipeline converts very large collection of documents from one format to another, e.g., from Word to PDF or encrypts the documents; they could also use OCR (Optical Character Recognition) to produce digital images of documents.

Batch processing systems also cover a broad spectrum of data-intensive applications in enterprise computing. Such applications typically have deadlines, and the failure to meet these deadlines could have serious economic consequences; security is also a critical aspect for many applications of batch processing. A nonexhaustive list of batch processing applications includes:

- Generation of daily, weekly, monthly, and annual activity reports for organizations in retail, manufacturing, and other economic sectors.
- Processing, aggregation, and summaries of daily transactions for financial institutions, insurance companies, and healthcare organizations.
- Inventory management for large corporations.
- Processing billing and payroll records.
- Management of the software development, e.g., nightly updates of software repositories.
- Automatic testing and verification of software and hardware systems.

Lastly, cloud applications in the area of web access are of increasing importance. Several groups of web sites have a periodic or temporary presence, e.g., web sites for conferences or other events. There are also web sites active during a particular season (e.g., the holiday season) or supporting a particular type of activity, such as US income tax reporting with the April 15 deadline each year. Other limited-time web site are used for promotional activities, or web sites that “sleep” during the night and auto-scale during the day.

It makes economic sense to store the data in the cloud close to where the application runs; as we have seen in Section 2.2, the cost per GB is low, and the processing is much more efficient when the data is stored close to the computational servers. This leads us to believe that several new classes of cloud computing applications could emerge in the years to come, for example, batch processing for decision support systems and other aspects of business analytics.

Another class of new applications could be parallel batch processing based on programming abstractions such as MapReduce, discussed in Section 11.5. Mobile interactive applications that process large volumes of data from different types of sensors and services that combine more than one data source, e.g., mashups,⁹ are obvious candidates for cloud computing.

Science and engineering could benefit from cloud computing as many applications are compute- and data-intensive. Similarly, a cloud dedicated to education would be extremely useful. Mathematical software, e.g., MATLAB[®] and Mathematica, could also run on the cloud.

⁹ A mashup is an application that uses and combines data, presentations, or functionality from two or more sources to create a service. The fast integration, frequently using open APIs and multiple data sources, produces results not envisioned by the original services; combination, visualization, and aggregation are the main attributes of mashups.

11.10 Clouds for science and engineering

For more than two thousand years, science was empirical. Several hundred years ago, theoretical methods based on models and generalization were introduced, and this propelled substantial progress in human knowledge. In the last few decades, we have witnessed the explosion of computational science based on the simulation of complex phenomena.

In a talk delivered in 2007 and posted on his web site just before he went missing in January 2007, Jim Gray discussed *eScience* as a transformative scientific method [235]. Today, *eScience* unifies experiment, theory, and simulation; data captured from measuring instruments, or generated by simulations is processed by software systems; data and knowledge are stored by computer systems and analyzed with statistical packages.

The generic problems in virtually all areas of science are: (1) collection of experimental data; (2) management of a very large volumes of data; (3) building and execution of models; (4) integration of data and literature; (5) documentation of the experiments; (6) sharing the data with others; and (7) data preservation for long periods of time. All these activities require powerful computing systems.

A typical example of a problem faced by agencies and research groups is data discovery in large scientific data sets. Examples of such large collections are the biomedical and genomic data at NCBI,¹⁰ the astrophysics data from NASA,¹¹ and the atmospheric data from NOAA¹² and NCAR.¹³ The process of online data discovery can be viewed as an ensemble of several phases: (i) recognition of the information problem; (ii) generation of search queries using one or more search engines; (iii) evaluation of the search results; (iv) evaluation of the web documents; and (v) comparing information from different sources. The web-search technology enables the scientists to discover text documents related to such data, but the binary encoding of many of them poses serious challenges.

High Performance Computing on AWS. Reference [256] describes the set of applications used at NERSC (National Energy Research Scientific Computing Center) and presents the results of a comparative benchmark of EC2 and three supercomputers. NERSC is located at Lawrence Berkeley National Laboratory and serves a diverse community of scientists; it has some 3 000 researchers and involves 400 projects based on some 600 codes. Some of the codes used are:

CAM (Community Atmosphere Model), the atmospheric component of CCSM (Community Climate System Model), is used for weather and climate modeling.¹⁴ The code developed at NCAR uses two two-dimensional domain decompositions; one for the dynamics and the other for re-mapping. The first is decomposed over latitude and vertical level and the second is decomposed over longitude–latitude. The program is communication-intensive; on-node/processor data movement and relatively long MPI¹⁵ messages that stress the interconnect point-to-point bandwidth are used to move data between the two decompositions.

¹⁰ NCBI is the National Center for Biotechnology Information, <http://www.ncbi.nlm.nih.gov/>.

¹¹ NASA is the National Aeronautics and Space Administration, <http://www.nasa.gov/>.

¹² NOAA is the National Oceanic and Atmospheric Administration, www.noaa.gov.

¹³ NCAR is the National Center for Atmospheric Research.

¹⁴ See <http://www.nersc.gov/research-and-development/benchmarking-and-workload-characterization>.

¹⁵ MPI, Message Passing Interface, is a communication library based on a standard for a portable message-passing system.

GAMESS (General Atomic and Molecular Electronic Structure System) is used for *ab initio* quantum chemistry calculations. The code developed by the Gordon research group at the Department of Energy's Ames Lab at Iowa State University has its own communication library, the Distributed Data Interface (DDI) that is based on the SPMD (Same Program Multiple Data) execution model. DDI presents the abstraction of a global shared memory with one-side data transfers, even on systems with physically distributed memory. On the cluster systems at NERSC, the program uses socket communication; on the Cray XT4, the DDI uses MPI and only one-half of the processors compute, while the other half are data movers. The program is memory- and communication-intensive.

GTC (Gyrokinetic¹⁶) is a code for fusion research.¹⁷ It is a self-consistent, gyrokinetic tri-dimensional Particle-in-cell (PIC)¹⁸ code with a non-spectral Poisson solver; it uses a grid that follows the field lines as they twist around a toroidal geometry representing a magnetically confined toroidal fusion plasma. The version of GTC used at NERSC uses a fixed, one-dimensional domain decomposition with 64 domains and 64 MPI tasks. Communication is dominated by the nearest neighbor exchanges that are bandwidth-bound. The most computationally intensive parts of GTC involve gather/deposition of charge on the grid and particle “push” steps. The code is memory intensive because the charge deposition uses indirect addressing.

IMPACT-T (Integrated Map and Particle Accelerator Tracking Time) is a code for the prediction and performance enhancement of accelerators; it models the arbitrary overlap of fields from beam-line elements and uses a parallel, relativistic PIC method with a spectral integrated Green function solver. This object-oriented *Fortran90* code uses a two-dimensional domain decomposition in the $y - z$ directions and dynamic load balancing based on the domains. Hockney’s FFT (Fast Fourier Transform) algorithm is used to solve Poisson’s equation with open boundary conditions. The code is sensitive to the memory bandwidth and MPI collective performance.

MAESTRO is a low Mach-number hydrodynamics code for simulating astrophysical flows.¹⁹ Its integration scheme is embedded in an adaptive mesh refinement algorithm based on a hierarchical system of rectangular nonoverlapping grid patches at multiple levels with different resolution; it uses a multigrid solver. Parallelization is via a tri-dimensional domain decomposition using a coarse-grained distribution strategy to balance the load and minimize communication costs. The communication topology tends to stress simple topology interconnects. The code has a very low computational intensity; it stresses memory latency, and the implicit solver stresses global communications; the message sizes range from short to relatively moderate.

MILC (MIMD Lattice Computation) is a QCD (Quantum Chromo Dynamics) code used to study “strong” interactions binding quarks in protons and neutrons and holding them together in the nucleus.²⁰ The algorithm discretizes the space and evaluates field variables on sites and the links of a

¹⁶ The trajectory of charged particles in a magnetic field is a helix that winds around the field line; it can be decomposed into a relatively slow motion of the guiding center along the field line and a fast circular motion called cyclotronic motion. Gyrokinetics describes the evolution of the particles without taking into account the circular motion.

¹⁷ See <http://www.scidareview.org/0601/html/news4.html>.

¹⁸ PIC is a technique to solve a certain class of partial differential equations; individual particles (or fluid elements) in a Lagrangian frame are tracked in continuous phase space, whereas moments of the distribution such as densities and currents are computed simultaneously on Eulerian (stationary) mesh points.

¹⁹ See <http://www.astro.sunysb.edu/mzingale/Maestro/>.

²⁰ See <http://physics.indiana.edu/~sg/milc.html>.

regular hypercube lattice in four-dimensional space-time. The integration of an equation of motion for hundreds or thousands of time steps requires inverting a large, sparse matrix. The CG (Conjugate Gradient) method is used to solve a sparse, nearly-singular matrix problem. Many CG iterations steps are required for convergence; the inversion translates into tri-dimensional complex matrix–vector multiplications. Each multiplication requires a dot product of three pairs of tri-dimensional complex vectors; a dot product consists of five multiply–add operations and one multiply. The MIMD computational model is based on a four-dimensional domain decomposition; each task exchanges data with its eight nearest neighbors and is involved in the *all-reduce* calls with very small payloads as part of the CG algorithm; the algorithm requires *gather* operations from widely separated locations in memory. The code is highly memory- and computational-intensive and is heavily dependent on prefetching.

PARATEC (PARAllel Total Energy Code) is a quantum mechanics code; it performs *ab initio* total-energy calculations using pseudopotentials, a plane wave basis set and an all-band (unconstrained) conjugate gradient (CG) approach. Parallel three-dimensional FFTs transform the wave functions between real and Fourier space. The FFT dominates the runtime; the code uses MPI and is communication-intensive. The code uses mostly point-to-point short messages. The code parallelizes over grid points, thereby achieving a fine-grain level of parallelism. The BLAS3 and one-dimensional FFT use optimized libraries, e.g., Intel’s MKL or AMD’s ACML, and this results in high cache reuse and a high percentage of per-processor peak performance.

The authors of [256] use the HPCC (High Performance Computing Challenge) benchmark to compare the performance of EC2 with the performance of three large systems at NERSC. HPCC²¹ is a suite of seven synthetic benchmarks: three targeted synthetic benchmarks that quantify basic system parameters that characterize individually the computation and communication performance; four complex synthetic benchmarks that combine computation and communication and can be considered simple proxy applications. These benchmarks are:

- DGEMM²²—the benchmark measures the floating point performance of a processor/core; the memory bandwidth does little to affect the results as the code is cache friendly. Thus, the results of the benchmark are close to the theoretical peak performance of the processor.
- STREAM²³—the benchmark measures the memory bandwidth.
- The network latency benchmark.
- The network bandwidth benchmark.
- HPL²⁴—a software package that solves a (random) dense linear system in double precision arithmetic on distributed-memory computers; it is a portable and freely available implementation of the High Performance Computing Linpack Benchmark.
- FFTE—measures the floating-point rate of execution of double precision complex one-dimensional DFT (Discrete Fourier Transform)

²¹ For more information, see <http://www.novellshareware.com/info/hpc-challenge.html>.

²² For more details, see <https://computecanada.org/?pageId=138>.

²³ For more details, see <http://www.streambench.org/>.

²⁴ For more details, see <http://netlib.org/benchmark/hpl/>.

Table 11.1 The results of the measurements reported in [256].

System	DGEMM Gflops	STREAM GB/s	Latency μs	Bndw GB/S	HPL Tflops	FFTE Gflops	PTRANS GB/s	RandAcc GUPS
Carver	10.2	4.4	2.1	3.4	0.56	21.99	9.35	0.044
Frankl	8.4	2.3	7.8	1.6	0.47	14.24	2.63	0.061
Lawren	9.6	0.7	4.1	1.2	0.46	9.12	1.34	0.013
EC2	4.6	1.7	145	0.06	0.07	1.09	0.29	0.004

- PTRANS—parallel matrix transpose; it exercises the communications where pairs of processors communicate with each other simultaneously. It is a useful test of the total communications capacity of the network.
- RandomAccess—measures the rate of integer random updates of memory (GUPS).

The systems used for the comparison with cloud computing are:

Carver—a 400-node IBM iDataPlex cluster with quad-core Intel Nehalem processors running at 2.67 GHz and with 24 GB of RAM (3 GB/core). Each node has two sockets; a single Quad Data Rate (QDR) IB link connects each node to a network that is locally a fat-tree with a global two-dimensional mesh. The codes were compiled with the Portland Group suite version 10.0 of and Open MPI version 1.4.1.

Franklin—a 9 660-node Cray XT4; each node has a single quad-core 2.3 GHz AMD Opteron “Budapest” processor with 8 GB of RAM (2 GB/core). Each processor is connected through a 6.4 GB/s bidirectional HyperTransport interface to the interconnect via a Cray SeaStar-2 ASIC. The SeaStar routing chips are interconnected in a tri-dimensional torus topology, where each node has a direct link to its six nearest neighbors. Codes were compiled with the Pathscale or the Portland Group version 9.0.4.

Lawrencium—a 198-node (1 584 core) Linux cluster; a compute node is a Dell Poweredge 1950 server with two Intel Xeon quad-core 64 bit and 2.66 GHz Harpertown processors with 16 GB of RAM (2 GB/core). A compute node is connected to a Dual Data Rate InfiniBand network configured as a fat tree with a 3 : 1 blocking factor. Codes were compiled using Intel 10.0.018 and Open MPI 1.3.3.

The virtual cluster at Amazon had four EC2 CUs (Compute Units), two virtual cores with two CUs each, and 7.5 GB of memory (an `m1.large` instance in Amazon parlance); a Compute Unit is approximately equivalent to a 1.0–1.2 GHz 2007 Opteron or 2007 Xeon processor. The nodes are connected via a gigabit Ethernet. The binaries were compiled on Lawrencium. The results reported in [256] are summarized in Table 11.1.

The results in Table 11.1 give us some idea about the characteristics of scientific applications likely to run efficiently on the cloud. Communication intensive applications will be affected by the increased latency (more than 70 times larger than *Carver*) and lower bandwidth (more than 70 times smaller than *Carver*).

11.11 Cloud computing and biology research

Biology is one of the scientific fields that needs vast amounts of computing power and was one of the first to take advantage of cloud computing. Molecular dynamics computations are CPU-intensive, while protein alignment is data-intensive.

An experiment carried out by a group from Microsoft Research illustrates the importance of cloud computing for biology research [316]. The authors carried out an “all-by-all” comparison to identify the interrelationships of the 10 million protein sequences (4.2 GB size) in NCBI’s nonredundant protein database using AzureBLAST, a version of the BLAST²⁵ program running on the Azure platform [316].

Azure offers VM with four levels of computing power depending on the number of cores: small (1 core), medium (2 cores), large (8 cores), and extra large (> 8 cores). The experiment used 8 core CPUs with 14 GB RAM and a 2-TB local disk. It was estimated that the computation would take six to seven CPU-years; thus, the experiment was allocated 3 700 weighted instances or 475 extra-large VMs from three data centers; each data center hosted three AzureBLAST deployments, each with 62 extra large instances. The 10 million sequences were divided into multiple segments; each segment was submitted for execution by one AzureBLAST deployment. With this vast amount of resources allocated, it took 14 days to complete the computations which produced 260 GB of compressed data spread across over 400 000 output files.

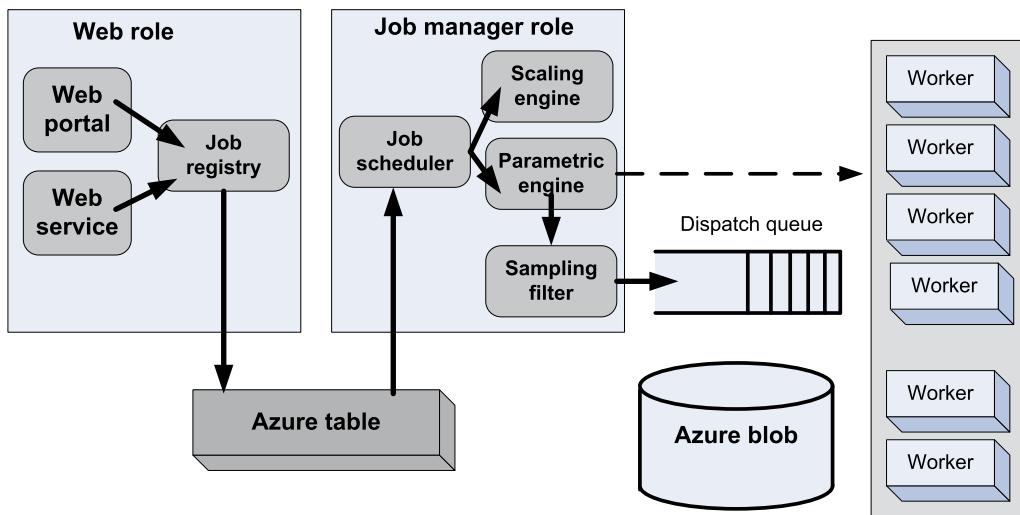
A post-experiment analysis led to a few conclusions useful for many scientific applications running on Azure. When a task runs for more than two hours, a message automatically reappears in the queue requesting the task to be scheduled, leading to repeated computations. The simple solution to this problem is to check if the result of a task has been generated before launching its execution. Many applications, including BLAST, allow for the setting of some parameters, but the computational effort for finding optimal parameters is prohibitive. To meet budget limitations each user is expected to decide on an optimal balance between cost and the number of instances.

A number of inefficiencies were observed: Many VMs were idle for extended periods of time. When a task finished execution, all worker instances waited for the next task. When all jobs use the same set of instances, resources are either under- or over-utilized. Load imbalance is another source of inefficiency; some of the tasks required by a job take considerably longer than others and delay the completion time of the job.

The analysis of the logs shows unrecoverable instance failures; some 50% of active instances lost connection to the storage service but were automatically recovered by the fabric controller. System updates caused several ensembles of instances to fail.

Another observation is that a computational science experiment requires the execution of several binaries, thus the creation of workflows, a challenging task for many domain scientists. To address this challenge, the authors of [303] developed a general platform for executing legacy Windows applications on the cloud. In the Cirrus system, a job has a description consisting of a prologue, a set of commands, and a set of parameters. The prologue sets up the running environment; the commands are sequences

²⁵ The Basic Local Alignment Search Tool (BLAST) finds regions of local similarity between sequences; it compares nucleotide or protein sequences to sequence databases and calculates the statistical significance of matches that can be used to infer functional and evolutionary relationships between sequences and to help identify members of gene families. More information is available at <http://blast.ncbi.nlm.nih.gov/Blast.cgi>.

**FIGURE 11.13**

Cirrus—a general platform for executing legacy Windows applications on the cloud.

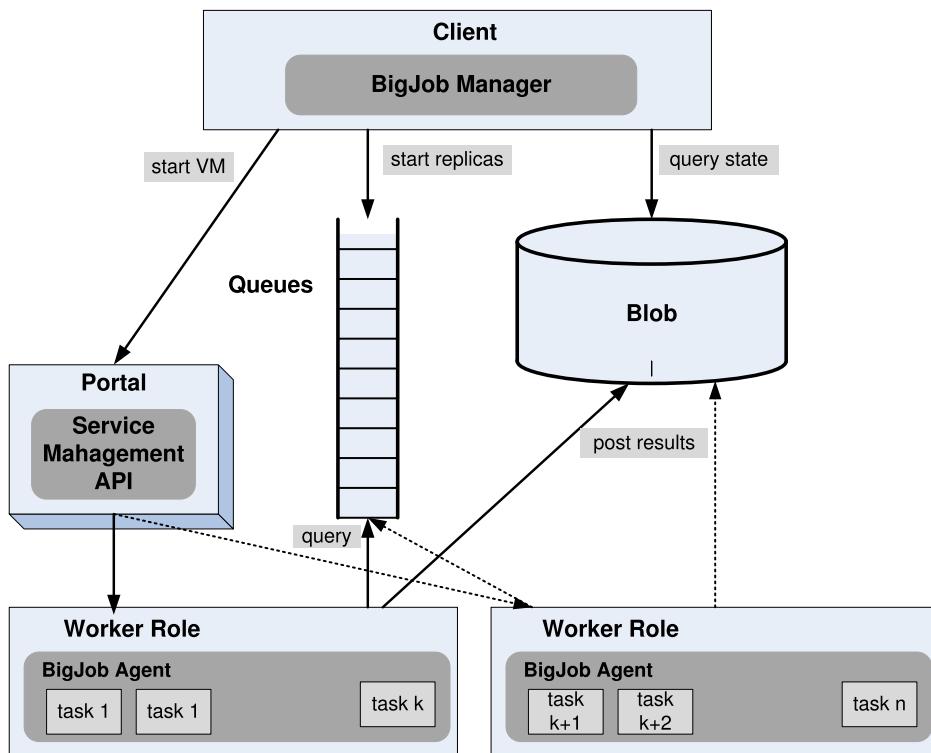
of shell scripts including Azure-storage-related commands to transfer data between Azure blob storage and the instance.

After the Windows Live ID service authenticates the user, it can submit and track a job through the portal provided by the web role; see Fig. 11.13; the job is added to a table called *job registry*. The execution of each job is controlled by a *job manager instance* that first scales the size of the worker based on the job configuration, and then, the parametric engine starts exploring the parameter space; if this is a test-run, the parameter sweeping result is sent to the sampling filter.

Each task is associated with a record in the task table, and this state record is updated periodically by the worker instance running the task; the progress of the task is monitored by the manager. The dispatch queue feeds into a set of worker instances. A worker periodically updates the task state in the task table and listens for any control signals from the manager.

A loosely coupled workload for an ensemble-based simulation on the Azure cloud is reported in [317]. A *role* in Azure is an encapsulation of an application; as noted earlier, there are two kinds of roles: (i) the web roles for web applications and front-end code; and (ii) the worker roles for background processing. Scientific applications, such as AzureBLAST, use worker roles for the compute tasks, and they implement their API that provides a run method and an entry point for the application and the state or configuration change notifications. The applications use the Blob Storage (ABS) for large raw data sets, the Table Storage (ATS) for semi-structured data, and the Queue Storage (AQS) for message queues; these services provide strong consistency guarantees, but the complexity is moved to the application space.

Fig. 11.14 illustrates the use of a software system called BigJob to decouple resource allocation from resource binding for the execution of loosely coupled workloads on an Azure platform [317]; this

**FIGURE 11.14**

The execution of loosely coupled workloads using the Azure platform.

software eliminates the need for the application to manage individual VMs. The results of measurements show a noticeable overhead for starting VMs and for launching the execution of an application task on a remote resource; increasing the computing power of the VM decreases the completion time for long-running tasks.

11.12 Social computing, digital content, and cloud computing

Social networks play an increasingly important role in people's lives; they have expanded in terms of the size of the population involved and function performed. A promising solution for analyzing large-scale social networks data is to distribute the computation workload over a large number of cloud servers. Traditionally, evaluating the importance of a node or a relationship in a network is done using sampling and surveying, but in a very large network, structural properties cannot be inferred by scaling up the results from small networks. The evaluation of social closeness is computationally intensive.

Social intelligence is another area where social and cloud computing intersect. Knowledge discovery and techniques based on pattern recognition demand high-performance computing and resources that can only be provided by computing clouds. Case-based reasoning (CBR), the process of solving new problems based on the solutions of similar past problems, is used by context-aware recommendation systems; it requires similarity-based retrieval. As the case base accumulates, such applications must handle massive amounts of history data, and this can be done by developing new reasoning platforms running on the cloud. CBR is preferable to rule-based recommendation systems for large-scale social intelligence applications. Indeed, the rules can be difficult to generalize or apply to some domains; all triggering conditions must be strictly satisfied, scalability is a challenge as data accumulate, and the systems are hard to maintain because new rules have to be added as the amount of data increases.

The *BetterLife 2.0* [246] a CBR-based system, consists of a cloud layer, a case-based reasoning engine, and an API. The cloud layer uses Hadoop clusters to store application data represented by cases, as well as social network information, such as relationship topology and pairwise social-closeness information. CBR engine calculates similarity measures between cases to retrieve the most similar ones and also stores new cases back to the cloud layer. The API connects to a master node that is responsible for handling user queries, distributes the queries to server machines, and receives results.

A case consists of a problem description, solution, and optional annotations about the path to derive the solution. CBR uses MapReduce; all the cases are grouped by their *userId*, and then a *breadth first search* (BFS) algorithm is applied to the graph where each node corresponds to one user. MapReduce is used to calculate the closeness according to pairwise relationship weight. A reasoning cycle has four steps: (a) retrieve the most relevant or similar cases from memory to solve the case; (b) reuse—map the solution from the prior case to the new problem; (c) revise—test the new solution in the real world or in a simulation and, if necessary, revise; and (d) retain—if the solution was adapted to the target problem, store the result as a new case.

In the past, social networks have been constructed for a specific application domain, e.g., MyExperiment and nanoHub for biology and nanoscience, respectively. These networks enable researchers to share data and provide a virtual environment supporting remote execution of workflows. Another form of social computing is the *volunteer computing* when a large population of users donate resources, such as CPU cycles and storage space for a specific project. The Mersenne Prime Search, initiated in 1996, followed in the late 1990s by the SETI@Home, the Folding@home, and the Storage@Home, a project to back up and share huge data sets from scientific research, are well-known examples of volunteer computing. Information about these projects is available online at: www.myExperiment.org, www.nanoHub.org, www.mersenne.org, setiathome.berkeley.edu, and at folding.stanford.edu.

Volunteer computing cannot be used for applications where users require some level of accountability. The PlanetLab project is a credit-based system in which users earn credits by contributing resources and then spend these credits when using other resources. Berkeley Open Infrastructure for Network Computing (BOINC) is the middleware for a distributed infrastructure suitable for multiple applications.

An architecture designed as a Facebook application for a social cloud is presented in [98]. Methods to get a range of data including friends, events, groups, application users, profile information, and photos are available through a Facebook API. The Facebook Markup Language is a subset of HTML with proprietary extensions, and Facebook JavaScript is a version of JavaScript. The prototype uses web services to create a distributed and decentralized infrastructure. There are numerous examples of cloud platforms for social networks. There are scalable cloud applications hosted by commercial clouds.

The new technologies supported by cloud computing favor the creation of digital content. *Data mashups* or *composite services* combine data extracted by different sources; *event-driven mashups*, also called Svc, interact through events rather than the request-response traditional method. A recent paper [454] argues that “the *mashup* and the cloud computing worlds are strictly related because very often the services combined to create new Mashups follow the SaaS model and more, in general, rely on cloud systems.” The paper also argues that the Mashup platforms rely on cloud computing systems, for example, the IBM Mashup Center and the JackBe Enterprise Mashup server.

There are numerous examples of monitoring, notification, presence, location, and map services based on the Svc approach, including: Monitor Mail, Monitor RSSFeed, Send SMS, Make Phone Call, GTalk, Fireeagle, and Google Maps. As an example, consider a service to send a phone call when a specific email is received; the Mail Monitor Svc uses input parameters such as: User ID, Sender Address Filter, Email Subject Filter, to identify an email and generate an event that triggers the *Make TTS Call* action of a *Text To Speech Call* Svc linked to it.

The system in [454] supports creation, deployment, activation, execution, and management of Event Driven Mashups; it has a user interface, a graphics tool called Service Creation Environment that supports easily the creation of new Mashups and a platform called *Mashup Container* that manages Mashup deployment and execution. The system consists of two subsystems, the *service execution platform* for Mashups execution and the *deployer* module that manages the installation of Mashups and Svcs. A new Mashup is created using the graphical development tool, and it is saved as an XML file; it can then be deployed into a *Mashup Container* following the Platform as a Service (PaaS) approach. The *Mashup Container* supports a primitive SLA allowing the delivery of different levels of service.

The prototype uses the JAVA Message Service (JMS), which supports an asynchronous communication; each component sends/receives messages, and the sender does not block waiting for the recipient to respond. The system’s fault tolerance was tested on a system based on the VMware vSphere. In this environment, the fault tolerance is provided transparently by the VMM, and neither the VMs nor the applications are aware of the fault tolerance mechanism; two VMs, a Primary and a Secondary one, run on distinct hosts and execute the same set of instructions such that, when the Primary fails, the Secondary continues the execution seamlessly.

11.13 Software fault isolation

Software fault isolation (SFI) offers a technical solution for sandboxing binary code of questionable provenance that can affect security in cloud computing. Insecure and tampered VM images comprise one of the security threats; binary codes of questionable provenance for native plugins to a web browser can pose a security threat because web browsers are used to access cloud services.

A recent paper [444] discusses the application of the sandboxing technology for two modern CPU architectures, ARM and x86-64. ARM is a load/store architecture with 32-bit instruction, 16 general-purpose registers. It tends to avoid multicycle instructions, and it shares many of the RISC architecture features, but: (a) it supports a “thumb” mode with 16-bit instruction extensions; (b) has complex addressing modes and a complex barrel shifter; and (c) condition codes can be used to predicate most instructions. In the x86-64 architecture, general-purpose registers are extended to 64-bits, with an *r* replacing the *e* to identify the 64 versus 32-bit registers, e.g., *rax* instead of *eax*; there are eight new

Table 11.2 The features of the SFI for the Native Client on the x86-32, x86-64 and ARM; ILP stands for Instruction Level Parallelism.

Feature/Architecture	x86-32	x86-64	ARM
Addressable memory	1 GB	4 GB	1 GB
Virtual base address	any	44GB	0
Data model	ILP32	ILP32	ILP 32
Reserved registers	0 of 8	1 of 16	0 of 16
Data address mask	None	Implicit in result width	Explicit instruction
Control address mask	Explicit instruction	Explicit instruction	Explicit instruction
Bundle size (bytes)	32	32	16
Data in text segment	forbidden	forbidden	allowed
Safe address registers	all	rsp, rbp	sp
Out-of-sandbox store	trap	wraps mod 4 GB	No effect
Out-of-sandbox jump	trap	wraps mod 4 GB	wraps mod 1 GB

general-purpose registers named $r8-r15$. To allow legacy instructions to use these additional registers, x86-64 defines a set of new prefix bytes to use for register selection.

This SFI implementation is based on the previous work of the same authors on Google Native Client (NC). This implementation assumes an execution model where a trusted runtime system shares a process with the untrusted multithreaded plugin. The rules for binary code generation of an untrusted plugin are:

1. The code section is read-only, and it is statically linked.
2. The code is divided into 32 byte *bundles*, and no instruction or pseudo-instruction crosses the bundle boundary.
3. The disassembly starting at the bundle boundary reaches all valid instructions.
4. All indirect flow control instructions are replaced by pseudo-instructions that ensure address alignment to bundle boundaries.

The features of the SFI for the Native Client on the *x86-32*, *x86-64*, and *ARM* are summarized in Table 11.2 [444]. The control flow and store sandboxing for the ARM SFI incur less than 5% average overhead, and the ones for x86-64 SFI incur less than 7% average overhead.

11.14 Further readings

MapReduce is discussed in [129], and [487] presents the GrepTheWeb application. Cloud applications in biology are analyzed in [316] and [317], and social applications of cloud computing are presented in [98], [246], and [454]. Benchmarking of cloud services is analyzed in [106], [256]. The use of structured data is covered in [321]. An extensive list of publications related to Divisible Load Theory is at <http://www.ece.sunysb.edu/~tom/dlt.html>.

There are several cluster programming models. Data flow models of MapReduce and Dryad [255] support a large collection of operations and share data through stable data. High-level programming languages, such as DryadLINQ [530] and FlumeJava [90], enable users to manipulate parallel collec-

tions of datasets using operators such as *map* and *join*. There are several systems providing high-level interfaces for specific applications such as *HaLoop* [77]. There are also caching systems including *Spark* [533] and *Tachyon* [302], discussed in Chapter 4, and *Nectar* [214].

Hive [475] was the first SQL over Hadoop to use another framework such as MapReduce or Tez to process SQL-like queries. *Shark* uses another framework, *Spark* [521] as its runtime. *Impala* [176] from Cloudera, LinkedIn *Tajo* (<http://tajo.incubator.apache.org/>), *MapR Drill* (<http://www.mapr.com/resources/community-resources/apache-drill>) and Facebook *Presto* (<http://prestodb.io/>) resemble parallel databases and use long-running custom-built processes to execute SQL queries in a distributed fashion. *Hadapt* [5] uses a relational database (PostgreSQL) to execute query fragments. Microsoft *PolyBase* [144] and Pivotal [97] use database query optimization and planning to schedule query fragments and read HDFS data into database workers for processing.

A discussion of cost-effective cloud-based high-performance computing and a comparison with supercomputers [478] is reported in [87]. [517] discusses scientific computing on clouds. Service level checking is analyzed in [99].

11.15 Exercises and problems

- Problem 1.** Download and install the *Zookeeper* from the site <http://zookeeper.apache.org/>. Use the API to create the basic workflow patterns shown in Fig. 11.3.
- Problem 2.** Use *AWS Simple Workflow Service* to create the basic workflow patterns in Fig. 11.3.
- Problem 3.** Use *AWS CloudFormation* service to create the basic workflow patterns in Fig. 11.3.
- Problem 4.** Define a set of keywords ordered based on their relevance to the topic of cloud security; then search the web using these keywords to locate 10–20 papers and store the papers in an *S3* bucket. Create a MapReduce application modeled after the one discussed in Section 11.6 to rank the papers based on the incidence of the relevant keywords. Compare your ranking with the rankings of the search engine you used to identify the papers.
- Problem 5.** Use the *AWS MapReduce* service to rank the papers in Problem 4.
- Problem 6.** The paper [84] describes the *elasticLM*, a commercial product that provides license and billing Web-based services. Analyze the merits and the shortcomings of the system.
- Problem 7.** Search the web for reports of cloud system failures and discuss the causes of each incident.
- Problem 8.** Identify a set of requirements you would like to be included in a SLA.
- Problem 9.** Consider the workflow for your favorite cloud application. Use XML to describe this workflow, including the instances and the storage required for each task. Translate this description into a file that can be used for the *Elastic Beanstalk AWS*.
- Problem 10.** In Section 11.10, we analyze cloud-computing benchmarks and compare them with the results of the same benchmarks performed on a supercomputer. This is not unexpected; discuss the reasons why we should expect the poor performance of fine-grained parallel computations on a cloud.
- Problem 11.** An IT company decides to provide free access to a public cloud dedicated to higher education. Which one of the three cloud computing delivery models, SaaS, PaaS, or IaaS, should it embrace and why? What applications would be most beneficial for the students? Will this solution have an impact on distance learning? Why?

Big Data, data streaming, and the mobile cloud 12

Advances in processor, storage, software, and networking technologies enable us to store and process massive amounts of data for the benefit of humanity, for profit, for entertainment, for nefarious schemes, or simply because we can. One can talk about “democratization of data” as scientists and decision makers, journalists and health care providers, artists and engineers, neophytes and domain experts attempt to extract knowledge from data.

This chapter covers three of the most exciting and demanding classes of cloud applications: Big Data, data streaming, and mobile cloud computing. Big Data is a reality; every day, we generate 2.5 quintillion, 2.5×10^{18} , bytes of data. This colossal volume of data is collected every day by devices ranging from inexpensive sensors in mobile phones to the detectors of the Large Hadron Collider, by online services offered by Google and Amazon, or by devices connected by the Internet of Things.

Big Data is a defining challenge for cloud computing; a significant amount of the data collected every day is stored and processed on computer clouds. Big Data and data streaming applications require low latency, scalability, versatility, and a high degree of fault tolerance. Achieving these qualities at scale is extremely challenging.

It makes sense to define a more comprehensive concept of scale for the applications discussed in this chapter. The “scale” in this context means millions of servers operating in concert in large data centers, a very large number of diverse applications, tens of millions of users, and a range of performance metrics reflecting the views of users, on one hand, and those of the CSPs, on the other hand. The scale has a *disruptive* effect: It changes how we design and engineer such systems and broadens the range of problems that can be solved and of the applications that can only run on computer clouds.

Scale amplifies the unanticipated benefits, as well as the dreaded nightmares, of system designers. Even a slight improvement of the individual server performance and/or of the algorithms for resource management could lead to huge cost savings and rave reviews. At the same time, the failure of one of the billions of hardware and software components can be amplified, propagate throughout the entire system, and have catastrophic consequences.

Several important lessons when engineering large-scale systems for Big Data storage and processing are: (a) Prepare for the unexpected because low probability events occur and can cause major disruptions; (b) It is utterly unreasonable to assume that strict performance guarantees can be offered at scale; and (c) It is unfeasible to build fault-free systems beyond a certain level of system complexity; understanding this truth motivated the next best solution, the development of design principle for fault-tolerant systems.

A realistic alternative to the applications discussed in this chapter is to develop tail-tolerant techniques for the distributions of performance metrics, such as response time latency. This means to understand why a performance metric has a heavy-tail distribution, detect the events leading to such an undesirable state as early as feasible, and take the necessary actions to limit its effects. Another design

principle is that a reasonably accurate result and fast response time are preferable to a delayed best result.

The defining attributes of Big Data are analyzed in Section 12.1. The next sections discuss how Big Data is stored and processed. High-capacity datastores and databases are necessary to store the very large volume of data. Scaling data warehouses and databases poses its own challenges. The Mesa datastore and Spanner and F1 databases developed at Google are discussed in Section 12.2, while another class of Big Data applications combining mathematical modeling with simulation and measurements, dynamic, data-driven applications (DDAS), are discussed in Section 12.3.

Clouds host several classes of data streaming applications, ranging from content delivery data streaming to applications consuming a continuous stream of events. Such applications are discussed in Sections 12.4, 12.5, and 12.6.

Mobile devices, such as smartphones, tablets, laptops, and wearable devices, are ubiquitous and indispensable for life in a modern society. Mobile devices are in a symbiotic relationship with computer clouds, and their users benefit from the democratization of data processing. The user of a mobile device has access to the vast amounts of computing cycles and storage available on computer clouds. Mobile cloud computing enables execution of mobile applications on mobile devices and on computer clouds. Mobile devices act as producers and consumers of data stored on clouds and shared with others.

Section 12.7 is an introduction to mobile computing and its applications, while Section 12.8 covers energy efficiency of mobile computing. Section 12.9 analyzes the effects of latency and presents alternative mobile computing models, including Cloudlets.

Scale allows us to add mission-critical applications demanding very high availability, a topic discussed in Section 12.10. Scale amplifies variability, often causing heavy-tail distributions of critical performance metrics as is the case of latency discussed in Section 12.11. Lastly, edge computing and Markov decision processes are analyzed in Section 12.12.

12.1 Big Data

Some of the defining characteristics of Big Data are the three “Vs,” volume, velocity, and variety, along with persistency. Volume is self-explanatory, and velocity means that responses to queries and data analysis requests have to be provided very fast. Variety recognizes the wide range of data sources and data formats. Persistency means that data has a lasting value; it is not ephemeral.

Big Data covers a wide spectrum of data including user-generated content and machine-generated data. Some of the data is highly structured, as in the case of patient records in health care, insurance claims, or mortgage documents. Others are raw data from sensors, log files, or data generated by social media.

Big Data has affected the organization of database systems. The traditional relational databases are unable to satisfy some of these requirements, and NoSQL databases proved to be better suited for many cloud applications. A database *schema* is a way to logically group objects such as tables, views, stored procedures, etc. A schema can be viewed as a container of objects. One can assign a user login permission to a single schema so that the user can only access the objects they are authorized to.

For decades, the database community used *schema-on-write*. First, one defines a schema, then writes the data, and data comes back according to the original schema when reading. The alternative, *schema-*

on-read loads data as-is, and a user-defined filter extracts data for processing. Schema-on-read has several advantages:

- Often, data is a shared asset among individuals with differing roles and differing interests who want to get different insights from that data. Schema-on-read can present data in a schema that is best adapted to the queries being issued.
- It is not necessary to develop a super-schema that covers all datasets when multiple data sets are consolidated.

An insightful discussion of the state of research on databases [3] starts by acknowledging that “Big Data requirements will cause massive disruptions to the ways that we design, build, and deploy data management solutions.” The report continues by identifying three major causes for these disruptions: “...it has become much cheaper to generate a wide variety of data, due to inexpensive storage, sensors, smart devices, social software, multiplayer games, and the emerging IoT ... it has become much cheaper to process large amounts of data, due to advances in multicore CPUs, solid state storage, inexpensive cloud computing, and open source software ... not just database administrators and developers, but many more types of people have become intimately involved in the process of generating, processing, and consuming data...”

Big Data revolutionized computing. Chapter 11 presented specialized frameworks for Big Data processing such as MapReduce and Hadoop. The myriad system software components reviewed in Chapter 4, including Pig, Hive, Spark, and Impala, are essential elements of a more effective infrastructure for processing unstructured or semi-structured data. The evidence of the effort to mold the cloud infrastructure to the requirements of Big Data is overwhelming. In spite of their diversity Big Data workloads have a few common traits:

- Data is immutable; consequently, widely used storage systems for Big Data such as HDFS only allow *append* operations.
- Jobs such as MapReduce are deterministic; therefore fault tolerance can be ensured by recomputations.
- The same operations are carried out on different segments of data on different servers; replicating programs is less expensive than replicating data.
- It is feasible to identify a *working set*, a subset of data frequently used within a time window and keep this working set in the memory of the servers of a large cluster. Systems such as Spark [533] and Tachyon [302] exploit this idea to dramatically improve performance.
- *Locality*, the proximity of the data to the location where it is processed is critical for the performance. Supporting data locality is a main objective of scheduling algorithms, including delay scheduling [532] and data-aware scheduling [492], as we have seen in Chapter 9.

The cloud hardware infrastructure also faces a number of challenges. Bridging the gap between the processor speed and the communication latency and bandwidth is a major challenge. This challenge is greatly amplified by response-time constraints when Big Data is processed on the cloud. Addressing this challenge requires prompt adoption of faster networks, such as InfiniBand, and of full bisection bandwidth networks between servers. Remote direct memory access capabilities can also help bridge this gap.

Nonvolatile random-access memories, specialized processors such as GPUs and field-programmable gate arrays (FPGAs), and application-specific integrated circuits (ASICs) contribute to the effort to de-

Table 12.1 Capacity and bandwidth of hard disks (HDD), solid-state disks (SSD), and memory of a modern server according to <http://www.dell.com/us/business/p/servers>. The network bandwidth is 1.25 GB/sec.

Media	Capacity (TB)	Bandwidth GB/sec
HDD (x12)	12–36	0.2–2
SDD (x4)	1–4	1–4
Memory	0.128–0.512	10–100

velop a scalable infrastructure. The storage technology has dramatically improved, as summarized in Table 12.1.

Some of the enduring challenges posed by Big Data are:

1. Develop a scalable data infrastructures capable to respond promptly to the timing constraints of an application.
2. Develop effective means to accommodate diversity in data management systems.
3. Support comprehensive end-to-end processing and knowledge extraction from data.
4. Develop convenient interfaces for a layman involved in data collection and data analysis.

The next sections address the scalability challenges faced by the development of large data storage systems and by processing very large volumes of data.

12.2 Data warehouses and Google databases for Big Data

A *data warehouse* is a central repository for the important data of an enterprise; it is a core enterprise software component required by a range of applications related to so-called *business intelligence*. Data from multiple operational systems are uploaded and used for predictive analysis and for the detection of hidden patterns in the data. Data models developed using Statistical Learning Theory allow enterprises to optimize their operations and maximize their profits.

Scaling data repositories is very challenging due to the low-latency, scalability, versatility, availability, and fault-tolerance requirements. Google realized early on the need to develop a coherent storage architecture for the massive amounts of data required by their cloud services and by AdWords.¹ This storage architecture reflects the realization that “developers do not ask simple questions of the data, change their data access patterns frequently and use APIs that hide storage requests while expecting uniformity of performance, strong availability and consistent operations, and visibility into distributed storage requests” [173].

Two of the most popular Google data stores, BigTable and Megastore, discussed in Sections 7.11 and 7.12, respectively, were designed and developed early on to support Google’s cloud computing services. A major complaint about BigTable is the lack of support for cross-row transactions. Megas-

¹ AdWords consists of hundreds of applications supporting Google’s advertising services.

tore supports schematized semi-relational tables and synchronous replication. At least 300 applications within Google use Megastore, including Gmail, Picasa, Calendar, Android Market, and App Engine.

From the experience with the BigTable storage system, the developers of Google's storage software stack learned that it is hard to share distributed storage. Another lesson is that distributed transactions are the only realistic option to guarantee the low latency required for a high-volume transaction processing system. They also learned that end-user latency matters and that the application complexity is reduced if an application is close to its data [173]. Another important lesson when developing a large-scale system is that making strong assumptions about applications is not advisable.

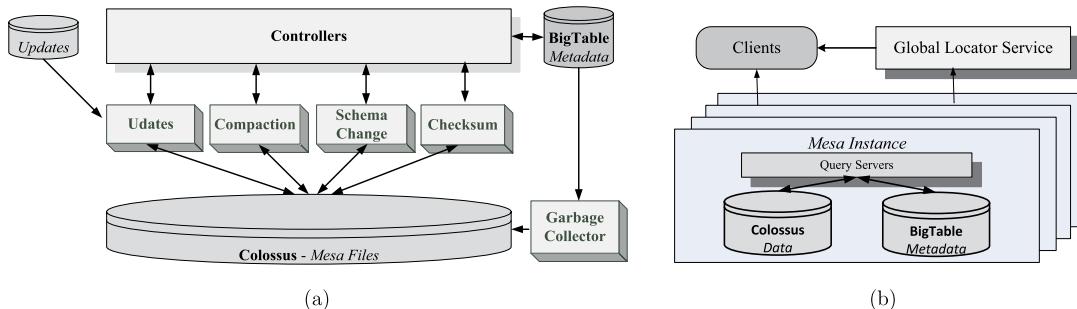
These lessons led to the development of Colossus, the successor of the GFS as a next-generation cluster-level file system, the Mesa warehouse, and Spanner and F1 databases, discussed in this section. Colossus is built for real-time services and supports virtually all web services, from Gmail, Google Docs, and YouTube, to Google Cloud Storage service offered to third-party developers. The system allows client-driven replication and encoding. The automatically shared metadata layer of Colossus enables availability analysis. To reduce cost Colossus data is typically encoded with Reed–Solomon codes.

Mesa—a scalable data warehouse. The system developed at Google [217] is an example of a data warehouse designed to support measurement data for the multibillion-dollar advertising business of the organization. The system is expected to support near-real-time data processing and be highly available and scalable. The extremely high availability is ensured by geo-replication.² Mesa is able to handle petabytes of data, respond to billions of queries accessing trillions of rows of data per day, and update millions of rows per second. The system complexity reflects the stringent requirements including the support for:

- Complex queries such as “How many ad clicks were there for a particular advertiser matching the keyword “fig” during the first week of December between 11:00 AM and 2:00 PM that were displayed on google.com for users in a specific geographic location using a mobile device?” [217].
- Multidimensional data with *dimensional* attributes, called keys, and *measure* attributes, called *values*.
- Atomic updates, consistency, and correctness. Multiple data views defined for different performance metrics are affected by a single user action and all must be consistent. The BigTable storage system does not support atomicity, while the Megastore system provides consistency across geo-replicated data; see Sections 7.11 and 7.12, respectively.
- Availability - planned downtime is not allowed, and unplanned downtime should never be experienced.
- Scalability - accommodates a very large volume of data and a large user population.
- Near-real-time performance - support live customer queries, reports, and updates; allows queries to multiple data views from multiple data centers.
- Flexibility and ability to support new features.

Logical and physical data organization in Mesa. The system stores data using tables with a very large key, K , and value, V , spaces. These spaces are represented by tuples of columns of items of identical

² The term *geo-replication* used in the title of the paper [217] means that the Mesa system runs at multiple sites concurrently. The term used in [218] for this strategy supporting high availability is *multihoming*.

**FIGURE 12.1**

Mesa instance. (a) The controller/worker subsystem. Controllers interact with four types of workers: update, compaction, schema change, and checksum. The data and metadata are stored on Colossus and BigTable, respectively. (b) The query subsystem of a Mesa instance interacts with the clients and with the Global Location Service.

data type, e.g., integers, strings, or floating-point numbers. The data is horizontally partitioned and replicated.

The table structure and the *aggregation function* $F : V \times V \mapsto V$ are specified by a *table schema*. Function F is associative and often commutative. To maximize the throughput, updates, each consisting of at most one aggregated value for every (table name, key) pair, are applied in batches.

Updates are applied by *version number*, are atomic, and the next update can only be applied after the previous one has finished. The time associated with a version is the time when the version was generated. A query has also a version number and a predicate P on key K .

A *delta* is a pre-aggregation of versioned data consisting of a set of rows corresponding to the set of keys for a range of versions $[V_1, V_2]$ with $V_1 \leq V_2$. Deltas can be aggregated by merging row keys and aggregating values accordingly, e.g.,

$$[V_1, V_2] \& [V_2 + 1, V_3] \rightarrow [V_1, V_3]. \quad (12.1)$$

Deltas are immutable, and the rows in one are stored in sorted order in files of limited size. A row consists of multiple row blocks; each row block is transposed and compressed. Each table has one or more *table indexes*, and each table index has a copy of the data sorted according to the index order.

Mesa limits the time a version can be queried to reduce the amount of space. Older versions can be deleted, and queries for such versions are rejected. For example, updates can be aggregated into base $B \geq 0$ with version $[0, B]$, and any updates with $[V_1, V_2]$ with $0 \leq V_1 \leq V_2 \leq B$ can be deleted.

Mesa instances. One Mesa instance runs at every site and consists of two subsystems, the *update/maintenance* and the *query* subsystem shown in Fig. 12.1. The pool of workers of the first subsystem operate on data stored in *Colossus*. Workers load updates, carry out table compaction, apply schema changes, and run table checksums under supervision of *controllers* that determine the work to be done and manage metadata stored on *BigTable*.

To scale, the controller is *sharded* by table; a shard is a horizontal partition of data in a data store and each shard is held on a separate physical storage device. The controller maintains separate queues of work for each type of worker and redistributes the workload of a slow worker to another worker in

the same pool. Workers poll the controller for additional work when idle and notify the controller upon completion of a task. Work requiring global coordination, including schema change and checksums, is initiated by components outside the controller.

The query subsystem includes a pool of query servers that process client queries. A query server looks up BigTable for metadata, determines the files where the data are stored, performs on-the-fly aggregation of data, and converts data from the internal format to the client protocol format. Multiple queries acting upon the same tables are assigned to a group of servers to reduce access time and optimize the system performance.

Mesa instances are running at multiple sites to support a high level of availability. A *committer* coordinates updates at all sites. The committer assigns a new version number to batches of updates and publishes the metadata for the update to the *versions* database, a globally replicated and consistent data store using the Paxos algorithm. Controllers detect the availability of new updates by listening to the changes in the *versions* database and then assign the work to update workers.

Mesa reads 30 to 60 MB compressed data, updates 3 to 6 million distinct rows, and adds some 3×10^5 thousand new rows per second. It executes more than 500 million queries and returns $(1.7\text{--}3.2) \times 10^{12}$ rows per day. Updates arrive in batches about every five minutes, with median and 95th-percentile commit times of 54 seconds and 211 seconds, respectively.

Spanner—a globally distributed database. Scaling traditional databases is not without major challenges. Spanner [117] is a distributed database replicating data on many sites across the globe. Some applications replicate their data across three to five data centers running Spanner in one geographic area. Other applications spread their data over a much larger area, e.g., F1, presented later in this section, maintains five replicas of data around the US. Spanner has been used by many Google applications since its release in 2011. The first user of Spanner was the F1 system.

The data model of each application using Spanner is layered on top of the directory-bucketed key-value mappings supported by the distributed database. Spanner supports consistent backups, atomic schema updates, and consistent MapReduce execution and provides externally consistent reads and writes and globally consistent reads across the database at a time stamp. This is possible because Spanner assigns globally meaningful commit time stamps to transactions reflecting the serialization order, even though transactions may be distributed. Serialization order satisfies *external consistency*: The commit time stamp of transaction T_1 is lower than that of transaction T_2 when T_1 commits before T_2 starts.

Spanner applications can control the number and the location of data centers where the replicas reside, as well as read and write latency. The read and write latencies are determined by how far data is from its users and, respectively, how far are the replicas from one another. The system is organized in *zones*, units of administrative deployment and of physical isolation similar with AWS zones.

Spanner organization. A *zonemaster* is in charge of several thousand *spanservers* in each zone. The clients use *location proxies* to find the spanservers able to serve their data. A *placement driver* handles the migration of data across zones with latencies of minutes, and the *universe master* maintains the state of all zones. A spanserver serves data to the clients. The spanserver implements a Paxos state machine per tablet and manages 100–1 000 tablets. A *tablet* is a data structure implementing a bag of the mapping

$$(key : string, aimestamp : int64) \rightarrow string. \quad (12.2)$$

Table 12.2 Spanner latency and throughput for write, read-only, and snapshot read. The mean and standard deviation over 10 runs reported in [117].

<i>Replicas</i>	Latency in ms.			Throughput in Kops/sec.		
	<i>Write</i>	<i>Read-only</i>	<i>Snapshot read</i>	<i>Write</i>	<i>Read-only</i>	<i>Snapshot read</i>
1	14.4 ± 1.0	1.4 ± 0.1	1.3 ± 0.1	4.1 ± 0.5	10.9 ± 0.4	13.5 ± 0.1
3	13.9 ± 0.6	1.3 ± 0.1	1.2 ± 0.1	2.2 ± 0.5	13.8 ± 3.2	38.5 ± 0.3
5	14.4 ± 0.4	1.4 ± 0.05	1.3 ± 0.4	2.8 ± 0.3	25.3 ± 5.2	50.0 ± 1.1

A great deal of attention is paid to replication and concurrency control. A set of replicas is called a *Paxos group*. The system administrators control the number and types of replicas and the geographic placement of them. Applications control the manner the data is replicated.

Each spanserver implements a *lock table* for concurrency control. Every replica has a leader. Every Paxos write is logged twice, once in the tablet's log and once in the Paxos log. The $(key, value)$ mapping state is stored in the tablet. At each site, the local spanserver software stack includes a site *participant leader* communicating with its peers at other sites. This leader controls a *transaction manager* and manages the lock table.

The system stores all tablets in Colossus. The implementation of the Paxos algorithm is optimized. The algorithm is pipelined to reduce latency. The leaders of the Paxos algorithm discussed in Section 10.13 are long-lived, about 10 seconds.

Spanner *directories*, also called *buckets*, are sets of contiguous keys sharing a common prefix. A directory is the smallest data unit whose placement can be specified by an application. A background task called *moved* moves data directory-by-directory between the Paxos groups. This task is also used to add or remove replicas to/from Paxos groups.

Spanner transactions. The database supports read-write transactions, read-only transactions, and snapshot reads. A *read-write* transaction implements a standalone *write*, and the concurrency control is pessimistic. A *read-only* transaction benefits from snapshot isolation.³ A *snapshot read* is a lock-free *read* in the past, at a time stamp specified by the client or at a time stamp chosen by the system before a time stamp upper bound specified by the client.

The system supports atomic schema-change transactions. The transaction is assigned a time stamp t in the future. The time stamp is registered during the prepare phase so that the schema changes on thousands of servers can complete with minimal disruption to other concurrent activity. Reads and writes, implicitly depending on the schema, are synchronize with any registered schema-change and may proceed if their time stamps precede time t ; otherwise, they must block until schema changes.

Some of the results of a micro-benchmark reported in [117] are shown in Table 12.2. The data was collected on time-shared systems with spanservers running in each zone on 4 cores AMD Barcelona 2200MHz servers with 4GB RAM. The two-phase commit scalability was also assessed: it increases from 17.0 ± 1.4 ms for one participant to 30.0 ± 3.7 for 10 participants, to 71.4 ± 7.6 for 100, and 150.5 ± 11.0 ms for 200.

TrueTime. A two-phase locking is used for transactional reads and writes. An elaborate process for time stamp management based on the TrueTime is used. The *TrueTime* API enables the system to

³ Snapshot isolation is a guarantee that all reads made in a transaction see a consistent snapshot of the database.

support consistent backups, atomic schema updates, and other desirable features. This API represents time as a *TTinterval* with the starting and ending times of type *TTstamp*. A *TTinterval* has bounded time uncertainty. Three methods using *t* of type *TTstamp* as argument are supported

TT.now()—returns a *TTinterval* : $[earliest, latest]$.

TT.after(t)—returns *true* if time *t* has definitely past.

TT.before(t)—returns *true* if time *t* has definitely not arrived.

TrueTime guarantees that for an invocation

$$tt = TT.now(), tt.earliest \leq t_{abs}(e_{now}) \leq tt.latest, \quad (12.3)$$

where *e_{now}* is the invocation event. Each data center will have one *time master*, and each server will have its own *timeslave* daemon interacting with several time masters to reduce the probability of errors.

TrueTime guarantees correctness of concurrency control and supports externally consistent transactions, lock-free read-only transactions, and nonblocking reads in the past. It guarantees that a whole-database audit at time *t* sees the effects of every transaction committed up to time *t*. The results show that refining clock uncertainty associated with jitter or skew allows building distributed systems with much stronger time semantics.

F1-scaling a traditional SQL database. A traditional database design shares the goals with the design of data warehouses, i.e., scalability, availability, consistency, usability, and latency hiding. Google system developers realized that scaling up their sharded MySQL implementation supporting online-transaction processing (OLTP) and online analytical processing (OLAP) systems was not feasible. Instead, a distributed SQL database called F1 was developed. F1 uses Spanner and stores data on Colossus File System (CFS).

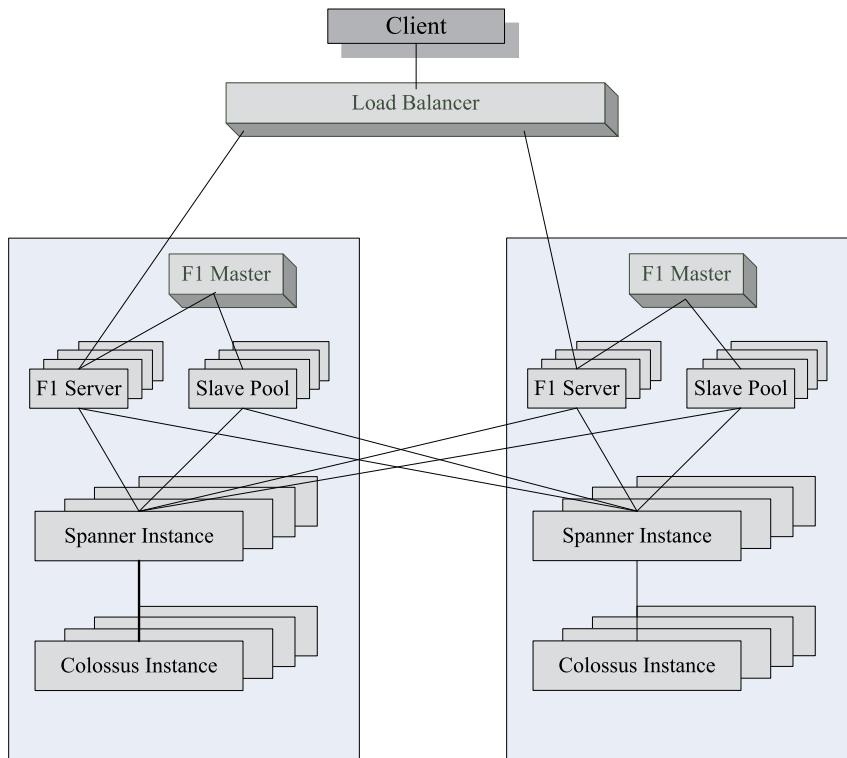
The F1 database [447] has been used since 2012 for the AdWords advertising ecosystem at Google. F1 inherits Spanner’s scalability, synchronous replication, strong consistency, and ordering properties and adds to them: (1) distributed SQL queries; (2) secondary indexes transactionally consistent; (3) asynchronous schema changes; (4) optimistic transactions; and (5) automatic change history recording and publishing.

Fig. 12.2 illustrates the F1 organization. Users interact with F1 using a *client library*. F1 servers are colocated in each data centers along with Spanner servers. Multiple Spanner instances run at each site along with multiple CFS instances. CFS is not a globally replicated service, and Spanner instances communicate only with local CFS instances. Storing data locally reduces the latency. The system is scalable, and its throughput can be increased by adding additional F1 and Spanner servers. The commit latencies are in the range 50–150 ms due to synchronous data replication across multiple data centers.

F1 has a logical Relational Database Management System schema with some extensions, including explicit table hierarchy and columns with Protocol Buffer data types. F1 stores each child table clustered with, and interleaved within, the rows from its parent table. Table columns contain structured data types based on the schema and binary encoding format of Google’s open-source Protocol Buffer library.

F1 supports *nonblocking schema changes* in spite of several challenges including:

- The scale of the system.
- High availability and tight latency constraints.
- The requirement to continue queries and transaction processing while schemas change.

**FIGURE 12.2**

F1 architecture. A load balancer distributes the workload to F1 server instances at every site which, in turn, interact with Spanner instances at all sites where F1 is running. F1 data is stored on Colossus at the same site. The shared slave pool execute parts of distributed query plans on behalf of regular F1 servers. F1 master monitors the health of slave pool processes and communicates to F1 server the list of available slaves.

- Each F1 servers has a copy of the schema in local memory for efficiency reasons; thus it is impractical to require atomical schema updates for all servers.

The schema change algorithm requires that at most two different schema are active at any time; one can be the current schema, the other the next schema. A server cannot use a schema after its lease expires. A schema change is divided in multiple phases such that consecutive pairs of phases are mutually compatible.

F1 transactions. F1 supports ACID transactions discussed in Section 7.3 and is required by systems such as financial systems and AdWords. The three types of F1 transactions built on top of Spanner transactions are:

1. Snapshot transactions—read-only transactions with snapshot semantics used by SQL queries and by MapReduce. Snapshot isolation is a guarantee that all reads made in a transaction see a consistent

snapshot of the database. The transaction itself will successfully commit only if no updates it has made conflict with any concurrent updates made since that snapshot.

2. Pessimistic transactions—map to the same Spanner transaction type.
3. Optimistic transactions—have an arbitrarily long lockless read phase, followed by a short write phase, and are the default transactions used by the clients.

Optimistic transactions have a number of benefits:

- Are long-lasting.
- Can be retried transparently by an F1 server.
- The state is kept on the client and is immune to server failure.
- The time stamp for a read can be used by a client for a speculative write that can only succeed if no other writes occurred after the read.

Optimistic transactions have two drawbacks, insertion phantoms and low throughput for high contention. Insertion phantoms occur when one transaction selects a set of rows, another transaction inserts rows that meet the same criteria, and, then, different results are produced when the first transaction re-executes the query.

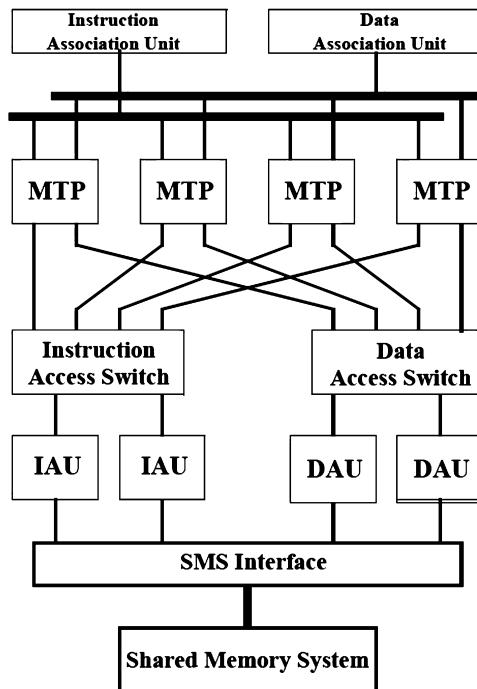
The default locking in F1 is at the row-level, though concurrency levels can be changed in the schema. By default, tables are change-tracked unless the schema shows that the some tables or columns have opted out. Every transaction creates one or more ChangeBatch Protocol Buffers, including the primary key and before and after values of changed columns for each updated row.

12.3 Dynamic data-driven applications

There are many applications, some involving Big Data, which combine mathematical modeling with simulation and measurements. For example, *large-scale-dynamic-data* refers to data captured by sensing instruments and control in engineered, natural, and societal systems. Some sensing instruments capture very large volumes of data; this is the case of the Large Hadron Collider at CERN, discussed in Section 3.15. A special case of such systems, Dynamic Data-Driven Application Systems (DDDAS), is discussed in this section. Such applications can benefit from dataflow architectures and programming models, such as FreshBreeze, developed by Professor Jack Dennis at MIT [140,141,305].

Dynamic Data-Driven Application Systems. Efforts to analyze, understand, and predict the behavior of complex systems often require a dynamic feedback loop. Sometimes, the simulation of a system uses measurement data to refine the system model through successive iterations or to speedup the simulation. Alternatively, the accuracy of a measurement process can be iteratively improved by feeding the data to the system model and then using the results to control the measurement process. In all these cases, the computation side and the instrumentation have to work in concert. The mathematical modeling, the simulation algorithms, the control system, the sensors, and the computing infrastructure of a DDDAS system should support an optimal logical and physical dataflow through a feedback loop.

Some DDDAS applications discussed in <http://www.1dddas.org/activities/2016-pi-meeting> are: adaptive stream mining, modeling of nanoparticle self-assembly processes, dynamic integration of motion and neural data to capture human behavior, real-time assessment and control of electric microgrids, optimization, and health monitoring of large-scale structural systems.

**FIGURE 12.3**

A FreshBreeze system consists of multithreaded processors (MTPs), a shared memory system (SMS), instruction and data access units (IAUs) and (DAUs), and instruction and data switches [140].

The FreshBreeze multiprocessor architecture and FreshBreeze model of thread execution. A chip architecture guided by modular programming principles described in [140] is well-suited for DDDAS applications, as we shall see in the discussion of the Mahali project. The system is designed to satisfy the software design principles discussed in Section 3.11. One of the distinctive features of the architecture is to prohibit memory updates and use a cycle-free heap. Objects retrieved from memory are immutable, and the allocation, release, and garbage collection of fixed-size chunks of memory are implemented by hardware mechanisms. This eliminates entirely the challenges posed by the cache coherence.

System organization is depicted in Fig. 12.3. The system consists of several multithreaded scalar processors (MTPs), a shared memory system (SMS) organized as a collection of fixed-size chunks of 1 024 bits, and an interconnection network. An MTP supports up to four execution threads involving integer and floating-point operations. The MTPs communicate with Instruction Access Units (IAUs) and Data Access Units (DAUs) to access chunks containing instructions and data, respectively. The access units have multiple slots of size equal to the size of memory chunks and maintain chunk usage data used by LRU algorithms for purging data to the SMS. The SMS performs automatically all memory updates including garbage collection.

An array is represented by a tree of chunks with element values at level 0. The number of tree levels is determined by the largest index value with a defined element, up to a maximum of eight levels. A collection of chunks containing pointers to other chunks forms a *heap*. The FreshBreeze execution model allows only the creation of *cycle-free* heaps. Many applications use also stream data types, unending series of values of uniform type. A Fresh Breeze stream is represented by a sequence of chunks, each chunk containing a group of stream elements. A chunk may include a reference to the chunk holding the next group of stream elements. As an alternative, an auxiliary chunk may contain “current” and “next” references.

The FreshBreeze execution model is discussed in [141]. A *master thread* spawns *slave threads* and initializes a *join point* providing the slave with a join ticket, similar to the return address of a method. Any slave thread may be a master for a group of slaves. The master does not continue after spawning the slaves, and there is no interaction between the master and the slaves or among the slaves other than the contribution of each to the continuation thread. A program can generate an arbitrary hierarchy of concurrent threads corresponding to the available application parallelism.

A join point is a special entry in the local data segment of the master thread that provides space for a record of master thread status and for the result produced by the slave. Only one slave can be given a join ticket to the same join point to avoid race conditions. Several instructions are used to access a join point:

1. *Spawn*—sets a flag, stores a joint ticket in the slave’s local data segment, and starts slave execution.
2. *EnterJoinResult*—allows the slave to save the result in the join point and then quit.
3. *ReadJoinValue*—returns the join value if it is available or suspends the master if the join value has not yet been entered.

The operation of a thread is deterministic because any heap operation either reads data or creates private data, and operations at a join point are independent of the order of slave thread arrival.

The parallelization of a dot product of two vectors discussed next illustrates an application of the FreshBreeze execution model [305]. First, the vectors are converted to tree-based memory chunks. The vectors are split into 16-element segments and organized as a tree structure. The leaf chunks hold the actual values, while the internal nodes of the tree store chunks holding handles of chunks that point to other chunks.

The *TraverseVector* thread takes the roots of two trees-of-chunks of vectors *A* and *B* as inputs and checks the depth of the tree to see if it is a leaf node. If not a leaf node, then it recursively spawns *TraverseVector* threads taking the root handles of the next level as inputs. At the leaf level, a *Compute* thread is spawned to compute the dot product of 16 elements and return the result sum to the *Sync* chunk. The continuation *Reduce* thread adds all partial results in the *Sync* chunk filled by lower level *Compute/Reduce* threads. Then, the *Reduce* thread returns the handle of the *Sync* chunk to the upper level *Sync* chunk until reaching the root level *Sync* chunk as the final result.

Space Weather Monitoring. The Mahali⁴ space weather monitoring project at MIT captures the quintessence of DDDAS applications. The project uses multicore mobile devices, such as phones and tablets, to form a global space weather monitoring network. Multifrequency GPS sensor data is processed on a cloud to reconstruct the structure of the space environment and its dynamic changes.

⁴ “Kila Mahali” means “everywhere” in the Swahili language; see <https://mahali.mit.edu/>.

The core ideas of the project are: (1) leverage the entire ionosphere as a sensor for ground-based and space-based phenomena; and (2) take advantage of mobile technology as a game changer for observatories.

12.4 Data streaming

Data streaming is *the transfer of data at a steady high-speed rate, with low and well-controlled latency*. The data volume in data streaming is high and decisions have to be made in real-time. High-definition television (HDTV) is a ubiquitous application of data streaming.

Cloud data streaming services support content distribution from many organizations, e.g., AWS hosts Netflix and Google cloud hosts YouTube. Clouds support a variety of other data streaming applications in addition to hosting content providers. For example, AWS supports several streaming data platforms, including Apache Kafka, Apache Flume, Apache Spark Streaming, and Apache Storm.

Data streaming versus batch processing. According to AWS: “Streaming data ... is generated continuously by thousands of data sources, which typically send in the data records simultaneously, and in small sizes (order of Kilobytes). Streaming data includes a wide variety of data such as log files generated by customers using your mobile or web applications, ecommerce purchases, in-game player activity, information from social networks, financial trading floors, or geospatial services, and telemetry from connected devices or instrumentation in data centers. This data needs to be processed sequentially and incrementally on a record-by-record basis or over sliding time windows, and used for a wide variety of analytics including correlations, aggregations, filtering, and sampling.”

There are important differences between streaming and batch data processing:

1. Streaming processes individual records or microbatches, rather than large batches.
2. Streaming processes only the most recent data or data over a rolling time window, rather than the entire, or a large segment of a, data set.
3. Streaming requires latency of milliseconds, rather than minutes or hours.
4. Streaming provides simple response functions, aggregates, and rolling metrics, rather than carrying out complex analytics.
5. It can be hard to reason about the global state in data streaming because different nodes may be processing data that arrived at different times, while in batch processing the system state is well defined and can be used to checkpoint and later restart the computation.

These differences suggest that the data streaming programming models and the APIs will be different from the ones for batch processing discussed in Chapters 4 and 11. It is therefore necessary to develop new data processing models for cloud data streaming services. Such models should be simple, yet effective.

Spark Streaming. Spark Streaming system along with a data streaming model, called D-Streams, are discussed in [534]. D-Streams model offers a high-level functional programming API, strong consistency, and efficient fault recovery. This model, prototyped as Spark Streaming, achieves fault tolerance by replication; there are two processing nodes, the upstream backup and a downstream node.

To address the latency concerns, the Spark Streaming system relies on a storage abstraction, *Resilient Distributed Dataset* (RDD), to rebuild lost data without replication. A *parallel recovery mech-*

anism involve all cluster nodes working in concert to reconstruct lost data. The system divides time in short intervals, stores the input data received during each interval in RDD, and then processes data via deterministic parallel computations that may involve MapReduce and other frameworks. A D-Stream is a collection of RDDs.

Spark Streaming provides an API similar to DryadLINQ in the Scala language and supports stateless operations acting independently in each time interval, as well as aggregation over time window. To recover in the case of node failures, D-Streams and RDDs track their lineage using a dependency graph. The system supports:

1. Stateless transformations available in batch frameworks, such as *map*, *reduce*, *groupBy*, and *join*, and provides two types of operators: (i) stateless and statefull transform operators, the first act independently on each time interval and the second share data among intervals; and (ii) output operators that save data, e.g., store RDDs on HDFS.
2. New statefull operations: (i) windowing - a window groups records from a range of past intervals into one RDD; (ii) incremental aggregation over a window; and (iii) time-skewed joins when a stream is combined with RDDs.

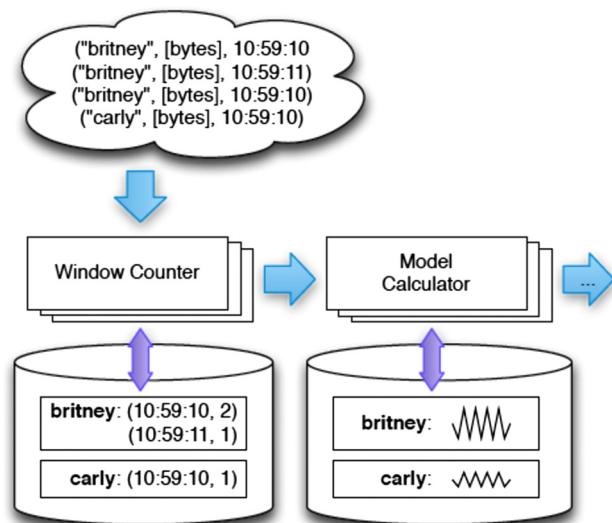
Spark Streaming uses an inefficient upstream backup approach, does not support finer checkpointing because it creates checkpoints every minute, and depends on application idempotency and system slack for recovery.

Zeitgeist and MillWheel. Google's Zeitgeist,⁵ a system used to track trends in web queries, is a typical application of data streaming. The system builds a historical model of each query and continually identifies queries that spike or dip. Zeitgeist processing pipeline includes a window counter with query searches as input, a model calculator and a spike/dip detector, and an anomaly notification engine. System buckets records arriving in one-second intervals and then compares the actual traffic for each time bucket to the expected traffic predicted by the model. An example of Zeitgeist aggregation of one-second buckets of (key, value, time stamp) query triplets given in [15] is shown in Fig. 12.4.

Google MillWheel framework developed for building fault-tolerant and scalable data streaming systems supports persistent state [15] and addresses some of the limitations of Spark Streaming and other data streaming systems. Zeitgeist system provided the initial requirements for MillWheel stream processing services. Some of these requirements are:

- Immediate data availability for data processing services.
- Exposure of persistent state abstraction to user applications.
- Graceful handling of out-of-order data.
- Exactly once delivery of records.
- Constant latency as the system scales up.
- Monotonically increasing low watermarking of data time stamps.

⁵ Zeitgeist is a German word literary translated as *time mind* and used with the sense of *the spirit of the time*. This concept, attributed to Georg Wilhelm Friedrich Hegel, an important figure of the German idealism philosophy school, constitutes the dominant ideas in the society at a particular time.

**FIGURE 12.4**

Zeitgeist aggregation of one-second buckets of (key, value, time stamp) query triplets [15].

MillWheel users express the logic of their data streaming applications as a directed graph where the stream of data records is delivered along the graph edges. A node or an edge in the arbitrary topology can fail at any time without affecting the correctness of the result. Record delivery is idempotent.⁶

The data structures used by MillWheel as input and output are (key, value, time stamp) triplets. The *key* and the *value* can be any string. The *time stamp* is typically close to the wall time of the event, though it can be assigned an arbitrary value by MillWheel. Input data trigger computations invoking either user-defined actions, or MillWheel primitives.

The key extraction function uses a user-assigned key for the aggregation and the comparison among records. For applications such as Zeitgeist, the key could be the text of the query. MillWheel uses the concept of *low-water mark* to bound the time stamp of future records. Waiting for the low-water mark allows computations to reflect a more complete picture of the data. *Timers* are programmatic hooks that trigger at a specific wall time or at a low-water mark value for a particular key.

A user-defined computation subscribes to input streams and publishes output streams, and the system guarantees delivery of both streams. Computation is run in the context of a specific key. The processing steps for an incoming record are:

1. Check for duplication and discard a previously seen record.
2. Run user code and identify pending changes to the timers, state, and production.
3. Commit pending changes to a backing store.

⁶ An idempotent operation can be carried out repeatedly without affecting the results.

4. Acknowledge senders.
5. Send pending downstream productions.

There are two types of productions, strong and weak. Strong productions support handling of inputs that are not necessarily ordered or deterministic. Checkpointing of *strong productions* is done before delivery as a single atomic write; the state modification is done at the same time as the atomic write. *Weak productions* are generated when the application semantics allows the user to disable the strong production option.

The computations are pipelined and may run on different hosts of the MillWheel cluster. To record the persistent state, the system uses either BigTable or other databases supporting single-row updates. A replicated master balances the load and distributes it based on lexicographic analysis of the record keys.

Measurements conducted on the system report low latency and scalability, well within the targets set at Google for stream processing. A median record delay of 3.6 milliseconds and a 95th-percentile latency of 30 milliseconds are reported for an experiment on 200 CPUs.

MillWheel is not suitable for monolithic computations whose checkpointing would interfere with the dynamic load balancing necessary to ensure low latency.

Caching strategies for data streaming. Routers of a Content-Centric Network (CCN)—see Section 6.12—use several replacement policies. LRU (Least Recently Used) keeps in cache *the most recently used content*, while LFU (Least Frequently Used) uses *cache request history to keep in cache highly used content*.

There are obvious limitations of both LRU and LFU. The former ignores the history, e.g., items highly used in the past are kept in cache in spite of a new patterns of access, while the latter ignores the popularity of an item. The combinations of the two, the Least Recently/Frequently Used replacement strategies, support tradeoffs by specifying a weight for each request that decays exponentially over time. The sLRFU (streaming Least Recently/Frequently Used) introduced in [419] aims to maximize cache. The novelties introduced by sLRFU are:

1. Partitions a cache of size C into an LFU-managed area using a sliding window of size k and an LRU-managed area of size $C - k$.
2. Estimates the top- k most popular data item in a sliding window of requests, keeps the frequently used ones in cache, and discards the old ones. k is set dynamically based on bounds given by the streaming algorithm.
3. For every request:
 - Increases the reference counter for the data item referenced by the request and decreases the counter for the request falling out of the window.
 - If the data requested is not in cache, it recovers the data and delivers it.
 - Using the last N requests, possibly rearranges the top- k most popular elements in cache (adding and/or removing content from the list) and complements the available spaces in cache with the $C - k$ most recently requested elements that are not already in cache.

Simulation results reported in [419] show that sLRFU has a hit rate of 70%, compared with a baseline LRU of 65%.

12.5 A dataflow model for data streaming

A more sophisticated programming model for data streaming was developed at Google. Alphabet, the holding company consisting of Google's core businesses and future-looking and finance activities, earns billions of dollars from advertising embedded in data streaming. That's why Google is interested to support an effective computational model for event-time correlations. The Dataflow model described in [16] is based on MillWheel [15] and FlumeJava [90].

Instead of the widely used terminology distinguishing batch processing from streaming, the two types of data processing are identified by a characterization of the input data: *bounded* for batch processing and *unbounded* for streaming. The implications of this dichotomy reflect what the execution engine is expected to do. In the bounded case, the engine processes a data set of known contents and size. In the unbounded case the engine processes a dynamic data set where one never knows if the set is complete because new records are continually added and old ones are retracted.

Dataflow SDK motivation for developing a new model is that grooming unbounded data sets to look like bounded ones does not allow an optimal balance among correctness, latency, and cost. It is thus necessary to refresh, simplify, and have flexible programming models for unbounded datasets.

The shortcomings of previous models are perfectly well illustrated by Google applications. For example, assume that a streaming video provider wants to bill advertisers placing their adds along with the video stream for the amount of advertising watched. This requires the ability to identify who watched a video stream and each add and for how long.

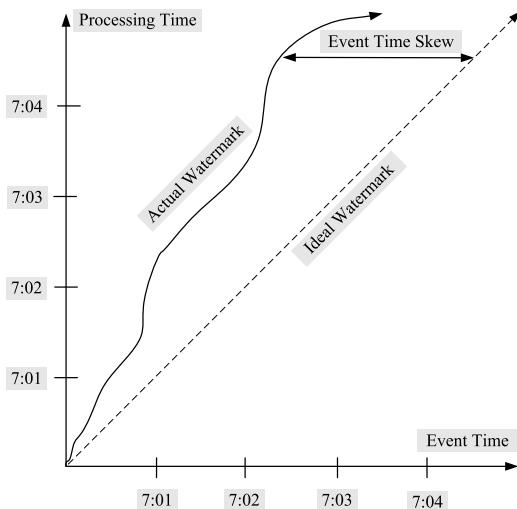
Existing systems did not provide exactly once semantics, do not scale well, are not fault tolerant, and have high latency, while others such as MillWheel or Spark Streaming do not have high-level programming support for event-time correlations. The model introduced in [16] "allows for calculation of event-time ordered results, windowed by features of the data themselves, over an unbounded, unordered data source, with correctness, latency, and cost tunable across a broad spectrum of combinations."

It also "Decomposes pipeline implementation across four related dimensions, providing clarity, composable, and flexibility: *What* results are being computed. *Where* in event time they are being computed. *When* in processing time they are materialized. *How* earlier results relate to later refinements." The design of the system was guided by several principles:

- Never rely on the notion of completeness.
- Encourage clarity of implementation.
- Support data analysis in the context it was collected.

The new model uses *windowing* to split a dataset into groups of events to be processed together, an essential concept for unbounded data processing. Windowing is almost always time-based and the windows can be static/fixed or sliding. A *sliding* window has a size and a sliding period, e.g., a two-hour window starting every 10 minutes. *Sessions* are sets of windows related by a data attribute, e.g., key for key-value datasets.

An *event* causes the change of the system state, for example, the arrival of a new record in case of data streaming. Typically, the *time skew*, the difference between the event time and the event processing time, varies during processing. The distinction between the event time and the processing time of the event is illustrated in Fig. 12.5.

**FIGURE 12.5**

Event time skew. *Event time* is the wall clock time when the event occurred; this time never changes. *Event processing time* is the time when the event was observed during the processing pipeline; this time changes as the event flows through the processing pipeline. For example, an event occurring at 7:01 is processed at 7:02.

Dataflow SDK uses the *ParDo* and *GroupByKey* operations of FlumeJava, discussed in Section 10.15, and defines a new operation *GroupByKeyAndWindow*. The entities flowing through the pipeline are four-tuples (*key, value, eventtime, window*).

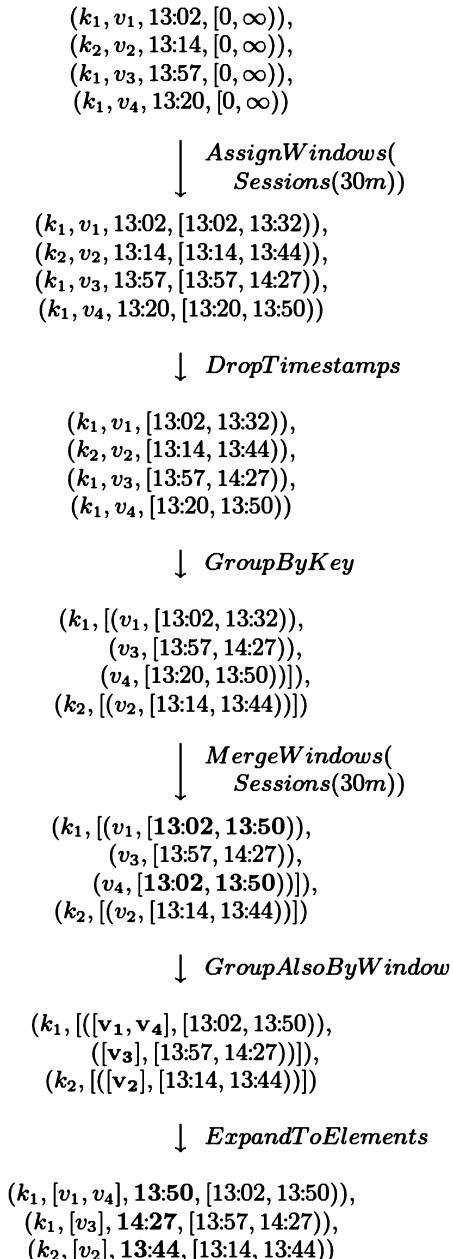
Several operations involving windows are defined. *Window assignment* replicates the object in every window; the (*key, value*) pairs are duplicated in windows overlapping the time stamp of an event. *Window merging* is a more complex operation because it involves the six steps in Fig. 12.6 [16]. In this example, the window size is 30 minutes, and there are four events, three with key k_1 and one with key k_2 .

In Step 2, four-tuples (*key; value; eventtime; window*) are transformed as three-tuples (*key; value; window*). The *GroupByKey* combines the three events with k_1 in Step 3. Overlapping windows $[13 : 02, 13 : 32]$ and $[13 : 20, 13 : 50]$ for k_1 are merged as $[13 : 02, 13 : 50]$ in Step 4. Then, in Step 5, the two k_1 events with values v_1 and v_4 in the same window $[13 : 02, 13 : 50]$ are grouped together. Finally, in Step 6, time stamps are added.

Watermarks used in MillWheel to trigger processing of the events in a window cannot by themselves ensure correctness; sometimes late data is missed. At the same time, a watermark may delay pipeline processing. The Dataflow system uses *triggers* to determine the time when the results of groupings are emitted as panes.⁷

The system has several refinement mechanisms to control how several panes are related to one another. Window contents are discarded once triggered, provided that later pipelines stages expect

⁷ A pane is a well-defined area within a window for the display of, or interaction with, a part of that window's application or output.

**FIGURE 12.6**

Six window-merging steps: 1. AssignWindow; 2. Drop time stamps; 3. GroupByKey; 4. MergeWindows—based on windowing strategy; 5. GroupAlsoByWindow—for each key group values by window; 6. ExpandToElements—expand per-key, per-window groups into $(key; value; eventtime; window)$ tuples, with new per-window time stamps.

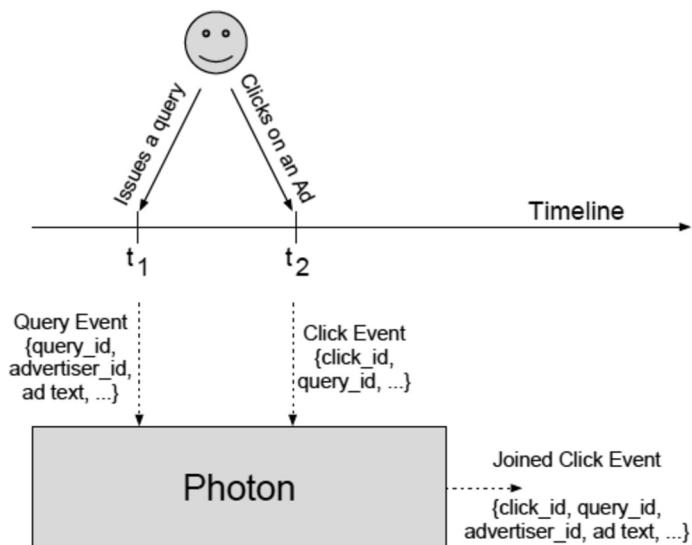


FIGURE 12.7

The primary and the secondary stream events in Photon are the query stream and ad clickstream events. At time t_1 , a web server responding to a query also serves an ad. The user clicks on the ad at time t_2 , and then the query event and the click event are combined into a joined click event [24].

the values of the triggers to be independent. The window contents are also discarded when the user expresses no interest in later events. The contents of a window can be saved in persistent storage when needed for refining the system state.

12.6 Joining multiple data streams

Applications such as advertising, IP network management, and telephone fraud detection require ability to correlate in real-time events occurring in separate high-speed data streams. For example, Google search engine delivers ads along with answers to queries. The Photon system [24] developed for Google Advertising System produces joint logs for the query data stream and the ad click data stream. The joint logs are used for billing the advertisers. A typical Photon scenario is illustrated in Fig. 12.7.

Photon processes millions of events per minute with an average end-to-end latency of less than 10 seconds. The system joins 99.9999% events within a few seconds and 100% events within a few hours. Photon designers had to address a set of challenges, including:

- Each click on an ad should be processed once and only once. This means: at-most-once semantics at any point of time, near-exact semantics in real-time, and exactly once semantics, eventually. If a click is missed, Google loses money; if the click is processed multiple times, the advertisers are overcharged.

- Automated data center-level fault tolerance. Manual recovery takes too long. Photon instances in multiple data centers will attempt to join the same input event, but must coordinate their output to guarantee that each input event is joined at-most-once.
- Latency constraints. Low latency is required by advertisers because it helps optimize their marketing campaigns. Load balancers redirect a user request to the closest running server, where it is processed without interacting with any other server.
- Scalability. The event rates are very high now and are expected to increase in the future; therefore, the system has to scale up to satisfy latency constraints.
- Combine ordered and unordered streams. While the query stream events are ordered by time stamps, the events in the clickstream may occur at any time and are not sorted. Indeed, the user may click on the ad long after the results of the query are displayed.
- Delays due to the system scale. The servers generating query and click events are distributed over the entire world. The volume of query logs is much greater than that of the click log; thus the query logs can be delayed relative to the click logs.

Photon organization and operation. *IdRegistry* stores the critical state shared among running Photon instances all over the world. This critical state includes the *eventId* of all events joined in the last N days. Before writing a joined event, each instance checks whether the *eventId* already exists in the *IdRegistry*; if so, it skips processing the event, otherwise it adds the event to *IdRegistry*.

IdRegistry is implemented using the Paxos protocol discussed in Section 10.13. An in-memory key-value store based on PaxosDB is consistently replicated across multiple data centers to ensure availability in case of the failure of one or more data centers. All operations are carried out as read-modify-write transactions to ensure that writing to the *IdRegistry* is carried out if and only if the event is not already in.

An *eventId* uniquely identifies the server, the process on that server that generated the event, and the time the event was generated

$$\text{EventId} = (\text{ServerIP}, \text{ProcessId}, \text{Timestamp}). \quad (12.4)$$

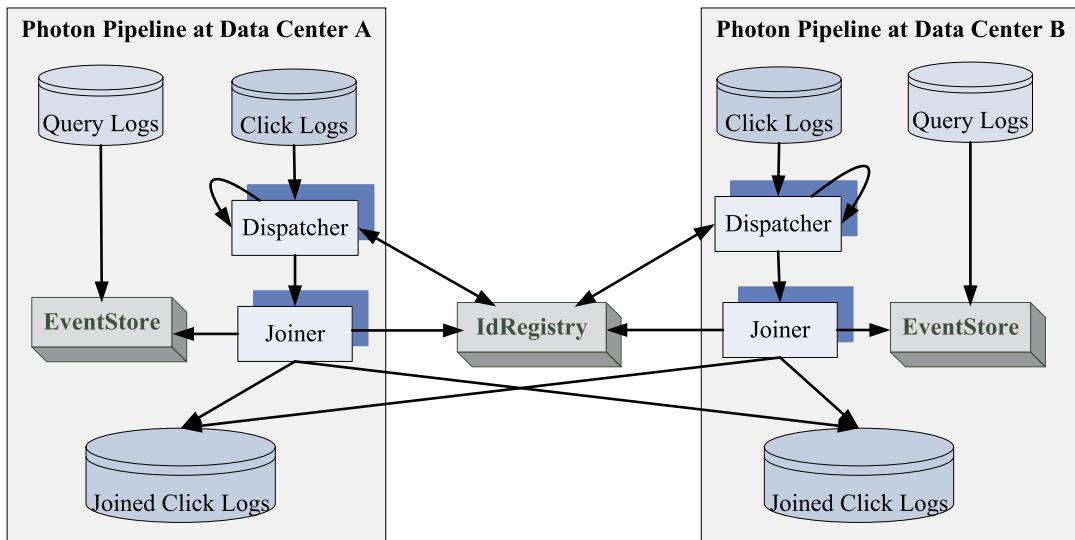
Events with different Ids can be processed independently of each other. This allows the *EventId space* of the *IdRegistry* to be partitioned into disjoint shards. *EventId*'s from separate shards are managed by separate *IdRegistry* servers.

Identical Photon pipelines run at multiple data centers around the world. A Photon pipeline has three major components, shown in Fig. 12.8:

1. *Dispatcher* - reads the stream of clicks and feeds them to the joiner.
2. *EventStore* - supports efficient query lookup.
3. *Joiner* - generates joined output logs.

The joining process involves several steps:

1. The *Dispatcher* monitors the logs, and, when new events are detected, it looks-up the *IdRegistry* to determine if the *clickId* is already recorded.
 - If already recorded, skip processing the click.
 - Else, the dispatcher sends the event to the joiner and waits for a reply. To guarantee at-least-once semantics, the dispatcher resends it until it gets a positive acknowledgment from the joiner.

**FIGURE 12.8**

Photon pipeline organization. The pipeline at each site includes a *Dispatcher*, an *EventStore*, and a *Joiner* operating on query logs, click logs, and joined click logs. *IdRegistry* stores the critical state shared among running Photon instances all over the world.

2. The *Joiner* extracts *queryId* and carries out an *EventStore* lookup to locate the corresponding query. If the query is

- Found: then the joiner attempts to register the *clickId* in the *IdRegistry*.
 - If the *clickId* is in the *IdRegistry*, the *Joiner* assumes that the join has already been done.
 - If not, the *clickId* is recorded in the *IdRegistry*, and the event is recorded in the joint event log.
- Not found: then the *Joiner* sends a failure response top the dispatcher; this will cause a retry.

In the US, there are five replicas of the *IdRegistry* in data centers in three geographical regions up to 100 ms apart in round-trip latency. The other components in the pipelines are deployed in two geographically distant regions on the east and west coasts. According to [24]: “during the peak periods, Photon can scale up to millions of events per minute,... each day, Photon consumes terabytes of foreign logs (e.g. clicks), and tens of terabytes of primary logs (e.g. queries), ...more than a thousand *IdRegistry* shards run in each data center,... each data center employs thousands of *Dispatchers* and *Joiners*, and hundreds of *CacheEventStore* and *LogsEventStore* workers.”

12.7 Mobile computing and applications

Mobile devices are in a symbiotic relationship with computer clouds. “Mobile cloud computing at its simplest, refers to an infrastructure where both the data storage and data processing happen

outside of the mobile device. Mobile cloud applications move the computing power and data storage away from mobile phones and into the cloud, bringing applications and mobile computing to not just smartphone users but a much broader range of mobile subscribers,” according to <http://www.mobilecloudcomputingforum.com/>.

Mobile devices are producers, as well as consumers, of data stored on the cloud. They also take advantage of cloud resources for computationally intensive tasks. Mobile devices benefit from this symbiotic relationship; their reliability is improved as data and the applications are backed up on the cloud. Computer clouds extend the limited physical resources of smartphones, tablets, and laptops: (i) The processing power and the storage capacity—an ubiquitous applications of mobile devices is to capture still images or videos, upload to a computer cloud, and make them available via cloud services such as Flickr, Youtube, or Facebook; and (ii) Battery life—migrating mobile game components to a cloud can save 27% of the energy consumption for computer games and 45% for a chess game [124]. Clouds also extend the utility of mobile devices providing access to vast amounts of information stored on their servers. There are also major benefits due to the integration of mobile devices in the cloud ecosystem [438]:

- Mobile devices capture images and videos rich in content, rather than simple scalar data such as temperatures. Such information can be used to understand the extent of catastrophic events, such as earthquakes or forest fires.
- There is power in numbers—cloud sourcing amplifies the power of sensing and can be used for applications related to security, rapid service discovery, etc.
- Near-real-time data consistency is important in disaster relief scenarios. For example, the aftershocks of an earthquake often trigger major structural damage to buildings. Images and/or videos taken before and after the event help assess the extendit of the damage and identify buildings requiring immediate evacuation.
- Opportunistic information gathering. For example, anti-lock braking devices on cars transmit their GPS coordinates on each activation, enabling maintenance crews to identify the slick spots on roads.
- Computer vision algorithms running on clouds can use images captured by mobile devices to locate lost children in crowds, estimate crowd sizes, etc.

Applications of mobile cloud computing. Many mobile applications consist of a lightweight front-end component running on the mobile device and a back-end data-intensive and computationally intensive component running on a cloud. There are countless examples of such ubiquitous applications of mobile cloud computing in areas as diverse as health care, learning, electronic commerce, mobile gaming, mobile location services, and search.

An important application of mobile healthcare is patient record storage and retrieval and medical image sharing using computer clouds and mobile devices. Other applications in this area are [148]: health monitoring services where patients are monitored using wireless services; emergency management involving the coordination of emergency vehicles; and wearable systems for monitoring the vital signs of outpatients after surgery.

There are many education tools helping students learn and understand subjects as diverse as anatomy, computer science, or engineering, or the arts. Though some progress has been made, the potential of mobile learning is not fully exploited at this time. Classroom tools using AI algorithms to identify relevant questions posed by students during lectures and enabling instructors to interact with the students in classes with large enrollments are yet to be developed.

Mobile gaming has the potential of generating large earnings for cloud service providers. The engine requiring large computing resources can be offloaded to a cloud, while gamers only interact with the screen interface on their mobile devices. For example, to maximize energy savings, the Maui system [124] partitions the application codes at runtime based on the costs of network communication and the CPU power and the energy consumption of the mobile device. Browsers running on mobile devices are often used for keyword-, voice-, or tag-based searching of large databases available on computer clouds.

All applications of mobile cloud computing face challenges related to communication and computing. The low bandwidth, the service availability, and the network heterogeneity pose serious problems on the communication side. The security, the efficiency of data access, and the offloading overhead are top concerns on the computing side.

Mobile devices are exposed to a range of threats, including viruses, worms, Trojan horses, and ransomware. Such threats are amplified because mobile devices have limited power resources, and it is impractical to continually run virus-detection software. Moreover, once files located on the mobile device are infected, the infection propagates to the cloud when the files are automatically uploaded. At the same time the integration of GPS is a source of concern for privacy as the number of location-based services (LBS) increases. Cloud data security, integrity, and authentication, along with digital rights management are also sources of concern.

Enhancing the efficiency of data access requires a delicate balance between local and remote operations and the amount of data exchanged between the mobile device and the cloud. The number of I/O operations executed on the cloud in response to a mobile device request should be minimized to reduce the access time and the cost. The memory and storage capacity of the mobile device should be used to increase the speed of data access, reduce latency, and improve energy efficiency of mobile devices.

The future of mobile cloud computing. As the technology advances, mobile devices will be equipped with more advanced functional units, including high-resolution cameras, barometers, light sensors, etc. Augmented reality and mobile gaming are emerging as important mobile cloud computing applications. Augmented reality could become a reality with new mobile devices and fast access to computer clouds. Composition of real-time traffic maps from collective traffic data sensing, monitoring environmental pollution, and traffic and pollution management in smart cities are only a few of the potential future applications of mobile cloud computing. A recent survey [505] analyzes applications, the solutions to the challenges they pose, and the future solutions:

- Code and computation offloading—currently based on static partitioning and dynamic profiling is expected to be automated in the future.
- Task-oriented mobile services—currently provided by Mobile-Data-as-a-Service, Mobile-Computing-as-a-Service, Mobile-Multimedia-as-a-Service, and Location-based Services—are expected to be replaced by human-centric mobile services.
- Elasticity and scalability—components of the resource allocation and scheduling are expected to be validated by algorithms using valid traffic VM migration models.
- Cloud service pricing—currently based on auctions, and bidding is expected to be replaced by empirical validation and optimization algorithms.

Unquestionably, inclusion of mobile devices into the cloud ecosystem has countless benefits. It opens new avenues for learning, increased productivity, and entertainment, but can also have less desirable effects since it can overwhelm us with information. Herb Simon reflected on the effects of

information overload [448] “What information consumes is rather obvious: it consumes the attention of its recipients. Hence a wealth of information creates a poverty of attention and a need to allocate that attention efficiently among the overabundance of information sources that might consume it.”

12.8 Energy efficiency of mobile computing

Mobile computing benefits from virtually infinite resources available on computer clouds provided that mobile devices have network access to these islands of plentiful computing and storage resources. All cloud users have to transfer data to/from the cloud, but there is a significant dissimilarity between the cloud users connected to computer clouds via landlines and using stationary devices and users of mobile devices connected via cellular or wireless local area networks (WLANs). The first are concerned with the time and the cost of transferring massive amounts of data, while for mobile cloud users the key issue is the energy consumption for communication.

Battery technology is lagging behind the needs of modern mobile devices. The amount of energy stored in a battery is growing by only about 5% annually. Increasing the battery size is not an option because devices tend to be increasingly lighter. Moreover, the power of small mobile devices without active cooling is limited to about three watts [367].

An analysis of the critical factors regarding mobile devices energy consumption reported in [347] is discussed in this section. As expected, the computing-to-communication ratio is the critical factor for balancing local processing and computation offloading. The analysis shows that not only the amount of transferred data, but also the traffic pattern is important. Indeed, sending a larger amount of data in a single burst consumes less energy than sending a sequence of small packets.

The analysis of energy used by mobile computing uses several parameters: E_{cloud} and E_{local} , the quantitative characterization of the energy for a computation on the cloud and locally, respectively; D - the amount of data to be transferred; C - the amount of computation expressed as CPU cycles. C_{eff} and D_{eff} are the efficiencies of computing and device specific data transfer, respectively. C_{eff} is measured in CPU cycles per Joule and represents the computations carried out with a given energy.

Dynamic voltage and frequency scaling affects only slightly the power and performance of the CPU during execution, thus C_{eff} . For example, the N810 Nokia processor requires 0.8 W and has $C_{eff} = 480$ when running at 400 MHz and needs only 0.3 W and has a $C_{eff} = 510$ at 165 MHz. This is also true for the more performant N900 Nokia processor; when running at 600 MHz, the power required and the C_{eff} are 0.9 W and 650 cycles/J, respectively, while at 250 MHz the same parameters are 0.4 W and 700 cycles/J, respectively.

D_{eff} is measured in bytes per Joule and represents the amount of data that can be transferred with a given energy. D_{eff} is affected by the traffic pattern. The time the network interface is active affects the energy consumption for communication. Typically, the power consumption for activation and deactivation of a cellular network interface is larger than that of a wireless interface. The larger the clock rate, the larger the power consumption and D_{eff} . For example, for N810, the power consumption and D_{eff} at 400 Mhz are 1.5 W and 390 KB/J, respectively, and decrease to 1.1 W and 310 KB/J, respectively, at 165 MHz.

It makes sense to offload a computation to a cloud if

$$E_{cloud} < E_{local}, \quad (12.5)$$

with

$$E_{cloud} = \frac{D}{D_{eff}} \quad \text{and} \quad E_{local} = \frac{C}{C_{eff}}. \quad (12.6)$$

The condition expressed by Eq. (12.5) becomes:

$$\frac{C}{D} > \frac{C_{eff}}{D_{eff}}. \quad (12.7)$$

According to [347], a rule of thumb is that “offloading computation is beneficial when the workload needs to perform more than 1 000 cycles of computation for each byte of data.”

The CPU cycles per byte ratios measured on a system with an ARM Cortex-A8 core running at 720 Mz for several applications are: 330 for gzip ACII compression; 1 300 for x264 VBR encoding; 1 900 for CBR encoding; 2 100 for htm12text on wikipedia.org; and 5 900 for htm12text on en.wikipedia.org.

Several conclusions were drawn from the experiments reported in [347]:

- Energy consumption of a mobile device is affected by the end-to-end chain involved in each transaction thus, the server-side resource management is important.
- Higher performance typically contributes to better energy efficiency.
- Simple models able to guide design decisions for increasing energy efficiency should be developed.
- Automated decisions on whether a computation should be uploaded to a cloud to maximize energy efficiency of mobile devices should be built into mobile cloud computing middleware.
- Latencies associated with wireless communication are critical for interactive workloads.

It is unlikely we will see dramatic improvements of the energy storage capacity for mobile devices batteries in the future. The need for increasingly more sophisticated software for energy optimization should motivate the research in this challenging area because new mobile cloud computing data and computation-intensive applications are developed every day. At the same time, the other critical limitations of mobile cloud computing applications, communication speed, and effectiveness are improving because WLAN speed is steadily increasing, along with the efficiency of mobile device antenna.

12.9 Alternative mobile cloud computing models

In Section 12.7, we have seen that mobile device access to the large concentration of resources in cloud computing data centers is supported by wide-area networks (WANs), cellular networks, and wireless networks. In this section, we examine the limitations due to the communication latency and, to some extent, of the bandwidth and discuss alternative mobile cloud computing architectures.

Latency effects. Network security, energy efficiency, and network manageability, along with bandwidth, are the main sources of concern for networking companies and for networking research. Unfor-

tunately, virtually all methods to address these concerns have a negative side effect, *they increase the end-to-end communication latency*. Indeed, the techniques to increase energy efficiency discussed in Section 12.8 include reduction of the time the network interface is active, delaying data transmission until large data blocks can be sent, and turning on the transceivers of mobile devices for short periods of time to receive and to acknowledge data packets buffered at a base station.

The speed of the fastest wireless LAN (802.11n) and of the wireless Internet HSPA (High-Speed Downlink Packet Access) technologies are 400 Mbps and 2 Mbps, respectively, and the corresponding transmission delays for a 4 Mbyte Jpeg image are 80 versus 16 msec, respectively. The range of Internet latencies is 30–300 msec. For example, the mean Berkeley to Trondheim-Norway and Berkeley to Canberra-Australia latencies are 197 and 174 msec, respectively, while Pittsburgh to Hong Kong and Pittsburgh to Seattle are 223 and 83.9 msec, respectively, as measured in 2008–2009 [437]. The current Internet latencies between selected pairs of points on the globe can be found at <https://www.internetweathermap.com/map>.

The latency effect differs from application to application. For example, the subjective effects of the latency, L , for a particular application, GNU Manipulation Program (GIMP) for Virtual Network Computing communication software, the graphical desktop sharing system that uses the Remote Frame Buffer protocol are: crisp when $L < 150$ msec; noticeable- to-annoying when $150 < L < 1\,000$ msec; annoying when $1 < L < 2$ sec; unacceptable when $2 < L < 5$ sec; and unusable when $L > 5$ sec [437].

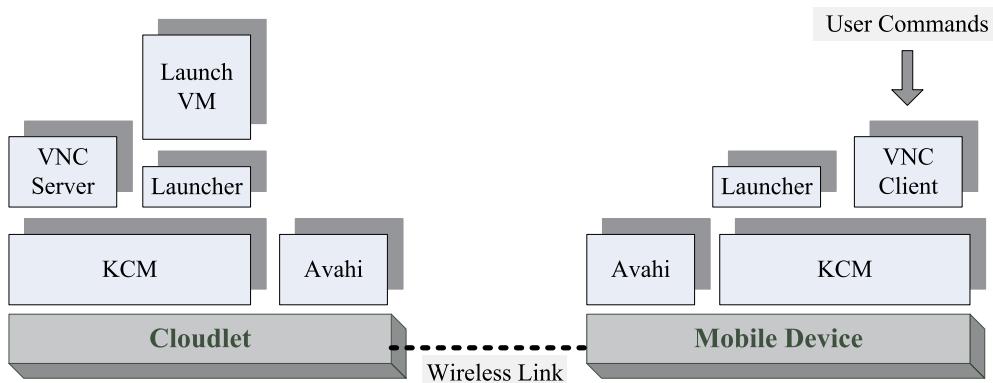
Cloudlets. A possible solution to reduce end-to-end latency mimics the model supporting wireless communication where access points are scattered throughout campuses of large organizations and in cities. Micro data centers or “clouds in a box,” called *cloudlets*, are placed in the proximity of regions with a high concentrations of mobile systems. A cloudlet could be a cluster of multicore processors with a fast interconnect and a high-bandwidth wireless LAN. The resources available on such cloudlets pale in comparison with the ones on a cloud; the cloudlets only assist the mobile devices for data and computationally intensive tasks since they are not permanent repositories of data. A mobile device could connect to a cloudlet and upload code to the cloudlet.

There are obvious benefits, as well as problems, with this solution proposed in [437]. The main benefit is the short end-to-end communication delay for mobile devices in the proximity of the cloudlet. On the other hand, the cost and the maintenance of cloudlets able to support a wide range of mobile users are of concern. It is safe to assume that hardware costs will decline as the processor technology evolves and that the processor bandwidth and reliability will increase. The software maintenance effort can be reduced if self-management techniques are developed. A range of business and technical challenges including cloudlet sizing must be addressed before this solution gets sufficient traction.

The solution for software organization proposed in [437] is to have a permanent cloudlet host software environment and a transient guest software environment. A transient customization configures the system for running an application, and, once the run is complete, a cleanup phase clears the way for running the next one. A VM encapsulates the transient guest environment.

The VM can be created on the mobile device, run locally until the VM needs additional resources, then stopped, its state saved in a file, and the file sent to a cloudlet in the vicinity of the mobile device where the VM is restarted. The problem with this straightforward solution is that the latency can be increased when the footprint of the VM migrated to the cloudlet is substantial.

Another solution is the *dynamic VM synthesis* in which the mobile device delivers to the cloudlet only a small overlay. This solution assumes that cloudlet already has the base VM used to derive the

**FIGURE 12.9**

Kimberley organization. Kimberley Control Manager (KCM) runs on both the cloudlet and the mobile device. The two communicate over a wireless link using the Virtual Network Computing (VNC) communication software and Avahi [437].

small overlay; therefore the cloudlet environment can start execution immediately without the need to contact a cloud or other cloudlets. The assumption that a relatively small set of base VMs will suffice for a large range of applications may prove to be overly optimistic.

Kimberley—a proof of concept cloudlet system [437]. The cloudlet infrastructure consists of a desktop running Ubuntu Linux and the mobile device is a Nokia N810 tablet running Maemo 4.0 Linux. The cloudlet uses a hosted hypervisor for Linux called VirtualBox and a tool with three components, *baseVM*, *instal-script*, and *resume-script* to create the VM overlays for any OS compatible with the components of the tool.

First, the *baseVM* is launched, then the *instal-script* in the guest OS is executed, and finally the *resume-script* in the guest OS launches the application. The VM encapsulates the application, the so-called *launchVM* can be activated without the need to reboot, and finally this VM is compressed and encrypted.

Kimberley organization, depicted in Fig. 12.9, shows the software components running on the cloudlet and the mobile device. Communication through the wireless link is supported by Virtual Network Computing (VNC) communication software and Avahi,⁸ a free zero-configuration⁹ for automatic networking implementation supporting multicast DNS/DNS-SD service discovery.

KCM abstracts the service discovery, including browsing and publishing in Linux using Avahi. A secure TCP tunnel using SSL is established between KCM instances. After the establishment of

⁸ Avahi is the Latin name of the woolly lemur primates indigenous to Madagascar.

⁹ Zero-configuration networking (zeroconf) creates a TCP/IP computer network based on automatic assignment of numeric network addresses for networked devices, automatic distribution and resolution of computer hostnames, and automatic location of network services. It is used to interconnect computers or network peripherals.

the secure TCP tunnel, the authentication using the Simple Authentication and Security Layer (SASL) framework is carried out. Then, the KCM on the cloudlet fetches the VM overlay from the mobile device KCM, decrypts and decompresses the VM overlay, and applies the overlay to the base VM.

12.10 System availability at scale (R)

As the cloud user population grows and the scale of the data center infrastructure expands, the concerns regarding system availability are amplified. Very high system availability and correctness are critical for data streaming processing systems. Such systems are particularly vulnerable because the likelihood of missed events is very high when the resources required for event processing suddenly become unavailable.

When we think about availability, a 99% figure of merit seems quite high. There is another way to look at this figure: *A 99% availability translates to 22 hours of downtime per quarter, and this is not acceptable for mission-critical systems.* Google aims to achieve 99.9999–99.99999% availability for their data streaming systems [218]. How can this be achieved? An answer to this question is the topic of this section.

Multiple failure scenarios are possible. Individual servers may fail, all servers in one rack may be affected by a power failure, or a power failure may force the entire data center infrastructure to shut down, though such events are rare. The interconnection network may fail and partition the network, or public networks connecting a data center to the Internet or private networks connecting the data center with other data centers of the same CSP may fail.

Workloads can be migrated to functioning systems in case of partial failures affecting a subset of resources, and then users may experience slower communication and extended execution times in such cases. Often, it takes some time before the cause of a partial failure is found and corrective measures are taken. Also, events are missed when partial failures affect data streaming.

Shutdown of servers, networks, or the entire data center may be planned for hardware or software updates or maintenance. Such events are announced in advance, and in such cases, the shutdown is graceful and without data loss or other unpleasant consequences for the cloud users. The unplanned shutdown are the ones the CSPs are concerned about. The most consequential are the data center-level failures.

While current systems manage quite well the failure of individual servers, data center-level failures, though seldom occurring, have catastrophic consequences especially for *single-homed* software systems, the ones running *only* at the site affected by the failure. Less affected are the *failover-based* software systems. Such systems only run at one site, but checkpoints are created periodically and sent to backup data centers. When the primary data center fails, the last checkpoints of the critical systems are used to restart processing at the backup site. Data streaming is affected in this case as events are likely to be missed.

Checkpointing can be done asynchronously or synchronously. In the first case, all events are processed exactly once during the restart if the events in progress are drained before checkpointing, and only when the shutdown was planned. Synchronous checkpointing can, in principle, capture all events even in case of unplanned shutdowns, but their processing is considerably more complex. For planned

shutdowns, they require each pipeline to generate a checkpoint and block until the checkpoint is replicated.

Virtually all CSPs have multiple data centers distributed across several regions, and a basic assumption for handling data center-level failures is that the probability of a simultaneous failure of data centers in multiple regions is extremely low. Critical software systems are *multi-homed*, the workload is dynamically distributed among these centers, and, when one of them fails, its share of the workload is reassigned to the ones still running.

The design of multi-homed systems is not without its own challenges. First, the global state of the system must be available to all sites. The size of the metadata reflecting the global state should be minimized. Communication delays of tens of milliseconds and limited global communication bandwidth do not make determination of the global state an easy task, even when a Paxos-based commit is used to update the global state.

Second, a data streaming pipeline is implemented as a network with multiple processing nodes. Checkpointing has to capture the state of all nodes, as well as metadata such as the state of the input and output queues of pending and completed node work. Clustering groups nodes together and recording only the global state of the cluster helps. Clustering increases the checkpointing efficiency and shrinks the checkpoint size.

Handling the same input events at several sites can be optimized so that the checkpoint includes only one copy of the event. This is actually done for the events recorded by logs in the system discussed next. Such primary events may require database lookup, and if the database data does not change, then the state of the primary event should not record the information from the database.

Lastly, guaranteeing exactly once semantics for multihomed systems is nontrivial because pipelines may fail while producing output; multiple sites may process the same events concurrently. When the global state is stored in a shared memory, updating global state can be atomic. Idempotence can help achieving the desired semantics. If multiple sites update the same record, this guarantees the desired semantics. If idempotence is not achievable, a two-phase commit can be used to write the results produced by the streaming system to the output.

Data streaming services collect data generated by user interactions and expect consistency. Data streaming consistency means that, if a state at time t includes an event e , then any future state at time $t + \tau$ must also include event e and an observer should see consistent outcomes if more than one are generated.

The concern for data center-level failures forced Google to run multiple copies of critical systems at multiple data centers. Several large-scale Google systems run hot in multiple data centers, including the F1 database [447], discussed in Section 12.2, and the Photon system [24], discussed in Section 12.6.

An infrastructure supporting ads management at Google is discussed in [218]. The strategy for supporting availability and consistency for the streaming system is to log the events caused by user interactions in multiple data centers. Local logs are then collected by a *log collection service* and are sent to several locations designated as *log data centers*. Consistency of these logs is critical and requires *exactly once* event processing semantics.

Once a system is developed for a single site, it is very hard to adapt to multisite operations, especially when consistency among sites is a strong requirement. An *ab initio* design as a multihome is the optimal solution for mission critical systems. Running a system at multiple sites transfers the burden to solve the very hard problem caused by data center-level failures to the infrastructure. This burden falls to

the users in traditional systems designed to run at one site only. Indeed, the users have to take periodic checkpoints and transfer them to backup sites.

Multihoming is challenging for system designers and adds to the resource costs. Additional resources are needed to process the workload normally processed by the failing data center and to catch up after delays. The spare capacity has to be reserved ahead of time and should be ready to accept the additional workload.

12.11 Scale and latency (R)

The “scale” of the cloud has altered our ideas on how to compute efficiently and has generated new challenges in virtually all areas of computer science and computer engineering, from computer architecture to security, via software engineering and machine learning. A fair number of models and algorithms for parallel processing have been adapted to the cloud environment. Frameworks for Big Data processing, resource management for large-scale systems, scheduling algorithms for latency-critical workloads discussed in Chapters 11, 4, 9, 5, and 12 offer only a limited window into the new computing landscape. Several critically important ideas for the future of cloud computing are discussed in this section.

Latency, the time elapsed from the instant when an action is initiated to its completion, has always been an important measure of quality of service for interactive applications, but its relevance and impact has been substantially amplified in the age of online searches and web-based electronic commerce. There are quite a few latency-critical applications of cloud computing now, and their number is likely to explode when a large number of cyber–physical systems of the IoT will invade our working and living space.

The latency has three major components: communication time, waiting time, and service time. Only the latter two are discussed in this section because the communication speed is not controlled by the cloud service provider, and in practice it has little impact on latency. A heavy-tail distribution of latency due to the latter two components is undesirable; it affects the user’s experience and, ultimately, the ability of the service provider to compete.

Heavy-tail distributions. A random variable X with the cumulative distribution function $F_X(x)$ is said to have a heavy right tail if

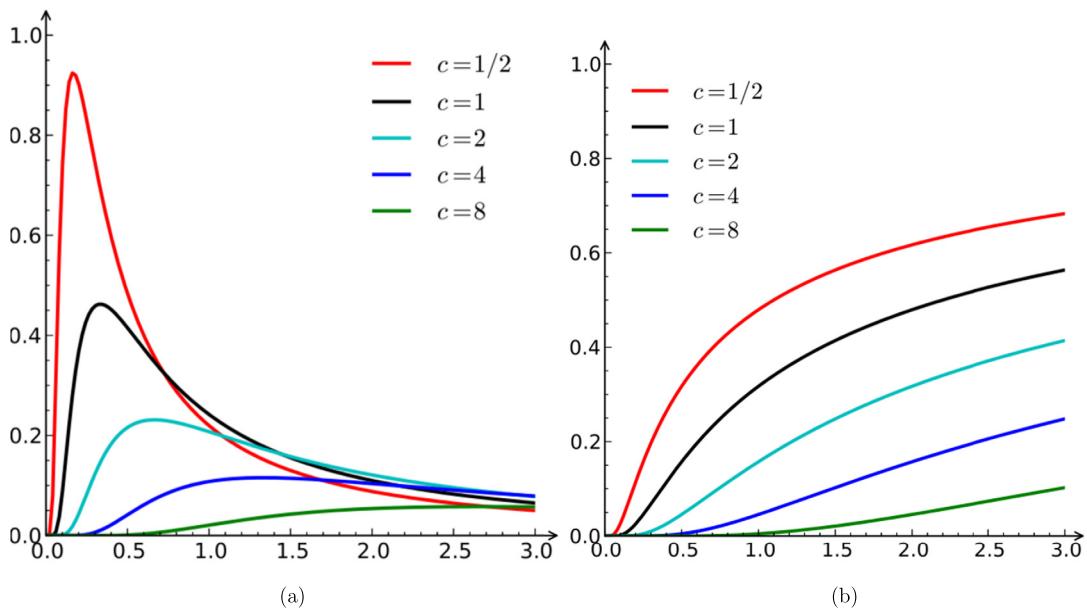
$$\lim_{x \rightarrow \infty} e^{\lambda x} Pr[X > x] = \infty, \quad \forall \lambda > 0. \quad (12.8)$$

This definition can also be expressed in terms of the tail-distribution function

$$\bar{F}_X(x) \equiv Pr[X > x] \quad (12.9)$$

and implies that $MF(t)$, the moment generating function of $F_X(x)$, is infinite for $t > 0$. The Lévy distribution is an example of a heavy-tail distribution. Its probability density function (PDF) over the domain $x \geq \mu$ and cumulative distribution function (CDF) are, respectively,

$$f_X(x; \mu, c) = \sqrt{\frac{c}{2\pi}} \times \frac{e^{-\frac{c}{2(x-\mu)}}}{(x-\mu)^{3/2}} \quad \text{and} \quad F_X(x; \mu, c) = erfc \left(\sqrt{\frac{c}{2(x-\mu)}} \right), \quad (12.10)$$

**FIGURE 12.10**

(a) The probability density function of the unshifted Lévy distribution, $\mu = 0$; (b) The cumulative distribution function.

with μ being the location parameter, c being the scale parameter, and $\text{erfc}(y)$ being the complementary error function. Fig. 12.10 displays the PDF and the CDF of the Lévy distribution. In probability theory and statistics, a scale parameter is a special kind of numerical parameter of a parametric family of probability distributions. The larger the scale parameter, the more spread out is the distribution function.

Query processing at Google. The analysis of query processing at Google gives us some insight into the reasons why latency has a heavy-tail distribution and how this can be managed. A query fans out into thousands of requests to a large number of servers where cached data resides. Some of these servers contain images, others videos, web data, blogs, books, news, and many other bits of data that could possibly answer the query.

The individual-leaf-request finishing times measured from the root node to the leafs of the query fan-out tree show the effect of the heavy-tail distribution [130]. As the limit x in the latency cumulative distribution function increases from 50 to 95 and then to 99 percentile, the latency increases:

- When a single randomly selected leaf node of the tree is observed the latency increases from 1 to 5, and then to 10 milliseconds, respectively.
- When we request that 95% of all leafs finish execution the latency increases from 12 to 32, and then to 70 milliseconds, respectively.
- When we request that all nodes finish execution the latency increases from 40 to 87, and then to 140 milliseconds, respectively.

These measurements show that the latency increases with the number of leafs of the fan-out tree; it also increases as we wait for more leaves to finish. The extent of the latency tail is significant, the difference between 95 and 99 percentile is dramatic, it doubles or nearly doubles for the three cases examined, 10 versus 5, 70 versus 32, and 140 versus 87, mimicking the law of diminishing returns.

Google is able to address the problems posed by the heavy-tail latency distribution [130]: “The system updates the results of the query interactively as the user types predicting the most likely query based on the prefix typed so far, performing the search and showing the results within a few tens of milliseconds.” How this is done is discussed in this section.

A brief analysis of the factors affecting the variability of the query response time helps understanding how this variability can be limited. Resource sharing is unavoidable and contention for system resources translates into longer response time for the losers. There are actors working behind the scenes carrying out system management and maintenance functions. For example, daemons are themselves users of resources and competitors of production workloads; they can occasionally cause an increase of the response time. Data migration, data replication, load balancing, garbage collection, and log compaction, are all activities competing for system resources.

Multiple layers of queuing in the hierarchical infrastructure increase the waiting time and this can be addressed by priority queuing and techniques discussed next. Optimization of resource utilization and energy saving mechanisms contribute to latency. For example, servers are switched to a power-saving mode when the workload dips. The latency can increase in this case because it takes some time before servers wakeup and are able to absorb a spike in demand. Dynamic frequency and voltage scaling, a mechanism to adapt the power consumption to the workload intensity, could also contribute to latency variability.

Several techniques can reduce the effects of component-level variability on the heavy-tail query latency distribution. These techniques are:

1. Define different service classes and use priority queuing to minimize waiting for latency-critical workloads.
2. Keep the server queues short allowing requests from latency-critical workloads to benefit from their high priority. For example, Google storage servers keep few operations outstanding, instead of maintaining a queue. Thus, high-priority requests are not delayed, while earlier requests from low-priority tasks are served.
3. Reduce head-of-line blocking. When a long-running task cannot be preempted, other tasks waiting for the same resource are blocked, a phenomenon called head-of-line blocking.¹⁰ The solution is to split a long-running task into a number of shorter tasks and/or to use time-slicing, to allocate the resource to task for brief periods of time. For example, Google’s Web search system uses time-slicing to prevent computationally expensive queries to add to the latency.
4. Limit the disruption caused by background activities. Some of the behind-the-scene activities such as garbage collection or log compaction require multiple resources. Their continuous running in the background could lead to increased latency of many high-priority queries over extended periods of time. It is more beneficial to allow such background activities to proceed in synchronous bursts of

¹⁰ For an intuitive scenario of head-of-line blocking, imagine a slow-moving truck on a one-lane, winding mountain road. It is very likely that a large number of cars will be stuck behind the truck.

Table 12.3 Read latency in ms with requests tied 1 ms [130] for a latency-critical application. (Left) Mostly idle system. (Right) System running a background job in addition to the latency-critical application.

Limit (%)	No hedge	Hedge-tied	No hedge	Hedge-tied
50	19	16	24	19
90	38	29	54	38
99	67	42	108	67
99.9	98	61	159	108

concurrent utilization of several servers, thus affecting only the latency-critical tasks over the time of the burst.

A very important realization is that *the heavy-tail distribution of latency cannot be eliminated; the only alternative is to develop tail-tolerant techniques for masking long latencies*. Two types of tail-tolerant techniques are used at Google: (i) *within request, short-term*, acting in milliseconds; and (ii) *cross-request, long-term*, acting at a scale of seconds.

The former technique works well for replicated data, e.g., for distributed file systems with data striped and replicated across multiple storage servers and for read-only datasets. For example, the spelling-correction service benefits from this mechanism because the model is updated once a day and handles a very high rate of requests, in the hundreds of thousands per second. Cross-request techniques aim to increase parallelism and prevent imbalance either by creating micro-partitions, by replicating the items likely to cause imbalance, or by latency-induced probation.

Hedged and *tied* requests are short-term tail-tolerant techniques. In both cases, the client issues multiple replicas of the request to increase the chance of a prompt reply. Hedged requests are separated by a short time interval; the client issues the requests, accepts the first answer, and then notifies the other servers canceling the request. The client should wait a fraction, say 90% to 95% of the average response time, before sending the replica of the request to limit the additional workload and to avoid duplication.

Requests are tied when each replica includes the address of other servers wishing to reply to the request. In this case, the servers receiving the request communicate with one another; the first server able to start processing the request sends a canceling message to the other servers at the time it starts execution. If the input queue of all recipients of the request is empty, then all can start processing at about the same time and canceling messages criss-cross the network that are unable to prevent work duplication. This undesirable situation is prevented if the client waits a time equal to twice the average message delay before sending a replica of the request.

Hedged requests are very effective. For example, a Google benchmark reads the key value of 1 000 key-value pairs stored in a large BigTable distributed over 100 servers. Hedged requests sent 10 ms after the original ones reduced the 99 percentile latency dramatically, from 1 700 to 74 ms, while sending only 2% more requests. Another Google benchmark shows the effect of the queries tied one millisecond to an idle cluster and the other query on a cluster running a batch job in the background. The BigTable data is not cached, and each file chunk has three replicas on different storage servers. Table 12.3 shows read latencies with no hedge and with hedged tied requests.

The results in Table 12.3 show first that hedged and tied requests work well not only when the system is lightly loaded but also at higher system loads. This also implies that the overhead of this mechanisms is low and adds only a minute workload to the system. These results also show the extent of the heavier tail; the difference in latency between 99% and 99.9% is significant for both the lightly and the heavily loaded system.

Micro-partitions, replication of items likely to cause imbalance, and latency-induced probation are long-term tail-tolerant techniques. Perfect load balancing in a large-scale system is practically unachievable for many reasons. These reasons include the difference in server performance, the dynamic nature of the workloads, and the impossibility of having an accurate picture of the global system state. A fairly intuitive approach is a fine-grained partitioning of resources on every server, thus the name micro-partitions.

These “virtual servers” are more nimble and able to process fine-grained units of work more quickly. Error recovery is also less painful since less work is lost in case of errors. For a long time, immovable data replication has been a method of choice for improved performance, and it is also widely used at Google. Intermediate servers in a hierarchically organized system can identify slow-responding servers and avoid sending them latency-critical tasks, while continuing to monitor their behavior.

12.12 Edge computing and Markov decision processes (R)

Edge computing is a distributed computing framework that brings enterprise applications closer to data sources, such as IoT devices or local edge servers. Edge computing harnesses growing in-device computing capability to provide deep insights and predictive analysis in near-real time. Edge computing and mobile edge computing on 5G networks enables more comprehensive data analysis, faster response times, and improved customer experiences.

Edge computing is effective as data is processed and analyzed closer to the point where it is created. Data does not traverse a network to be processed and latency is significantly reduced. Proximity to data source has obvious benefits: improved response times and lower bandwidth availability. It is estimated that, in 2025, 75% of enterprise data will be processed at the edge, instead of only 10% today.

The *follow-me cloud* [468] and the *mobile edge-cloud* are variations on the theme of cloudlets discussed in Section 12.9 that were conceived to reduce the end-to-end latency for cloud access. The mobile edge-cloud concept allows mobile devices to carry out computationally intensive tasks on stationary servers located in the small data centers distributed across the network and connected directly to base stations at the edge of the network.

The new twist is to support *dynamic service placement*, in other words, to allow computations initiated by a mobile device in motion to migrate from one mobile edge-cloud server to another following the movement of the mobile device [483,506]. Optimal service migration policies pose very challenging problems due to the uncertainty of the mobile device movements and the likely nonlinearity of the migration and communication costs. One method to address these challenges is to formulate the migration problem in terms of the Markov decision process discussed next.

Markov Decision Processes. Markov decision processes are *discrete-time* stochastic control processes used for a variety of optimization problems where the outcome is partially random and partially under the control of the decision maker. Markov decision processes extend Markov chains with choice and motivations; *actions* allow choices, and *rewards* provide motivation for actions.

The state of the process in slot t is s_t , and the decision maker may choose any action $a_t \in \mathcal{A}(s_t)$ available in that state. As a result of action a_t , the system moves to a new state s' and provides the reward $\mathcal{R}_{a_t}(s_t, s')$. *The next state s' depends only on the current state s_t and the action taken.* The probability that the system moves to state s' is given by the transition function $p_{a_t}(s_t, s')$.

A Markov decision process is a 5-tuple: $(\mathcal{S}, \mathcal{A}, \mathcal{P}(\cdot, \cdot), \mathcal{R}(\cdot, \cdot), \gamma)$ with

- \mathcal{S} —a finite set of system states;
- \mathcal{A} —a finite set of actions; $\mathcal{A}_{s_t} \in \mathcal{A}$ is the finite set of actions available in state $s_t \in \mathcal{S}$ in time slot t .
- \mathcal{P} —the set of transition probabilities; $p_{a_t}(s_t, s') = Pr(s(t+1) = s' | s_t, a_t)$ —the probability that action a_t in state s_t in time slot t will lead to state s' in time slot $t+1$.
- $\mathcal{R}_{a_t}(s_t, s')$ —the immediate reward after the transition from state s_t to state s' .
- $\gamma \in [0, 1]$ —a discount factor representing the difference in importance between present and future rewards.

The goal is to optimize a policy π maximizing a cumulative function of the random rewards, e.g., the expected discounted sum of rewards over an infinite time horizon is

$$\sum_{t=0}^{\infty} \gamma^t \mathcal{R}_{a_t}(s_t, s_{t+1}). \quad (12.11)$$

To calculate the optimal policy given \mathcal{P} and \mathcal{R} , the state transitions and, respectively, the rewards, one needs two arrays $\mathcal{V}(s)$ and $\pi(s)$ indexed by state, the value, and policies, respectively. The first array contains values and the second actions:

$$\pi(s) = \arg \max_{\alpha} \left\{ \sum_{s'} \mathcal{P}_{\alpha}(s, s') (\mathcal{R}_{\alpha}(s, s') + \gamma \mathcal{V}(s')) \right\} \quad (12.12)$$

$$\mathcal{V}(s) = \sum_{s'} \mathcal{P}_{\pi(s)}(s, s') (\mathcal{R}_{\pi(s)}(s, s') + \gamma \mathcal{V}(s')). \quad (12.13)$$

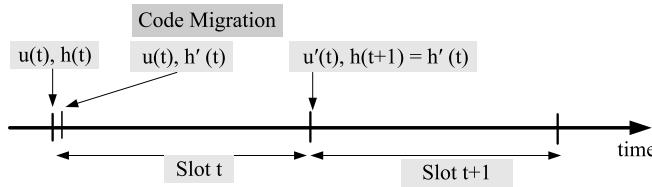
In the value induction proposed by Bellman, the calculation of $\pi(s)$ is substituted in the calculation of $\mathcal{V}(s)$:

$$\mathcal{V}_{i+1}(s) = \max_{\alpha} \left\{ \sum_{s'} \mathcal{P}_{\alpha}(s, s') (\mathcal{R}_{\alpha}(s, s') + \gamma V_i(s')) \right\}. \quad (12.14)$$

A Markov decision process can be solved by linear programming or by dynamic programming.

Migration decisions and cost in mobile cloud edge. The solution proposed in [483,506] assumes that:

1. The user's location in each slot is the same and changes from one slot to the next are according to the Markovian model discussed in this section; see Fig. 12.11.
2. The set of all possible locations, \mathcal{L} , can be represented as a 2D vector, and the distance between two locations $l_1, l_2 \in \mathcal{L}$ is given by $\|l_1 - l_2\|$.
3. The migration time is very small and will be neglected.

**FIGURE 12.11**

Time slots for the edge cloud. At the very beginning of slot t the user and service locations are $u(t)$ and $h(t)$, respectively. The service migrates to $h'(t)$, as soon as slot t starts and during the slot t the system operates with $u(t)$ and $h'(t)$. At the beginning of slot $t + 1$, the service location is $h(t + 1) = h'(t)$.

The following notations are used in [506]:

- $u(t)$ —users's location in slot t .
- $h(t)$ —service location in slot t .
- $d(t)$ —the distance between the user and the service; $d(t) = \|u(t) - h(t)\|$ in slot t .
- N —the maximum allowed distance between the user and the service, $N = \max d(t)$.
- $s(t) = (u(t), h(t))$ —the initial system state at the beginning of slot t . The initial state is $s(0) = s_0$.
- $c_m(x)$ —the migration cost is a nondecreasing function of x , the distance between the two edge cloud servers, $x = \|h(t) - h'(t)\|$.
- $c_d(x)$ —the transmission cost is a nondecreasing function of x , the distance between the edge cloud server and the user x , with $x = \|u(t) - h'(t)\|$. Initially, $c_d(0) = 0$.
- π —the policy used for control decision based on $s(t)$.
- $a_\pi(s(t))$ —the control action taken under policy π when the system is in state $s(t)$.
- \mathcal{C}_{a_π} —the sum of migration and transmission costs incurred by a control $a_\pi(s(t))$ in slot t . $\mathcal{C}_{a_\pi} = c_m(\|h(t) - h'(t)\|) + c_d(\|u(t) - h'(t)\|)$.
- γ —discount factor, $0 < \gamma < 1$.
- \mathcal{V}_π —expected discount under policy π .

The Markov decision process controller makes a decision at the beginning of each slot. The decision could be:

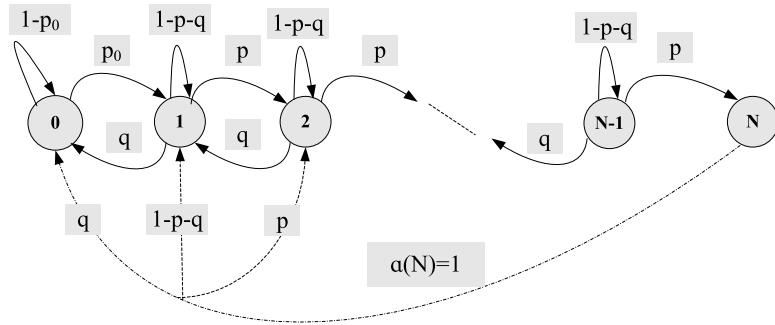
1. Do not migrate the service; then the cost is $c_m(x) = 0$.
2. Migrate the service from location $h(t)$ to location $h'(t) \in \mathcal{L}$; $c_m(x) > 0$.

Given a policy π , its long-term expected discounted sum cost is

$$\mathcal{V}_\pi(s_0) = \lim_{t \rightarrow \infty} \mathbb{E} \left\{ \sum_{\tau=0}^t \gamma^\tau \mathcal{C}_{a_\pi}(s(\tau)) \mid s(0) = s_0 \right\}. \quad (12.15)$$

An optimal control policy is one that minimizes $\mathcal{V}_\pi(s_0)$ starting from any initial state

$$\mathcal{V}^*(s_0) = \min_{\pi} \mathcal{V}_\pi(s_0), \quad \forall s_0. \quad (12.16)$$

**FIGURE 12.12**

Markov decision process state transitions assuming a 1-D mobility model for the edge cloud. The mobile device moves one step to the left or right with probability r_1 and stays in the same location with probability $1 - 2r_1$, thus $p = q = r_1$ and $p_0 = 2r_1$. Migration occurs in slot t if and only if $d(t) \geq N$. The actions in slot t with distance $d(t)$ are $\alpha(d(t)) = \alpha(N)$ for $d(t) > N$. This implies that we need only to study the states where $d(t) \in [0, N]$. After an action $\alpha(N)$, the system moves to states 0, 1, and 2 with probabilities q , $1 - p - q$, and p , respectively.

A stationary policy is given by Bellman's equation written as

$$\mathcal{V}^*(s_0) = \min_{\alpha} \left\{ \mathcal{C}_{\alpha}(s_0) + \gamma \sum_{s_1 \in \mathcal{L} \times \mathcal{L}} \mathcal{P}_{\alpha(s_0, s_1)} \mathcal{V}^*(s_1) \right\}, \quad (12.17)$$

with $\mathcal{P}_{\alpha(s_0, s_1)}$ being the transition probability from state $s'(0) = s_0 = \alpha(s_0)$ to $s(1) = s_1$. The intermediate state $s'(t)$ has no randomness when $s(t)$ and $\alpha(\cdot)$ are given.

A proposition reflecting the intuition that it is not optimal to migrate the service to a location that is farther away from the current location of the mobile device has the following intuitive implication, which simplifies the search for optimal policy [506]:

Proposition. *If $c_m(x)$ and $c_d(x)$ are constants and*

$$c_m(0) < c_m(x) \text{ and } c_d(0) < c_d(x) \text{ for } x > 0, \quad (12.18)$$

then the current user location is not optimal.

Fig. 12.12 shows the system transition model when the transition probabilities are p_0 , p , and q , assuming a uniform 1-D mobility model. In this model $u(t)$, $h(t)$, and $h'(t)$ are scalars. h' is the new service location chosen such that:

$$||h(t) - h'(t)|| = |d(t) - d'(t)| \quad \text{and} \quad ||u(t) - h'(t)|| = d'(t). \quad (12.19)$$

The migration occurs along the shortest path connecting $u(t)$, $h(t)$, and $h'(t)$.

An operating procedure is discussed in [506], and workload scheduling for edge clouds is presented in [483].

12.13 Bootstrapping techniques for data analytics (R)

The size of the data sets used for data analytics, as well as the complexity and the diversity of queries posed by impatient users, often without training in statistics, continually increase. In many instances, e.g., in case of exploratory queries, it is desirable to provide good enough, yet prompt answers, rather than perfect answers after a long delay.

This is only possible by limiting the search to a subset of data, but in such instances, the user expects an estimation of answer's quality along with the answer. The quality of the answer is context dependent, and a general solution is far from trivial. Bootstrapping techniques discussed in this section can be applied to a broad range of applications for estimating the quality of such approximations. Given a set F and a random variable U , the bootstrapping methods are based on the *bootstrap substitution principle*

To determine the probability distribution of $U \equiv u(Y, F)$ with $Y = \{Y_1, Y_2, \dots, Y_n\}$ random samples of F are taken and then F is replaced by a fitted model \hat{F} .

Thus, we make the approximation

$$\Pr\{u(Y, F) \leq u \mid F\} \approx \Pr\{u(Y^*, \hat{F}) \leq u \mid \hat{F}\}. \quad (12.20)$$

The superscript * distinguishes random variables and related quantities sampled from a probability model that has been fitted to the data. Sometimes, $u(Y, F) = T - \theta$ with T as the estimator of the parameter $\theta \equiv t(F)$, and more sophisticated cases involve transformations of T . Often, $T = t(\tilde{F})$, and \tilde{F} is an empirical distribution function of the data values.

While bootstrapping can produce reasonably accurate results, there are also instances when the results are unreliable. Some of such instances discussed in depth in [86] are:

1. Inconsistency of the bootstrap method when the model, the statistic, and the resampling fail to approximate the required properties, regardless of the sample size.
2. Incorrect resampling model in case of nonhomogeneous data when the random variation of data is incorrectly modeled.
3. Nonlinearity of statistic T as the good properties of the bootstrap are associated with an accurate linear approximation $T \approx \theta + n^{-1} \sum_i l(T_i)$, with $l(y)$ the influence function of $t(\cdot)$ at (y, F) .

A key idea is to construct a proxy to ground truth for a sample smaller than that of the complete observed data set and compare the bootstrap's results with this proxy. The bootstrapping techniques discussed in this section are based on [275], which states “existing diagnostic methods target only specific bootstrap failure modes, are often brittle or difficult to apply, and generally lack substantive empirical evaluations ...this paper presents a general bootstrap performance diagnostic which does not target any particular bootstrap failure mode but rather directly and automatically determines whether or not the bootstrap is performing satisfactorily.”

The bootstrap method. Let P be an unknown distribution, $\theta(P)$ be some parameter of P , and \mathcal{D} the set of n independent identically distributed (i.i.d.) sampled data points $\mathcal{D} = \{X_1, X_2, \dots, X_n\}$ from P .

Let $\mathbb{P} = n^{-1} \sum_{i=1}^n \delta_{X_i}$ be the empirical distribution of data. We wish to construct the estimate $\hat{\theta}(\mathcal{D})$ of $\theta(P)$ and then create $\xi(P, n)$, an assessment of the quality of $\hat{\theta}(\mathcal{D})$, consisting of a summary of the distribution Q_n of some quantity $u(\mathcal{D}, P)$.

Both P and Q_n are unknown; thus the estimate $\xi(P, n)$, called the *ground truth* in this discussion, cannot be computed directly; it can be approximated by $\xi(\mathbb{P}_n, n)$ using Monte Carlo procedure. The following steps are carried out repeatedly:

1. Form simulated data sets \mathcal{D}^* of size n consisting of i.i.d. sampled points from \mathbb{P}_n ;
2. Compute $u(\mathcal{D}^*, \mathbb{P}_n)$ for the simulated data set \mathcal{D}^* ;
3. Form the empirical distribution Q_n of the computed values of u ;
4. Return the desired summary of this distribution.

The final bootstrap output will be a real-valued $\xi(Q_n, n)$. The assessment $\xi(P, n)$ could compute:

1. The bias, the expectation of

$$u(\mathcal{D}, P) = \hat{\theta}(\mathcal{D}) - \theta(P). \quad (12.21)$$

2. A confidence interval based on the distribution of

$$u(\mathcal{D}, P) = n^{1/2}[\hat{\theta}(\mathcal{D}) - \theta(P)]. \quad (12.22)$$

3. Simply

$$u(\mathcal{D}, P) = \hat{\theta}(\mathcal{D}). \quad (12.23)$$

Given the estimator, the data generating distribution P , and n , the size of the data set, we want to determine if the output of the bootstrap procedure is *sufficiently close* to the ground truth. The formulation *sufficiently close* avoids a precise expression of accuracy, giving the procedure a degree of generality and allowing its use for a range of applications with own accuracy requirements.

In practice, we can observe only one set with n data points rather than many independent sets of size n . The solution is to randomly sample $p \in \mathbb{N}$ disjoint subsets of the dataset \mathcal{D} , each of size $b \leq \lfloor n/p \rfloor$. Then, to approximate the distribution of Q_b , we use the set of values of u calculated for each subset. This distribution yields an approximation of the ground truth, $\xi(P, b)$ for data sets of size b . Then, to determine if the bootstrap performs as expected on sample size b , we run the bootstrap on each of the p subsets and compare the p bootstrap outputs to the ground truth.

Carrying out this procedure for a single sample size b is insufficient, the bootstrap performance may be acceptable for small sample size, but it may get worse as the sample size increases or, conversely, be mediocre for small sizes, but improve as the sample size increases. Thus, it is necessary to compare the distribution of bootstrap outputs for a range of sample sizes, b_1, b_2, \dots, b_k , $b_k \leq \lfloor n/p \rfloor$. If the distribution of the bootstrap outputs converges monotonically for all smaller sample sizes b_1, b_2, \dots, b_k , we conclude that the bootstrap is performing satisfactorily for size n .

Convergence criteria are based on the relative deviation of the absolute value and the size of the standard deviation of the bootstrap output from the ground truth. The pseudocode of the Bootstrap Performance Diagnostic (BPD) algorithm illustrates these steps.

BPD Algorithm—Bootstrap Performance Diagnostic [275]

Input: $\mathcal{D} = \{X_1, \dots, X_n\}$: observed data

- u : quantity whose distribution is summarized to yield estimator quality assessments
- ξ : assessment of estimator quality
- p : number of disjoint subsamples used to compute ground truth approximations
- b_1, \dots, b_k : increasing sequence of subsample sizes for which ground truth approximations are computed with $b_k \leq \lfloor n/p \rfloor$ (e.g., $b_i = \lfloor n/(p2^{k-i}) \rfloor$ with $k = 3$)
- $c_1 \geq 0$: tolerance for decreases in absolute relative deviation of mean bootstrap output
- $c_2 \geq 0$: tolerance for decreases in relative standard deviation of bootstrap output
- $c_3 \geq 0, \alpha \in [0, 1]$: desired probability that bootstrap output at sample size n has absolute relative deviation from ground truth less than or equal to c_3 (e.g., $c_3 = 0.5; \alpha = 0.95$)

Output: *true* if the bootstrap is deemed to be performing satisfactorily, and *false* otherwise

$\mathbb{P}_n \rightarrow n^{-1} \sum_{i=1}^n \delta_{X_i}$

for $i \leftarrow 1$ **to** k **do**

- $\mathcal{D}_{i1}, \dots, \mathcal{D}_{ip} \rightarrow$ random disjoint subsets of \mathcal{D} , each containing b_i data points
- **for** $j \leftarrow 1$ **to** p **do**
- $u_{ij} \leftarrow u(\mathcal{D}, \mathbb{P}_n)$
- $\tilde{\xi}_{ij}^* \leftarrow \text{bootstrap}(\xi, u, b_i, \mathcal{D}_{ij})$
- **end**
- // Compute ground truth approximation for sample size b_i
- $\mathbb{Q}_{b_i} \leftarrow \sum_{j=1}^p \delta_{u_{ij}}$
- $\tilde{\xi}_i \leftarrow \xi(\mathbb{Q}_{b_i}, b_i)$
- // Compute absolute relative deviation of mean and relative standard deviation
- // of bootstrap outputs for sample size b_i
- $\Delta_i \leftarrow \left| \frac{\text{mean}\{\tilde{\xi}_{i1}^*, \dots, \tilde{\xi}_{ip}^*\} - \tilde{\xi}_i}{\tilde{\xi}_i} \right|$ | $\sigma_i \leftarrow \left| \frac{\text{stddev}\{\tilde{\xi}_{i1}^*, \dots, \tilde{\xi}_{ip}^*\} - \tilde{\xi}_i}{\tilde{\xi}_i} \right|$

end

return true if all of the following hold, and false otherwise

$$\Delta_{i+1} < \Delta_i \quad \text{OR} \quad \Delta_{i+1} \leq c_1, \forall i = 1, \dots, k \quad (12.24)$$

$$\sigma_{i+1} < \sigma_i \quad \text{OR} \quad \sigma_{i+1} \leq c_2, \forall i = 1, \dots, k \quad (12.25)$$

$$\frac{\# \left(j \in \{1, \dots, p\} : \left| \frac{\tilde{\xi}_{kj}^* - \tilde{\xi}_k}{\tilde{\xi}_k} \right| \leq c_3 \right)}{p} \geq \alpha \quad (12.26)$$

This algorithm generates a confidence interval with coverage $\alpha \in [0, 1]$.¹¹ A false positive, deciding that an approximation is satisfactory when it is not, is less desirable than a false negative, rejecting

¹¹ Confidence intervals offer a guarantee of the quality of a result. A procedure is said to generate confidence intervals with a specified coverage $\alpha \in [0, 1]$ if, on a proportion exactly α of the set of experiments, the procedure generates an interval that includes the answer. For example, a 95% confidence interval $[a, b]$ means that in 95% of the experiments, the result will be in $[a, b]$.

an approximation when in fact it is a satisfactory one. Eq. (12.26) reflects this conservative approach. The absolute value of the deviation of a quantity γ from γ_0 is defined as $|\gamma - \gamma_0| / |\gamma_0|$, and an approximation is satisfactory if the output of the bootstrap run on a data set of size n has an absolute value of the relative deviation from the ground truth of at most c_3 with a probability $\alpha \in [0, 1]$.

Choosing the sample sizes close together or choosing c_1 or c_2 too small will cause a larger number of false negatives. It is recommended to use an exponential distribution of sample sizes to ensure a meaningful comparisons of bootstrap performance for consecutive values b_i, b_{i+1} in the set $\{b_1, \dots, b_i, b_{i+1}, \dots, b_k\}$.

The process discussed in this section requires a substantial amount of data, but this does not seem to be a problem in the age of Big Data. For example, according to [275] when $p = 100$ and $b_k = 1000$ then $n \geq 10^{15}$ and if b_k increases, $b_k = 10000$, $n \geq 10^6$. Processing such large data sets requires significant resources.

Simulation experiments for several distributions including $Normal(0; 1)$, $Uniform(0; 10)$, $StudentT(1.5)$, $StudentT(3)$, $Cauchy(0; 1)$, $0.95Normal(0; 1) + 0.05Cauchy(0; 1)$, and $0.99Normal(0; 1) + 0.01Cauchy(104; 1)$ are reported in [275]. The results of these experiments show that the diagnostic performs well across a range of data generating distributions and estimators. Moreover, its performance improves with the size of the sample data sets.

In summary, given \mathcal{S} a random sample from \mathcal{D} , the subsamples generated by disjoint partitions of \mathcal{S} are also mutually independent random samples from \mathcal{D} . The procedure must be carried out using a sequence of samples of increasingly larger size, $b_1, \dots, b_i, \dots, b_k$. It is necessary to ensure that the error decreases while increasing the sample size and that the error is sufficiently small for the largest sample.

12.14 Approximate query processing (R)

The advantages of executing a query on a data sample of rather than on an entire dataset are quite perspicuous, and this idea has been applied to sampling relational databases since the 1970s. Different versions of this technique have been investigated since its first use. Approximate query processing and sampling-based approximate query processing have become popular enough to warrant the acronyms, AQP and S-AQP, respectively, along with the realization that approximate answers are most useful if accompanied by accuracy estimates.

Let θ be a query on a dataset \mathcal{D} and the desired answer to it be $\theta(\mathcal{D})$. Simple random sampling with plug-in estimation is often used for approximate query processing. This method generates a sample with replacement $\mathcal{S} \subset \mathcal{D}$ of cardinality $n = |\mathcal{S}| \leq |\mathcal{D}|$ and produces a *sample estimate result* $\theta(\mathcal{S})$ instead of computing $\theta(\mathcal{D})$ with the *sampling error* $\epsilon = \theta(\mathcal{S} - \mathcal{D})$ and the *sampling error distribution* $Dist(\epsilon)$.

The emergence of Big Data processed on large clusters in the cloud and the requirement of near real-time response have increased the demand for high-quality error estimates. Such error estimates can be reported to the users enabling them to judge the impact of errors on their specific applications and/or to application developers to decide if the sampling methods are adequate. These estimates can also be used to correlate the errors with the sample size necessary for the accuracy–response time tradeoffs.

Two methods for producing close-form estimates error bars are based on the Central Limit Theorem (CLT) and on Hoeffding bounds [220]. An error bar is a line segment through a point on a graph,

parallel to one of the axes, that represents the uncertainty or error of the corresponding coordinate of the point. Informally, CLT states that the sum of a large number of independent random variables has a normal distribution. More precisely, if $\{X_1, \dots, X_n\}$ is a sequence of i.i.d. random variables with $E[X_i] = \mu$ and $Var[X_i] = \sigma^2 < \infty$ and $S_n = 1/n \sum_{i=1}^n X_i$ is the sample average, then the random variables $\sqrt{n}(S_n - \mu)$ converge to a normal distribution, $N(0, \sigma^2)$, as n approaches infinity:

$$\sqrt{n} \left[\left(\frac{1}{n} \sum_{i=1}^n X_i \right) - \mu \right] \xrightarrow{d} N(0, \sigma^2). \quad (12.27)$$

The derivation of Hoeffding bounds starts with the estimation of the mean

$$\mu = \frac{1}{m} \sum_{i=1}^m v(i), \quad (12.28)$$

with v as a real-valued function defined on the set $\mathcal{S} = \{1, 2, \dots, m\}$ and $m > 1$ as a fixed integer. Let L_1, L_2, \dots, L_n and L'_1, L'_2, \dots, L'_n be random samples from the set \mathcal{S} with and without replacement, respectively. Given $n > 1$, let \bar{Y}_n and \bar{Y}'_n be two estimators for μ be defined as

$$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n v(L_i) \quad \text{and} \quad \bar{Y}'_n = \frac{1}{\min(n, m)} \sum_{i=1}^{\min(n, m)} v(L'_i). \quad (12.29)$$

These estimators are unbiased if $E[\bar{Y}_n] = E[\bar{Y}'_n]$. Hoeffding showed that, for any $n \geq m$ and convex function f , $E[f(\bar{Y}'_n)] \leq E[f(\bar{Y}_n)]$. It follows that when $f(x) = x^2 - \mu$,

$$Var[f(\bar{Y}'_n)] \leq Var[f(\bar{Y}_n)]. \quad (12.30)$$

These results are useful to bound the probability that the estimates deviate from μ more than a given amount. If the only information available before sampling is that

$$a \leq v(i) \leq b, \quad 1 \leq i \leq m, \quad (12.31)$$

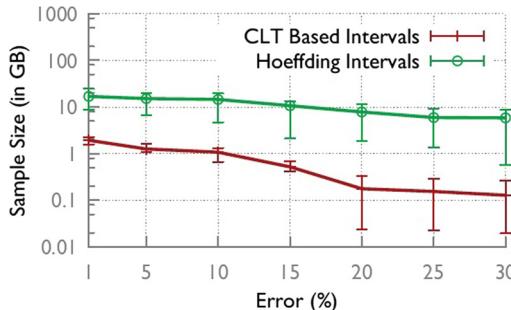
Hoeffding established the following bounds for the estimation error:

$$P\{|\bar{Y}_n - \mu| \geq t\} \leq 2e^{-2nt^2/(b-a)^2} \quad \text{and} \quad P\{|\bar{Y}'_n - \mu| \geq t\} \leq 2e^{-2n't^2/(b-a)^2} \quad (12.32)$$

for $t > 0$, $n \geq 1$ and

$$n' = \begin{cases} n & \text{if } n < m; \\ +\infty & \text{if } n \geq m. \end{cases} \quad (12.33)$$

Hoeffding bounds overestimate the error and increase the computational effort. The sample size for a system using Hoeffding bounds are one to two orders of magnitude larger than for CLT-based or bootstrap-based methods, but their accuracy is significantly higher. Experiments to estimate the sample sizes for achieving different levels of relative errors carried out on tens of terabytes of data are reported in [12] and reproduced in Fig. 12.13.

**FIGURE 12.13**

Sample size for CLT- and Hoeffding-based estimation for error bars [12].

Investigation of nonparametric bootstrap, closed-form estimation of variance, and large deviation bounds confirms that all the techniques have different failure modes. A benchmark, consisting of 100 different samples of some 10^6 rows, reported by [12], shows that only queries with COUNT, SUM, AVG, and VARIANCE aggregates are amenable to closed-form error estimation. All aggregates are amenable to the bootstrap, and 43.21% of the queries over one dataset and 62.79% of the queries over another can only be approximated using bootstrap-based error estimation methods. As expected, queries involving MAX and MIN are very sensitive to rare large or small values, respectively. In one data set, these two functions involved 2.87% and 33.35% of all queries, respectively. Bootstrap error estimation fails for 86.17% of these queries.

BEC: Bootstrap Error Computation for a SELECT query on \mathcal{S}

```

SELECT foo(col_<math>\mathcal{S}</math>), <math>\hat{\xi}(\text{resample\_error})</math> AS error
FROM (
    . . . .
    . . . .
    . . . .
    . . . .
    . . . .
)
    UNION ALL
    SELECT foo(col_<math>\mathcal{S}</math>) AS resample_answer
    FROM <math>\mathcal{S}</math> TABLE_SAMPLE POISSONIZED (100)
    UNION ALL
    SELECT foo(col_<math>\mathcal{S}</math>) AS resample_answer
    FROM <math>\mathcal{S}</math> TABLE_SAMPLE POISSONIZED (100)
    UNION ALL
    . . .
    . . .
    . . .
    . . .
    . . .
    . . .
    . . .
    . . .
    . . .
    . . .
)
    UNION ALL
    SELECT foo(col_<math>\mathcal{S}</math>) AS resample_answer
    FROM <math>\mathcal{S}</math> TABLE_SAMPLE POISSONIZED (100)
)
```

The diagnostic algorithm discussed in Section 12.13 is computationally intensive, thus impractical in many cases. A technique to quickly determine if a technique will work well for a particular query based on symmetrically centered confidence intervals was proposed in [12]. The workflow for a large-scale distributed approximate query processing proposed has several steps:

Logical Plan (LP)—a query is compiled into an LP consisting of three procedure to compute: (1) an approximative answer $\theta(\mathcal{S})$; (2) error $\bar{\xi}$; (3) diagnostic tests;

Physical Plan (PP)—initiates a DAG of tasks involving these procedures;

Data Storage Layer—distributes samples to a set of servers and manages cached data.

The system supports Poissonized resampling, thus allowing a straightforward implementation of the bootstrap error. For example, for a simple query of the form *SELECT foo(col_S) FROM S*, the bootstrap error is computed by using the BEC pseudocode in the box.

An important source of inefficiency of the bootstrap method is the execution of the same query on various samples of the data. *Scan consolidation* can eliminate this source of inefficiency. As a first step of this process, the Logical Plan is optimized by extending the resampling operations. Each tuple in the sample \mathcal{S} is extended with a set of 100 independent weights w_1, w_2, \dots, w_{100} drawn from a Poisson(1) distribution. The sample \mathcal{S} is partitioned into multiple sets of $a = 50$, $b = 100$ and $c = 200$ MB, and three sets of weights, D_{a1}, \dots, D_{a100} , D_{b1}, \dots, D_{b100} and D_{c1}, \dots, D_{c100} , are associated with each row to create 100 resamples of each set.

The logical plan is rewritten for further optimization. After finding the longest set of consecutive operators that do not change the statistical properties of the set of columns that are being finally aggregated,¹² the custom Poissonized resampling operator is inserted right before the first nonpassthrough operator in the query graph. The subsequent aggregate operators are modified to compute a set of resample aggregates by appropriately scaling the corresponding aggregation column with the weights associated with every tuple. Further optimization of the cache management is used to improve the performance of the procedure discussed in [12].

12.15 Further readings

Several references, including [3,302,532,533], discuss defining characteristics of Big Data applications. Insights into Google’s storage architecture for Big Data can be found in [173], and several systems, including Mesa, Spanner, and F1, are presented in [217], [117], and [447].

Data analytics is the subject of [528], and [104] analyzes interactive analytical processing in Big Data systems. [208] covers in-memory performance of Big Data. Continuous pipelines are discussed in [142]. The Starfish system for Big Data analytics is covered in [234], and [252] analyzes enterprise use of Big Data. Several papers including [86] and [275] cover bootstrapping techniques, and [12] analyzes approximate query processing. Hoeffding bounds are discussed in [220]. Several references such as [140,141,305] cover Dynamic Data-Driven Application Systems (DDDAS) and the FreshBreeze system.

Spark Streaming system is discussed in [534], and [15] covers the MillWheel framework developed at Google for building fault-tolerant and scalable data streaming systems. [419] presents caching strategies for data streaming. [24] covers Google’s Photon system and system scalability, and the performance of large-scale system is the topic of [218]. The problems posed by the heavy-tail distribution of latency are analyzed in [130].

¹² These so-called *passthrough operators* could be scans, filters, projections, and so on.

Mobile devices and applications are covered in the literature, including [124,148,437,438,505]. Energy efficiency of mobile computing is analyzed in [347]. The use of mobile devices for space weather monitoring is discussed in [387]. The Follow-Me cloud and edge cloud computing are the subjects of [468] and [483,506].

12.16 Exercises and problems

- Problem 1.** Read [183] and analyze the benefits and the problems related to *dataspaces*. Research the literature for potential applications of dataspaces to data management of data-intensive applications in science and engineering.
- Problem 2.** Discuss the possible solution for stabilizing cloud services mentioned in [182] inspired by BGP routing [209,491].
- Problem 3.** Discussing the bootstrap method presented in Section 12.13, reference [275] states: “Ideally, we might approximate $\xi(P, n)$ for a given value of n by observing many independent datasets, each of size nin practice we only observe a single set of n data points, rendering this approach an unachievable ideal. To surmount this difficulty, our diagnostic, the BPD Algorithm, executes this ideal procedure for dataset sizes smaller than n . That is, for a given $p \in \mathbb{N}$ and $b \leq \lceil n/p \rceil$ we randomly sample p disjoint subsets of the observed dataset \mathcal{D} , each of size b . For each subset, we compute the value of u ; the resulting collection of u values approximates the distribution Q_b , in turn yielding a direct approximation of $\xi(P, b)$ the ground truth value for the smaller dataset size b . Additionally, we run the bootstrap on each of the p subsets of size b , and comparing the distribution of the resulting p bootstrap outputs to our ground truth approximation, we can determine whether or not the bootstrap performs acceptably well at sample size b .” (1) Estimate the bootstrap speedup function of n, p, b ; ignore the time to compute Δ_i, σ_i . (2) Examine the simulation results in [275] and discuss the impact of the sample size.
- Problem 4.** Read [12] and discuss the merits and the shortcomings of the three techniques for estimating the sample distribution using only a single sample: nonparametric bootstrap, close-form estimation, and large deviation bounds.
- Problem 5.** What is a key extraction function, and what role does it play in MillWheel? Give an example of such a key extraction function in Zeitgeist.
- Problem 6.** In systems with very high fan-out, a request may exercise an untested code path, causing crashes or extremely long delays on thousands of servers simultaneously. How can this problem be prevented?
- Problem 7.** Consider a system where each server typically responds in 10 ms but with a 99th-percentile latency of one second. If a user request must collect responses from 100 such servers in parallel, then 63% of user requests will take more than one second [130]. Assuming that there are 1 000 servers and the user request needs responses from 1 000 servers running in parallel, what is the probability that response latency will be larger than one second? What if, instead of 1 000 individual requests, the user request needs data from 2 000 servers running in parallel?

Problem 8. Research the power consumption of processors used in mobile devices and their energy efficiency. Rank the components of a mobile device in terms of power consumption. Establish a set of guidelines to minimize the power consumption of mobile applications.

Problem 9. What is the main benefit of Algorithm 1 for finding the optimal policy of a Markov Decision Process in [505] compared with the standard approaches?

Emerging clouds

13

In this chapter, we explore several areas of cloud computing research impacting the clouds of the future. Past history is often a good predictor of the future, and the past shows that cloud computing has promptly absorbed and exploited new developments in science and engineering for the benefit of an increasingly larger user community. At the same time, cloud computing has forced significant revisions of traditional concepts and ideas in many areas of computer science and computer engineering. Thus, it seems safe to expect that these processes will continue in predictable, as well as unpredictable, ways.

Arguably, Google Research is the most reliable and active source of cloud computing research publications. To identify the most disruptive ideas affecting cloud computing, we compared the publications of Google researchers over the last five years in two traditional areas, Hardware and Architecture (HA) and Data Mining (DM) with three emerging areas, Machine Intelligence (MI), Machine Perception (MP), and Quantum Computing (QC). Table 13.1 indicates that the interests of Google researchers have evolved with fewer publications about traditional areas and an increasing number of publications in AI/ML and QC.

The chapter starts with a brief discussion of the obvious challenges cloud computing is expected to address in Section 13.1. Artificial Intelligence (AI) and Machine Learning (ML) are the most impactful areas of research for virtually all aspects of human activities requiring computing. Quantum computing provides a glimmer of hope and, hopefully, a practical alternative to silicon as the effective end of Moore's law is in sight. The data in Table 13.1 justify why AI and ML applications and quantum computing on the cloud are covered next, in Sections 13.2 and 13.3.

Thousands of cars in parking garages or airport parking lots can contribute to edge cloud computing if their owners are offered adequate incentives. Vehicular clouds is a topic of practical interest because powerful systems on a chip are already used in cars to support advanced safety features and Level 0 ("basic safety"), 1 ("hands on"), 2 ("hands off"), using the SAE classification of driving automation

Table 13.1 The number of papers published by Google researchers in two traditional areas, Hardware and Architecture (HA) and Data Mining (DM) and three emerging areas, Machine Intelligence (MI), Machine Perception (MP), and Quantum Computing (QC), during the period January 2015–July 2020 and ALL, i.e., since Google was founded in 1998.

Area	2015	2016	2017	2018	2019	2020	All
HA	6	7	5	13	7	3	99
DM	22	21	18	24	32	14	284
MI	100	169	270	386	498	189	2 051
MP	40	57	99	158	151	58	807
QC	3	7	8	14	11	6	56

levels. Prof. Stefan Olariu, one of the pioneers of this field, has graciously consented to overview this topic in Section 13.4. Vehicular clouds support an alternative type of edge computing.

13.1 A short-term forecast

Cloud functionality is evolving as new services and more diverse and powerful instance types are released every year. For example, Amazon introduced Lambda, a service in which applications are triggered by user-defined conditions and events. In late 2016, AWS added P2, a powerful, scalable instance with GPU-based parallel compute capabilities and a range of services supporting machine learning and AI. Google has been adding to the software stack for coarse-grained parallelism based on MapReduce and adding TPUs (Tensor Processing Units) to its cloud infrastructure. IBM's efforts target computational intelligence as evidenced by the success of Watson in healthcare and data analytics. Microsoft is attempting to extend the range of commercial applications supported by its growing cloud infrastructure.

The cloud research community will most likely be faced with new and challenging problems in the future. *Variability* and *conflicting requirements* are disruptive qualities of cloud computing demanding new thinking in system design. The physical infrastructure consisting of servers with different architecture and performance is constantly updated as solid-state technologies and processor architecture evolve. New software orchestrating a coherent system view is developed every month, and the depth of the software stack is continually increasing. If cloud computing is to continue to be successful, it is very likely that new applications will emerge. The size of the cloud user population and the diversity of their requirements will grow.

Conflicting system design requirements are not a novelty, but the depth and the breadth of such conflicts is unprecedented and qualitatively different due to scale of the cloud infrastructure and application diversity. Many contradictory requirements in the design of computer clouds have to be carefully balanced. For example, resource sharing is a basic design principle, yet strict performance and security isolation of cloud application are also critical. To deliver cheap computing cycles, the cloud infrastructure should maximize resource utilization, while sufficient resources should be kept in reserve to respond to large workload spikes. The system as a whole should present itself as flawless and extremely reliable though the failure rate of cheap, off-the-shelf hardware components can be fairly high. Performance guarantees should be provided, while a mix of workloads with very different requirements will continue to dynamically share system resources.

Unquestionably, computer clouds will continue to evolve, but how? A parallel of clouds with the Internet is unavoidable. Initially, the Arpanet, the precursor of the Internet, was a network used to transfer data files from one location to another. It was designed as a best-effort network doing its best to transport data packets without providing end-to-end delivery guarantees. Support for communication with real-time delivery constraints was not foreseen. The Internet's success forced changes, and, as a result, today's Internet supports low-latency and high-bandwidth data streaming. Data-streaming traffic is *shaped*¹ to guarantee that routers have enough resources to transmit a continuous stream of data with low jitter.

¹ Traffic shaping is a bandwidth management technique used on datagram networks that delays some or all datagrams to bring them into compliance with a desired traffic profile.

How will computer clouds simultaneously provide QoS guarantees, increase resource utilization, support elasticity, and be more secure? This is the question. Cloud evolution poses fundamental questions that deserve further research. A first question is whether applications with special requirements, such as real-time constraints or applications exhibiting fine-grained parallelism, could migrate to the cloud. Such a migration requires changes in software, in particular, in resource management and scheduling components of the software stack. At the same time, the hardware and, in particular, the cloud interconnection networks have to offer lower latency and higher bandwidth.

Another question is how to reduce cost by increasing resource utilization, without affecting the QoS promises of cloud computing. It is self-evident that cloud elasticity cannot be supported without some form of over-provisioning which implies lower average resource utilization. In conclusion, alternative means to reduce cost should be considered.

A solution practiced by AWS and others is to combine a reservation system with spot allocations. Spot allocations are designed to consume excess resources if and when such resources are available. Cloud users with a good understanding of the resources needed and the time required by their applications should use the reservation system. They will benefit from QoS guarantees and pay more for cloud services. The other cloud users should compete for lower-cost spot allocations.

An alternative that may be explored in the future is to profile the cloud users using large databases of historic data and using machine learning algorithms and data analytics to predict the resource needs and the time required to carry out the computation. Once this information is available, a virtual private cloud that best fits the profile of the application should be configured. Another alternative is to support cloud self-organization and self-management. Market-based resource allocation could be at the heart of such an approach in spite of the problems it poses [446], including auctions as suggested in [333,334].

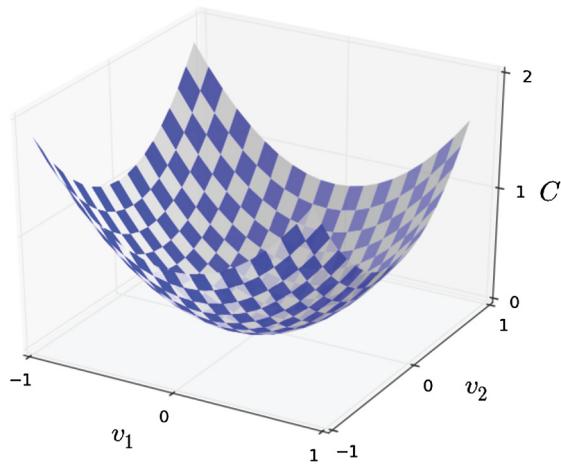
Homogeneity of the cloud computing infrastructure was one of the fundamental design tenets. The obvious advantages of infrastructure homogeneity are simplification of resource management and lowering hardware and software maintenance costs. Also, acquiring large volumes of identical hardware components lowers the infrastructure cost.

In the last few years CSPs realized why cloud users are clamoring for heterogeneity. As a result, today's clouds have multiple types of processors and co-processors, e.g., GPUs and TPUs and clouds use solid-state disks for real-time applications. The cloud infrastructure may include data flow engines in the near future. It is also likely that islands of systems communicating through InfiniBand, or other high-performance networks will be part of the cloud computing landscape. Such cloud islands will allow fine-grained parallel scientific and engineering applications to perform well on computer clouds.

It is necessary to address the question of how to accommodate cloud heterogeneity, while preventing a dramatic increase of infrastructure cost and an increase of resource management policies and mechanisms implementing these policies. Market mechanisms, which proved so successful in dealing with a diverse set of goods and a large consumer population in a free-market economy, could provide the answer. Cloud computing will continue to have a profound influence on the very large number of individuals and institutions who are now empowered to process huge amounts of data.

13.2 Machine learning on clouds

Machine learning (ML) uses algorithms that learn to discover how to perform a specific task without being explicitly programmed to do so. ML is considered by some an area of artificial intelligence (AI),

**FIGURE 13.1**

Graduate descent used to find a global minimum.

while others believe that ML is a field in itself. There is general agreement that ML algorithms refine a mathematical model of a system/environment using *training data* in several modes:

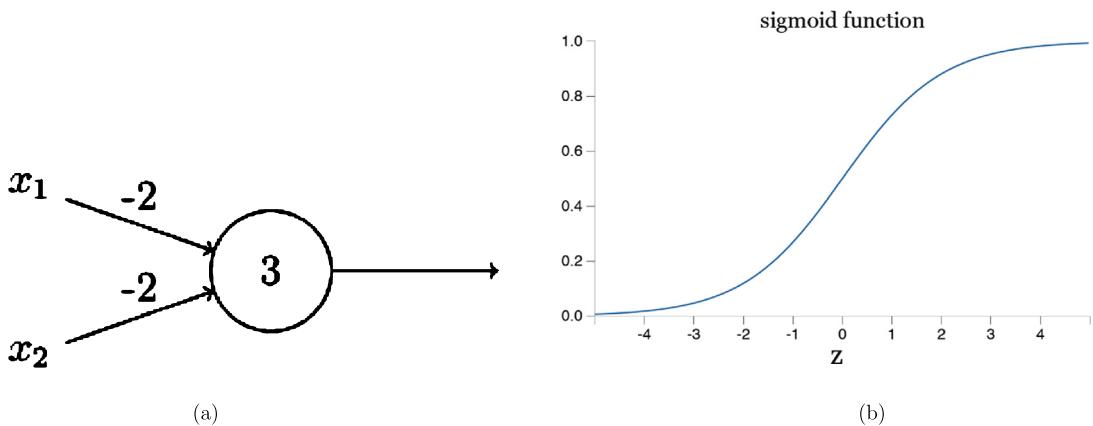
1. Supervised learning—the goal is to learn general rules that map inputs to outputs given example inputs and the desired output for each input provided by a “supervisor.”
2. Unsupervised learning—the goal is to discover hidden patterns in input data supplied without labels or the use the unstructured input as means to achieve an end, the so-called feature learning.
3. Reinforcement learning—the goal is to perform a specific task through interactions with an environment that provides feedback after each attempt and, in the process, to maximize a rewards function.

Supervised learning uses stochastic *gradient descent* and adjusts the weights to maximize the descent of the gradient obtained by back-propagation, as shown in Fig. 13.1. *Back-propagation* repeatedly adjusts the weights of the connections in the network to minimize a measure of the difference between the actual output vector of the net and the desired output vector [425].

Given a cost function $C(v_1, v_2, \dots)$, the change of the cost, ΔC , produced by a small $\Delta V = (\Delta V_1, \Delta V_2, \dots)^T$ is $\Delta C \approx \nabla C \times \Delta V$, with the gradient ∇C given by the vector

$$\nabla C \equiv \left(\frac{\partial C}{\partial v_1}, \dots, \frac{\partial C}{\partial v_m} \right)^T. \quad (13.1)$$

Neural networks. ML uses *neural networks* (NNs), inspired by networks of human brain neurons, to learn how to accomplish a desired task. NNs were developed in the 1950s and 1960s by Frank Rosenblatt, who introduced the concept of *perceptron* [422]. Rosenblatt was inspired by earlier work by Warren McCulloch and Walter Pitts [342]. A perceptron takes several binary inputs x_1, x_2, \dots and produces a single binary output. Rosenblatt introduced weights, w_1, w_2, \dots , real numbers expressing the importance of the respective inputs to the output. The binary output of a perceptron is determined

**FIGURE 13.2**

(a) A perceptron with two inputs, each with weight -2 and a bias of 3 , can simulate a NAND gate. (b) A sigmoid function.

by the *weights* and a real number, the *threshold*, as follows:

$$\text{output} = \begin{cases} 0 & \text{if } \sum_i w_i x_i \leq \text{threshold} \\ 1 & \text{if } \sum_i w_i x_i > \text{threshold}. \end{cases} \quad (13.2)$$

The *perceptron rule* can be expressed using a *bias*, b :

$$\text{output} = \begin{cases} 0 & \text{if } (\sum_i w_i \times x_i) + b \leq \text{threshold} \\ 1 & \text{if } (\sum_i w_i \times x_i) + b > \text{threshold}. \end{cases} \quad (13.3)$$

A perceptron can be used to implement a NAND gate; NAND are *universal gates*; therefore the perceptron is also universal, i.e., one can implement any Boolean function using only perceptrons. To convince ourselves, consider a perceptron with two inputs, each with weight -2 and a bias of 3 , as in Fig. 13.2(a). The input 00 produces output 1 as $0 \times (-2) + 0 \times (-2) + 3 = 3$ is positive. Similar calculations show that the inputs 01 and 10 produce output 0 . But the input 11 produces output 0 as $1 \times (-2) + 1 \times (-2) + 3 = -1$ is negative.

A *sigmoid neuron* (SN) is a perceptron with the weights vector $w = (w_1, w_2, \dots)$, the input vector $x = (x_1, x_2, \dots)$, and a bias b . The output of an SN is $\sigma = (w \times x + b)$, rather than 0 or 1 , where σ is a nonlinear activation function, the sigmoid function depicted in Fig. 13.2(b), and defined as:

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad \text{or} \quad \sigma(z) = \frac{1}{1 + \exp(-\sum_i w_i \times x_i - b)}. \quad (13.4)$$

There are several types of neural networks, including:

1. *Multi-Layer Perceptrons (MLP)*. Each new layer is a set of nonlinear functions of the weighted sum of all outputs (fully connected) from a prior one; the weighted sums reuse the weights of the previous layer.
2. *Convolutional Neural Networks (CNN)*. Each layer is a set of nonlinear functions of weighted sums of spatially nearby subsets of outputs from the prior layer; the weighted sums reuse the weights of the previous layer.
3. *Recurrent Neural Networks (RNN)*. Each subsequent layer is a collection of nonlinear functions of weighted sums of outputs and the previous state.
4. *Long Short-Term Memory (LSTM)*. The decision what to forget and what to pass on as state to the next layer is critical. Weights are reused across time steps. Most popular RNN.

Convolutional Neural Networks are widely used multilayer networks where each layer raises the level of abstraction. For example, in computer vision, the first CNN layer recognizes horizontal and vertical lines, the second layer recognizes corners, the third layer recognizes shapes, a fourth layer recognizes features, such as tree leaves, and higher layers recognize multiple types of trees. A class project to study ML scalability is summarized in Section A.7.

Deep Neural Networks (DNNs) have multiple neuron layers and nonlinear activations functions, such as sigmoids. For example, Fig. 13.3 shows a DNN for digit recognition using a training dataset of $28 \times 28 = 784$ pixel images. There are 784 neurons in the input layer, one for each image in the training set [203].

DNN development consist of two phases: (a) definition of architecture, including: number of DNN layers, number of neurons in each layer, topology, and size of training data; (b) training and learning phase in which a back-propagation algorithm is used to determine the weights for each edge of the DNN. Once developed, a DNN enters a production phase involving testing, scoring, and running.

Google realized early on the potential impact of DNN applications and how computationally expensive DNN training is. For example, the training set size for activity recognition for an 8-layer DNN requires 10^6 videos, and the training time using 10 NVIDIA GPUs was reported to take some 30 days [232]. Training a CNN, internally labeled CNN1, with 89 layers, 100×10^6 weights, and 1750 operations/weight requires 17.5×10^9 operations. Another CNN, CNN0 with 16 layers, 8×10^6 weights requires 2888 operations. In 2016, Google decided to develop tensor processing units (TPUs) and attach them to some of its instances. The design goal for TPUs was to achieve at least one order of magnitude improvement versus GPUs as discussed in Section 3.6.

Clouds ML services. Amazon, Google, IBM, Microsoft, and virtually every other CSP offer a palette of ML services. There are countless applications of ML in areas as diverse as medicine, resource management, network optimizations, robotics, manufacturing, and marketing, among others.

AWS offers an impressive number of ML-related services. *SageMaker* is a fully managed service that facilitates building, training, and deploying machine learning (ML) models. *SageMaker Ground Truth* helps improve the quality of labels through annotation consolidation and audit workflows. *Kendra* is an enterprise search service powered by machine learning. *Forecast* uses machine learning to automatically discover how time series data and variables, such as product features and store locations, affect each other. *Augmented AI* makes it easy to build the workflows required for human review of ML predictions. *Polly* turns text into lifelike speech for applications that build entirely new types of speech-enabled products. *Elastic Inference* supports TensorFlow, Apache MXNet, PyTorch, and ONNX

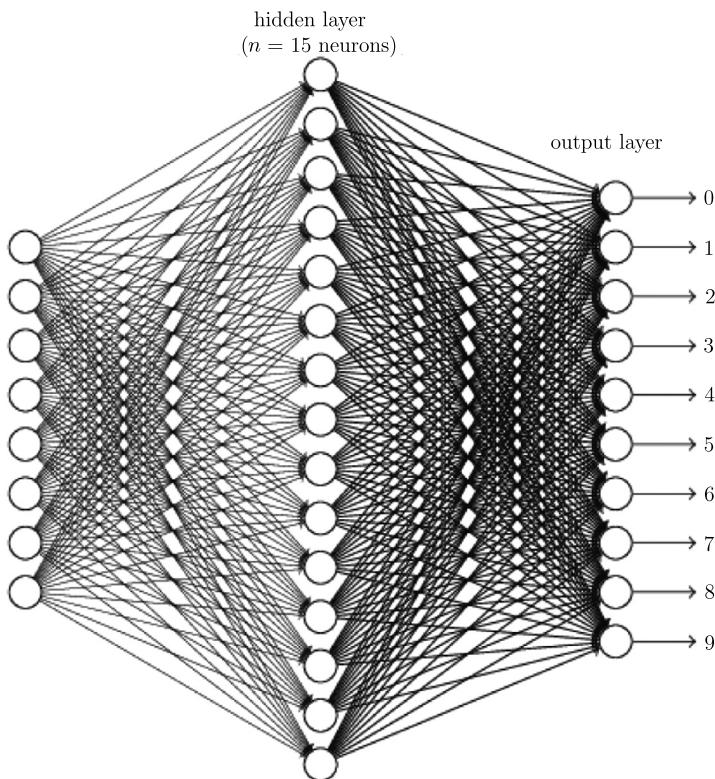


FIGURE 13.3

A DNN used to recognize digits with three layers: an input layer with 784 neurons, a hidden layer with 15 neurons, and an output layer with 10 neurons.

models; it attaches low-cost GPU-powered acceleration to Amazon EC2 and Sagemaker instances or Amazon ECS tasks and reduces the cost of running deep-learning inference by up to 75%.

AWS Deep Learning AMIs provide the infrastructure and the tools to accelerate deep learning on the cloud, at any scale. *Lex* can be used for building conversational interfaces into any application using voice and text; it provides the advanced deep-learning functionalities of automatic speech recognition (ASR) for converting speech to text, and natural language understanding (NLU) to recognize the intent of the text. *Rekognition* makes it easy to add image and video analysis to applications using deep-learning technology. Applications using this service can identify objects, people, text, scenes, and activities in images and videos, as well as detect any inappropriate content. *Comprehend* is a natural language processing (NLP) service using machine learning to discover insights and relationships in text. *Amazon Comprehend Medical* can be used to extract complex medical information from unstructured text.

The impact of AI and ML on computer architecture and implicitly on cloud computing is quite remarkable. David Patterson, the recipient of the 2017 Turing Award, shared his thoughts on Domain Specific Architectures (DSA) and the evolution of Tensor Processing Units (TPUs) at Google: “In 2013

Google calculated that if 100 million users started doing DNN three minutes per day on CPUs they would need to double the size of their data centers...Within 15 months they went from ideas to working hardware and software. The TPUv1 that Google designed had around a 80X performance per Watt of the 2015 Intel CPU and a 30X performance per Watt of the NVIDIA CPU because they were using 8-bit integer data rather than 32-bit floating point data and they dropped general purpose CPU/GPU features, which saves area and energy. Next Google created the TPUv2 that was designed to do ML training, which requires more computation, more memory, and bigger data. Google decided to build into the TPUv2 chips four Inter-Core Interconnect (ICI) links that each runs at 500 gigabits per second. Thus the links are approximately five times faster than those in a classic data center network at only one tenth of the costs. Eventually they created TPUv3 which further improved the system performance." Half a century of wisdom in designing general-purpose processors is now competing with an idea that would have been regarded an architectural blasphemy a decade ago.

13.3 Quantum computing on clouds

Most of the time, we are oblivious to the world of atomic and subatomic particles with strange, yet remarkable properties governed by the laws of quantum mechanics. The word "quantum" seems to always promise something extraordinary. Indeed, the quantum world tempts with the immense computing power generated by reversible circuits consuming little or no energy and with secure communication where an intrusion can be detected with very high probability.

Quantum theory and quantum mechanics. All started in 1900 when Max Plank, professor at the Friedrich-Wilhelms-Universität in Berlin, postulated that electromagnetic energy could be emitted only in quantized form $E = h \times \nu$, with h Planck's constant and ν the frequency of the radiation. The discovery of energy quanta won him the 1918 Nobel Prize in Physics.

Albert Einstein's elucidation of the photoelectric effect in 1905 confirmed Plank's hypothesis, and *quantum theory* was universally accepted. Einstein was rewarded with the 1921 Nobel Prize in Physics. Louis de Broglie developed the theory of electron waves in 1924. The *wave-particle duality principle* of de Broglie was soon followed by *quantum mechanics*, which describes the transformation and evolution of quantum systems.

Werner Heisenberg, Max Born, and Pasqual Jordan proposed the matrix formulation of quantum mechanics in 1925, and Erwin Schrödinger introduced the wave function that carries his name later the same year. In 1927, Heisenberg formulated the *Uncertainty Principle* asserting a fundamental limit to the precision the values of certain pairs of physical quantities characterizing a particle, such as position, x , and momentum, p , that can be predicted from initial conditions. Louis de Broglie, Werner Heisenberg, Erwin Schrödinger, and Max Born received the Nobel Prize in Physics in 1929, 1932, 1933, and 1954, respectively.

Quantum Information Processing. In 1982, Richard Phillips Feynman, who received the 1965 Nobel Prize in Physics for contributions to quantum electrodynamics, envisioned the development of a quantum computer and argued that only a quantum computer could *exactly simulate a quantum system*. In spite of the tremendous progress in quantum information processing during the past three decades, access to quantum computers is limited and simulation of quantum physical processes can only be

done using powerful classical computers, for example simulation of quantum many-body systems on the Amazon cloud [414].

The November 20, 2020, issue of *Physics Today*, a publication of the American Institute of Physics, includes the article *Quantum computing ramps up in the private sector*, which reviews the state of the art of quantum computing: “...hurdles remain to achieving useful quantum computers. The number of qubits needs to be scaled up. The qubits are needed not only for computations but also for correcting errors due to decoherence of the fragile quantum state. Engineering infrastructure must be designed and built. Algorithm must be created.” The article mentions the shift towards commercialization of Quantum Information Processing (QIP),² likely to create vast opportunities for young researchers in this field.

Quantum Computing. To compute in a digital world, we transform the state of solid-state devices capable to perform Boolean operations as directed by a program stored in the memory of the system. This process implies the ability to store information as the initial state of a physical system governed by the laws of classical physics and then to control the transformation of this information, i.e., the change of the system state, until we get the desired result in the shortest possible time and dissipating the least amount of energy.

If information can be stored as the state of a quantum system, e.g., the spin of electrons or the polarization of photons, then we can exploit the strange and often counterintuitive properties of quantum systems to process and communicate information faster and with little or no energy dissipation. There is a rich population of quantum particles, and several quantum particles and processes are being investigated for quantum computing; the list includes quantum dots³ based on spin or charge, trapped ions in a cavity, superconducting qubits, and liquid NMR (Nuclear Magnetic Resonance). At this time, ion traps and superconducting qubits are widely considered the leading candidates for quantum computers. Quantum communication only deals with information carried by photon polarization or other photon modes.

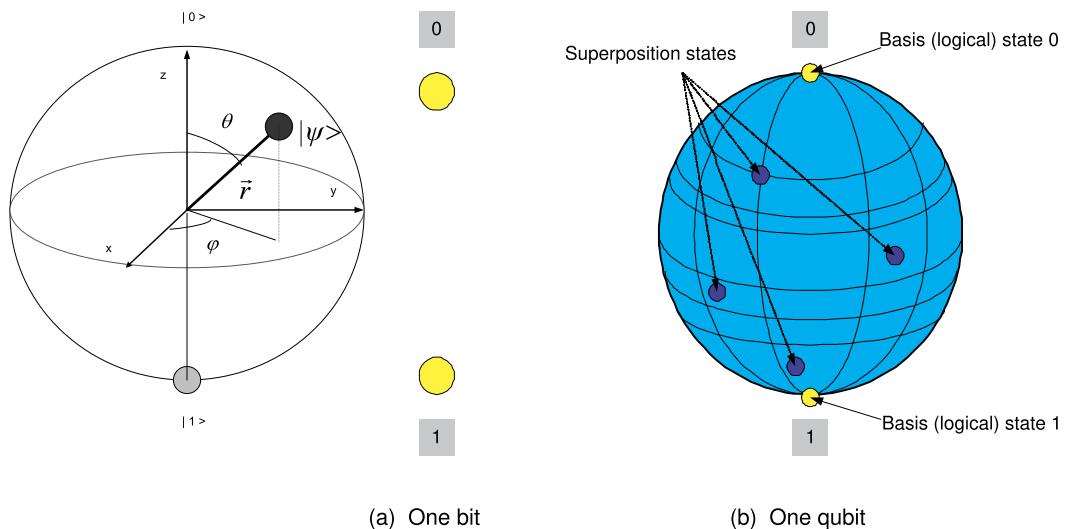
The physical processes in each embodiment of information as the state of a quantum system are investigated by various fields of physics, and for this reason we only discuss a unifying framework, the mathematical model of a quantum system used by quantum mechanics. In quantum mechanics, all state transformations happen in a finite-dimensional Hilbert space, a vector space over complex numbers. A set of 2^n complex numbers $\{\alpha_0, \alpha_1, \dots, \alpha_{2^n-1}\}$ with the property that the sum of their moduli is 1, i.e., $|\alpha_0|^2 + |\alpha_1|^2 + \dots + |\alpha_{2^n-1}|^2 = 1$, describes the state of a quantum system in an n -dimensional Hilbert space.

Bits, qubits, and the quantum circuit model. Quantum magic starts with the *qubit*, the quantum correspondent of a bit of classical information. A qubit embodies the state of the simplest quantum system. In the ket-bra notation introduced by Dirac, the state of a qubit is a vector in a two-dimensional Hilbert space $|\varphi\rangle = \alpha_0 |0\rangle + \alpha_1 |1\rangle$.

This possibly intimidating expression simply reflects the fact that a qubit is in a *superposition state*; $|0\rangle$ and $|1\rangle$ are the two basis vectors, and α_0 and α_1 are complex numbers and $|\alpha_0|^2 + |\alpha_1|^2 = 1$. The projections of the qubit state, $|\varphi\rangle$, on the two basis vectors $|0\rangle$ and $|1\rangle$ are $|\alpha_0|^2$ and $|\alpha_1|^2$, respectively.

² QIP refers to quantum computing and quantum communication.

³ Quantum dots are man-made nanoscale crystals that can transport electrons.

**FIGURE 13.4**

(Left) A qubit $|\psi\rangle$ is represented as a vector r from the origin to a point with angular coordinates $(\theta, \varphi, \gamma)$ on the Bloch sphere. (Right) A bit can only be in one of two states, 0 or 1, while a qubit can be anywhere on the surface of the Bloch sphere, a sphere of radius 1, i.e., can be in one of the two basis states $|0\rangle$ or $|1\rangle$ or in a superposition state $|\psi\rangle = \alpha_0|0\rangle + \alpha_1|1\rangle$ with $|\alpha_0|^2 + |\alpha_1|^2 = 1$.

To fully appreciate the impressive power of this model, consider the image of the qubit on a sphere of radius one called the *Bloch sphere*; see Fig. 13.4. A qubit can be anywhere on the surface of the Bloch sphere so a qubit can represent infinitely many bits (no pun intended) of information for different combinations of α_0 and α_1 [326].

When we measure a qubit, we project its state on the two basis vectors $|0\rangle$ and $|1\rangle$, and we obtain classical information: a bit with value 0 with probability $|\alpha_0|^2$ and one with value 1 with probability $|\alpha_1|^2$. In other words, when we observe or measure a qubit, we transform quantum information into classical information and force the qubit to either the south pole of the Bloch sphere corresponding to a classical bit with value 1 or to its north pole corresponding to a classical bit with value 0.

This is pure heresy: “God does not play dice” is the famous pronouncement of Albert Einstein reflecting his view about nondeterministic models of physical reality!! Fortunately, all experiments carried out during the past 100 years or so are consistent with the prediction of quantum mechanics, so it seems that Einstein and others questioning the nondeterminism of quantum mechanics are wrong.

Quantum parallelism and quantum algorithms. So far, we have looked at a single qubit, but imagine the immensity of a 2^n -dimensional Hilbert space and its potential to host quantum state transformations for a quantum computation. The implication of using qubits to process information brings us to the concept of *quantum parallelism*. Given a function $f(x)$, $x \in \mathcal{I}$, assume there are $N = p^q$ elements in \mathcal{I} , the domain of the function. A classical computer can compute a particular value of the function one at a time, i.e., $f(a)$ for $x = a$. A quantum system computing $f(a)$ will result in a superposition state including all values in the codomain of function f , i.e., $f(x) \forall x \in \mathcal{I}$.

Regardless of the cardinality N of the function domain, all values of the function are computed at once, while a classical computer would need a time T to compute one value; a parallel computer would need N functional units to do so in T/N time. This sounds wonderful, but how could we find the value we are interested in, $f(a)$, from this superposition? A process called *amplitude amplifications* enables a particular waveform of the superposition to stand out.

Another important observation is that quantum parallelism is most useful for computations requiring the values of a function for its entire domain, as is the case of algorithms using the Fast Fourier Transform (FFT). The quantum factoring algorithm published by Peter Shor in 1994 cleverly exploits quantum parallelism using QFFT, the quantum version of FFT. When a quantum computer with a large number of qubits will be available, this algorithm will allow us to decrypt any messages encrypted since April 21, 753 AD,⁴ and possibly earlier.

Another class of quantum algorithms was proposed by Lov Grover in 1996. Grover's algorithms find with high probability the unique input to a function that produces a particular output value, using just $\mathcal{O}(\sqrt{N})$ function evaluations, where N is the size of function's domain. Algorithms for quantum simulation of both classical and quantum systems are yet another important class of applications of quantum computers.

Quantum gates and quantum circuits. Quantum gates are the building blocks of a quantum computer based on the quantum circuit model. A gate carries out a unitary transformation U of its input qubits in state $|I\rangle$ into a set of qubits in state $|O\rangle$, such that $|O\rangle = U|I\rangle$. Single qubit gates, such as Pauli X , Y , Z gates and the Hadamard gate, rotate one qubit on the Bloch sphere. The target qubit of a CNOT gate is flipped when the control qubit is in state $|1\rangle$ and left unchanged otherwise. A Toffoli gate is a three-qubit gate with two control qubits and one target qubit; the target qubit is flipped when both control qubits are in state $|1\rangle$. Figs. 13.5(a) and (b) show one-qubit full adder circuit and reversible circuit with CNOT and Toffoli gates, respectively.

A set G of quantum gates is *universal* if, for any $\epsilon > 0$ and any unitary transformation U on n qubits, there is a sequence of gates $g_1, g_2, \dots \in G$ such that the following norm satisfies the condition:

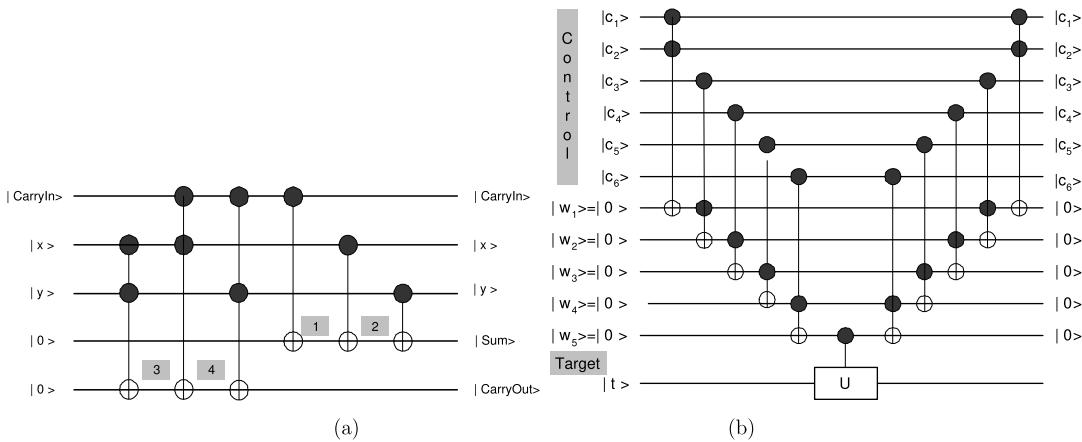
$$\|U - U_{g_1}U_{g_2}\dots U_{g_k}U_{g_1}\| \leq \epsilon. \quad (13.5)$$

Here, $U_g = V \otimes I$, where V is the unitary transformation on k qubits operated on by the quantum gate g , I is the identity operator acting on the remaining $n - k$ qubits, and \otimes means the tensor product. Unitary transformations carried out by quantum gates are reversible, and Eq. (13.5) illustrates that a sequence of transformations carried out to compute the result should then be executed in reverse order to bring return the system to its original state. Examples of universal gate sets include: (a) CNOT and all single qubit gates; (b) CNOT, Hadamard, and phase flip gates; and (c) Toffoli and Hadamard gates.

Quantum computers. In mid-2020 there were several quantum processors with a small number of qubits, the most powerful being the 53 qubit IBM Q quantum computer operational since October 2019.

What makes building a quantum computer so challenging, why does it take so much effort to have a useful quantum computer when there are several embodiments of quantum information? In a seminal paper on the physical implementation of quantum computers, David DiVincenzo, at IBM at that time,

⁴ According to a myth, on April 21, 753 B.C., Romulus and his twin brother, Remus, founded Rome.

**FIGURE 13.5**

(a) One-qubit full adder circuit with 3 Toffoli gates and 3 CNOT gates. (b) Quantum circuit with 10 Toffoli gates, 6 control qubits, $|c_1\rangle$, $|c_2\rangle$, $|c_3\rangle$, $|c_4\rangle$, $|c_5\rangle$, and $|c_6\rangle$, one target qubit $|t\rangle$, and 5 work qubits $|w_1\rangle$, $|w_2\rangle$, $|w_3\rangle$, $|w_4\rangle$, and $|w_5\rangle$ initially in state $|0\rangle$. The first stage produces the logical product, AND, of all six control qubits in $|w_5\rangle$. The second stage performs a single qubit U-controlled transformation of the target qubit. The last stage returns the work qubits to their original state, $|0\rangle$.

formulated five criteria for building a useful quantum computer, plus two additional ones required for transfer of information. These criteria are:

1. A scalable physical system with well-characterized qubits.
2. The ability to initialize the qubits to a pure state $|000\dots0\rangle$.
3. Long decoherence times, much longer than the gate operation times.
4. Existence of a “universal” set of quantum gates.
5. A qubit-specific measurement capability.
6. Ability to interconvert “stationary” (memory) and “flying” (communication) qubits.
7. Ability to faithfully transmit “flying” qubits between specified locations.

The last two requirements are related to communication. There are obvious similarities between these requirements and the ones for a classical computing engine. A recent paper [307] compares two architectures, superconducting transmon qubits⁵ and trapped ions. The two fully programmable multi-qubit machines provide the user with the flexibility to implement arbitrary quantum circuits with high-level interface. This makes it possible, for the first time, to test quantum computers irrespective of their particular physical implementation.

A research group at Google claimed quantum supremacy according to an October 2019 issue of the journal *Nature*. *Quantum supremacy*, a long awaited milestone for quantum computing, is achieved

⁵ A transmon is a superconducting charge qubit designed to have reduced sensitivity to charge noise. It was described in 2007 by Robert J. Schoelkopf, Michel Devoret, Steven M. Girvin, and their colleagues at Yale University.

when a quantum system solves a computational problem that is beyond the practical capabilities of “classical” computers. Google’s claim that the problem solved using the Sycamore chip would take 10 000 years to complete on the fastest classical supercomputer was contested by IBM researchers who argued that the problem solved by Sycamore could in fact be solved in 2.5 days on a classical supercomputer.

Decoherence and quantum error correction. Quantum systems are affected by *decoherence*, the loss of information from the quantum system into the environment the system is loosely coupled with. The decoherence time, the time elapsed since the state is set until the decoherence alters it, is very short. For example, the decoherence time for quantum dots-based charge and quantum dots-based spin are 10^{-9} and 10^{-6} seconds, respectively; the decoherence time for trapped indium atoms in quantum dots is 10^{-1} seconds [329]. This implies that error correction should be done periodically at intervals shorter than the decoherence time. “The speed at which decoherence occurs can make or break a qubit” summarizes the thinking of QIT researchers. Quantum error correction faces multiple challenges:

1. Only orthogonal quantum states can be distinguished with absolute certainty; if a quantum state is transformed by two distinct quantum errors into nonorthogonal states, then neither error can be corrected because the two states are not distinguishable.
2. A quantum measurement projects the state, thus, in the general case, it alters the state of a quantum system; we have to devise ingenious means to carry out the measurements necessary to identify a qubit in error.
3. A qubit may be affected by bit-flip, phase-flip, or by both bit- and phase-flip errors.
4. Entanglement adds a new dimension to quantum error correction; an error may propagate to some or all qubits entangled with one another, and we can expect non-Markovian quantum errors, errors correlated in time and/or space.

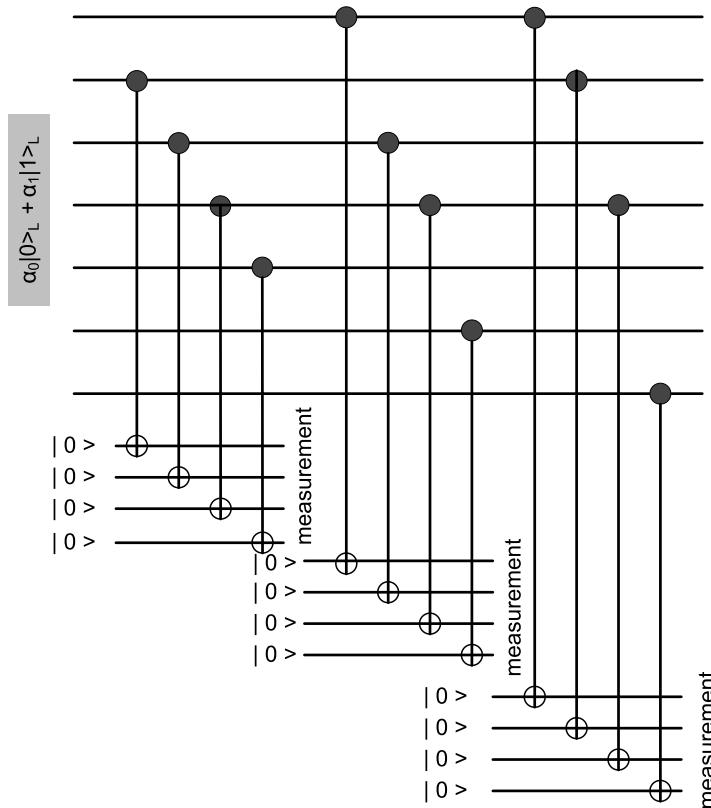
A major discovery is due to Peter Shor, who in 1995 showed that it is possible to restore the state of a quantum system using only partial information about the state. Until then, it was not known how to perform an error correction step without a measurement that would alter the state. Shor’s idea was *to fight entanglement with entanglement*.

The intuition that a classical binary $[n, k]$ linear code, with n the number of bits of every codeword and $k < n$ the number of information bits, is able to correct t errors iff (if and only if) Hamming spheres of radius t about each codeword are disjoint and exhaust the entire space of n -tuples extends to quantum codes. A quantum error correcting code encodes one logical qubit L to n physical qubits in a 2^n Hilbert space with basis vectors $\{|0\rangle, |1\rangle, \dots, |i\rangle, \dots, |2^n - 1\rangle\}$ with

$$|0\rangle_L = \sum_{i=0}^{2^n-1} \alpha_i |i\rangle \quad \text{and} \quad |1\rangle_L = \sum_{i=0}^{2^n-1} \beta_i |i\rangle$$

The 2^n -dimensional Hilbert space should accommodate orthogonal subspaces for possible bit-flip, phase-flip, and bit- and phase-flip errors affecting each one of the n physical qubits.

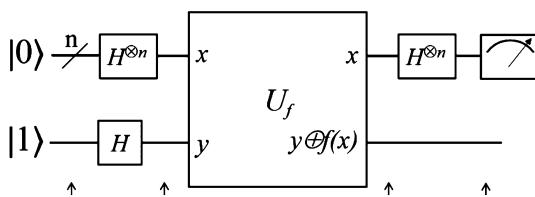
The quantum error correcting code must map coherently the two-dimensional Hilbert space spanned by $|0\rangle_L$ and $|1\rangle_L$ to multi-dimensional Hilbert spaces to ensure that the code is capable of correcting the three types of errors for each of the n qubits of the basis vectors $|0\rangle_L$ and $|1\rangle_L$, and, in addition, one subspace for entangled state $|0\rangle_L$ and $|1\rangle_L$.

**FIGURE 13.6**

A circuit for calculating the syndrome for bit-flip errors for the 7-qubit Steane code. The input is the logical qubit $|\psi\rangle = \alpha_0|0\rangle_L + \alpha_1|1\rangle_L$. The seven physical qubits act as the control qubits of the 12 CNOT gates that set the ancilla qubits; a nondemolition measurement of the syndrome reveals if one of the seven physical qubits is in error.

The Quantum Hamming Bound, $2(3n + 1) \leq 2^n$, reflects these requirements; it shows that the smallest number of physical qubits used to encode the superposition states $|0\rangle_L$ and $|1\rangle_L$ and then recover them regardless of the qubit in error and the type of error is $n = 5$.

The *syndrome* Hv^T of a linear code with the parity check matrix H uniquely identifies the bit(s) in error of the n -tuple v (v^T is the transpose of vector v). We apply the same basic strategy to identify a qubit in error, and we carry out a *nondemolition measurement* to determine the syndrome. Computation of the syndrome for a quantum code requires a number of ancilla (auxiliary) qubits; each ancilla qubit is entangled with multiple qubits of the encoded quantum word—see, for example, Fig. 13.6. The degree of entanglement of an ancilla qubit with each one of the n physical qubits is relatively weak; the non-demolition measurement of ancilla qubits allows us to determine the syndrome without altering the state of the n physical qubits.

**FIGURE 13.7**

Quantum circuit for the Deutsch–Jozsa problem.

Implementation of quantum error-correcting codes significantly complicates the design of quantum processors. Quantum error correction increases by one or more orders of magnitude the number of qubits and the number of elementary steps required by a quantum computation. The overhead of the error correction is measured by: (i) *scale-up*—the ratio of the number of physical to logical qubits; and (ii) *slow-down*—the ratio of the number of gates to the number of the elementary steps of computation.

Quantum circuit simulation with Qiskit. Qiskit is an open-source simulation software for experimentation with quantum computers at the level of circuits, pulses, and algorithms; see https://qiskit.org/documentation/intro_tutorial1.html. Qiskit use involves several steps: (i) design the quantum circuit(s); (ii) compile the circuit(s); (iii) run the compiled circuits; and (iv) analyze, i.e., compute summary statistics and visualize the results of the experiments. We illustrate now the use of this wonderful tool for a simple but insightful problem named after quantum computing researchers who proposed the problem and its solution. The Deutsch–Jozsa problem was conceived to illustrate the power and elegance of quantum algorithms, rather than solve a practical problem.

Deutsch–Jozsa problem. Assume that Alice plays the following game with Bob: (i) Bob chooses an either constant or balanced function $f(x)$. Function $f(x)$ is constant if it has the same value for all arguments x , i.e., $f(x) = 0 \forall x$ or $f(x) = 1 \forall x$; $f(x)$ is balanced if $f(x) = 1$ for half of values x and $f(x) = 0$ for the other half; (ii) Alice sends an n bit integer N ($0 \leq N \leq 2^n - 1$) to Bob; and (iii) Bob calculates $f(N)$ and sends Alice the one bit of information 0 if f is constant and 1 if f is balanced.

The question is how fast can Alice succeed, i.e., how many times does she need to send an n -bit integer to Bob before she is certain whether Bob had chosen a constant or a balanced function $f(x)$? Classically, in the worst case, Alice has to send 2^{n-1} values and receive a 0 before receiving a 1 or vice versa. So, she needs to query Bob $2^{n-1} + 1$ times. If Alice is very lucky, she can guess after two trials; if she gets different values, she will know that $f(x)$ is balanced. Fortunately, she can do much better using the quantum circuit discussed next.

Fig. 13.7 shows the quantum circuit for the Deutsch–Jozsa problem. The input consists of an n -qubit query register to store the value N and a one qubit answer register for Bob to store $f(N)$. An n -dimensional Hadamard Transform $H^{\otimes n}$ is applied to the first register, and a Hadamard Transform is applied to the one qubit answer register. The operator $H^{\otimes n} \otimes H$ transforms the input state at the leftmost uparrow \uparrow , and the resulting state at the next \uparrow is the input to a unitary transformation U_f carried out by an *oracle*. Call x and y the two inputs of U_f and x and $y \oplus f(x)$ with \oplus the binary addition, at the output, at the third \uparrow in the figure. Then the top register is transformed again by $H^{\otimes n}$ and then measured.

The quantum solution exploits quantum parallelism. Alice measures the query register in state

$$|\varphi\rangle = \sum_z \sum_x \frac{(-1)^{x \cdot z + f(x)}}{\sqrt{2^n}}$$

with $x \cdot z$ being the bitwise inner product of x and z modulo 2. The amplitude of state $|0^{\otimes n}\rangle$ is $A = \frac{\sum_x (-1)^{f(x)}}{2^n}$. If f is constant, A is either $+1$ or -1 , depending on $f(x)$ and the measurement of the query register will be 0 because $|\varphi\rangle$ is of unit length and all other amplitudes must be 0. If f is balanced, the positive and negative contributions to A cancel, thus $A = 0$; a measurement must show a result different from 0 for at least one qubit of the query register. In conclusion, Alice gets the result 0 if and only if function f is balanced after only one evaluation of function f , regardless of how large n is.

Qiskit simulation is presented in Fig. 13.8. Fig. 13.8a displays the initialization code for a simulation with a 4-qubit query register. Fig. 13.8b displays the code for Hadamard Transformations, oracle invocation, and measurement; the code includes drawing the circuit for each stage, but omitted in this presentation. The oracle implements the unitary transformation U_f . Fig. 13.8c shows the code for running the simulation, the full circuit, and the state at each stage of the circuit.

Cloud access to quantum simulation software; experimenting with quantum systems. To spark interest in quantum computing among computer programmers, venture capitalists, and students wishing to explore this fascinating field, several companies support cloud access to platforms for experimenting with quantum computing. Multiple sources for information about quantum software are also available; Quantum Open Source Foundation maintains a GitHub open-source quantum software projects repository; see <https://github.com/qosf/awesome-quantum-software>. An overview of open source quantum computing software can be found in [174].

Several organizations offer access to a set of quantum programming tools and to quantum systems with a small number of qubits, while others are only involved in quantum programming. Quantum programming is the process of assembling sequences of instructions, called quantum programs, capable of running on a quantum computer. Quantum programming languages help express quantum algorithms using high-level constructs.

Arguably, the most sophisticated quantum environment is *IBM Q Experience*, an online platform launched in 2016. Q Experience allows general public access to a set of IBM's prototype quantum computers including two 5-qubit processors and one 16-qubit processor. A quantum system with 20 qubits is also available through the IBM Q Network; see <https://www.ibm.com/quantum-computing/>.

IBM's Q simulators and quantum devices can be accessed over the cloud using an open-source quantum programming framework launched in 2017. Qiskit integrates quantum development into regular workflows using common programming languages and standard development tools.

Quantum Computing Playground is a WebGL Chrome Experiment offered by Google. The experiment uses a GPU-accelerated quantum computer with an IDE interface and a scripting language with debugging and 3D quantum state visualization support. The system has several quantum gates built into the scripting language, can simulate quantum registers of up to 22 qubits, and runs Shor factorization and Grover database search algorithms.

Microsoft released recently the *Quantum Development Kit*, replacing an earlier version called LIQUi|. This system includes: (i) a quantum programming language, Q#, integrated with Visual Studio development environment; (ii) simulators that run on a local system or on the Azure cloud platform;

```

# importing Qiskit
from qiskit import *
%matplotlib inline

# Set number of bits of input
n = 4
# Create circuit object
qc = QuantumCircuit(n+1,n)

# Create the oracle
b = '0101'
oracle = QuantumCircuit(n+1)

# Invert where we set the b string to 1
for i in range(n):
    if b[i] == '1':
        oracle.x(i)

# CNOT with the target bit for each bit in the input
for i in range(n):
    oracle.cx(i,n)

# Revert any inverted bits from the b string
for i in range(n):
    if b[i] == '1':
        oracle.x(i)

# Set the working bit to 1 (initial state)
qc.x(n)
qc.barrier()
qc.draw(output='mpl')

```

```

# Apply Hadamard gates
for i in range(n+1):
    qc.h(i)
qc.barrier()
qc.draw(output='mpl')

# Apply the oracle
qc += oracle
qc.barrier()
qc.draw(output='mpl')

# Repeat hadarmard gate
for i in range(n):
    qc.h(i)
qc.barrier()
qc.draw(output='mpl')

# Measure
for i in range(n):
    qc.measure(i,i)
qc.draw(output='mpl')

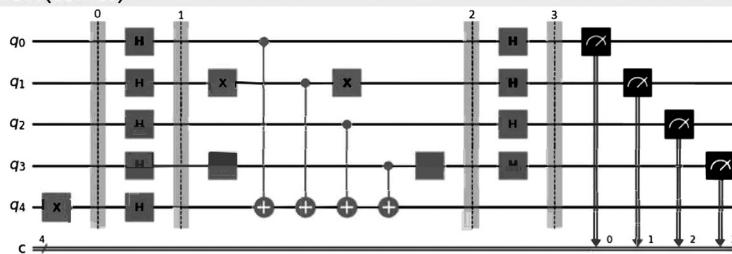
```

(a) Qiskit initialization; 4-qubit query register.

```

# Run a simulation
results = Aer.get_backend('qasm_simulator').run(assemble(qc, shots=10000)).result()
counts = results.get_counts()
print(counts)

```



Stage	Step	State
0	Input	00001
1	Hadamard	$\frac{1}{\sqrt{2^{n+1}}} \sum_{x=0}^{2^n-1} x\rangle (0\rangle - 1\rangle)$
2	Oracle	$\frac{1}{\sqrt{2^{n+1}}} \sum_{x=0}^{2^n-1} (-1)^{f(x)} x\rangle (0\rangle - 1\rangle)$
3	Hadamard	$\frac{1}{2^n} \sum_{y=0}^{2^n-1} \left[\sum_{x=0}^{2^n-1} (-1)^{f(x)} (-1)^{xy} \right] y\rangle$

(c) Run simulation, the circuit, and state.

FIGURE 13.8

Qiskit simulation: (a) and (b) the simulation code; (c) the circuit and the state.

and (iii) libraries and code samples. The software can be downloaded from <https://www.microsoft.com/en-us/quantum/development-kit>.

Quantum Inspire (<https://www.quantum-inspire.com/>) is the first European platform providing online access to a fully programmable 2-qubit electron spin quantum processor, *Spin-2*, and to a quantum processor with five superconducting qubits, *stramon-5*. *Spin-2* is a quantum processor hosting two single electron spin qubits in a double quantum dot in isotopically purified Si (Silicon). *Starmon-5* consists of five superconducting transmon qubits in an X configuration.

The QX emulator from Quantum Inspire simulates up to 26 qubits on a commodity cloud-based server and up to 31 qubits using “fat” nodes on the Cartesius supercomputer. Circuit-based quantum algorithms can be created using either a GUI or the Python-based Quantum Inspire SDK. A knowledge base is also available at <https://www.quantum-inspire.com/kbase/introduction-to-quantum-computing/>.

In 2018, QC Ware, a quantum-computing-as-a-service company from Palo Alto, launched *Forge*, a cloud platform allowing developers to run algorithms on hardware provided by multiple vendors. QC Ware is focused on software development rather than building its own quantum computers. For now, *Forge* offers access to a D-Wave quantum computer for simulations of Google and IBM quantum devices. D-Wave One, a quantum computer operating on a 128-qubit chipset became operational in 2011, D-Wave 2X, a 1000+ qubit quantum computer and D-Wave 2000Q and an open source repository containing software tools for quantum annealers became available in 2015 and 2017, respectively. D-Wave quantum computers use quantum annealing for solving optimization problems. Quantum annealing is a method for finding the global minimum of a function using quantum fluctuations.

Rigetti Computing is a Berkeley-based developer of quantum integrated circuits used for quantum computers. The company also develops a cloud platform called *Forest* that enables programmers to write quantum algorithms. The platform is based on a custom instruction language the company developed called Quantum Instruction Language (Quil) using a set of Python tools. The system allows users to write quantum algorithms for a simulation of a quantum chip with 36 qubits.

IonQ, a company founded in 2015, aims to take trapped-ion quantum computing out of the lab and into the market. In August 2020 IonQ opened an 11-qubit quantum computer to the public via the Amazon Braket cloud platform, and in October unveiled a 32-qubit version of it. Palo Alto-based PsiQuantum, also founded in 2015, patterns silicon wafers into thousands of photonic components containing waveguides for carrying the single photons encoding the qubits. PsiQuantum founders predict that they will have a useful quantum computer with a million qubits in a few years.

13.4 Vehicular clouds

By way of motivation. Present-day vehicles feature powerful on-board computers, ample storage, sophisticated transceivers, and an impressive panoply of sensing devices. In the near future, we expect to see millions of these advanced vehicles crisscrossing our roadways and city streets. We hold the view that it is only a matter of time before this fleet of sophisticated vehicles will be recognized as an abundant and underutilized compute resource that can be tapped into for the purpose of providing third-party or community services.

One of the characteristics that distinguishes vehicles from servers in a conventional cloud is mobility. Indeed, vehicles participating in traffic are involved, on a daily basis, in various dynamically

changing situations, ranging from normal traffic to congestion, to accidents, to other similar traffic-related events.

Under the present-day state of the practice, our vehicles are mere spectators that witness traffic-related events without being able to participate in the mitigation of their effects. We suggest that in such situations the vehicles have the potential to cooperate with various authorities to help solve problems that otherwise either take an inordinate amount of time to solve (e.g., dissipate traffic jams) or be solved in a timely manner for lack of adequate computing resources and accurate traffic information. In our vision, vehicles stuck in a traffic jam will be able to let municipal Traffic Management Centers use their on-board compute resources and sensors to run complex simulations designed to help alleviate the effects of congestion by city-wide rescheduling of traffic lights.

We also anticipate that, in planned evacuation situations, the combined compute power of the vehicles will enable evacuation management centers to compute and disseminate information about the safest evacuation routes.

Vehicular Clouds—How it all started. A decade ago, inspired by the promise of conventional clouds and by the realization that modern vehicles can act as servers in a vehicular data center, Abuelela and Olariu [7] and Eltoweissy et al. [163], introduced Vehicular Clouds (VCs) as: “*A group of vehicles whose corporate computing, sensing, communication and physical resources can be coordinated and dynamically allocated to authorized users*”.

Although the authors of [7] and [163] have contemplated from the start VCs set up on moving vehicles, it is interesting to note that the first papers discussing VCs considered the simpler case of parked vehicles. Thus, it has become customary in the VC literature to distinguish between *static* VCs involving stationary (e.g., parked) vehicles and *dynamic* VCs that are harnessing the compute resources of moving vehicles. We believe that this distinction is largely artificial and that, with the advent of ubiquitous network coverage available in the Smart Cities of the (near) future, connectivity will be ubiquitous and pervasive, and the on-board compute resources of vehicles will be available 24/7, whether these vehicles are parked or not.

Vehicular clouds are a nontrivial variant of conventional clouds. One of the key distinguishing characteristics of VCs is the dynamically changing amount of compute resources available. This characteristic is a direct consequence of mobility and of the distributed ownership of vehicles. As vehicles enter the VC, fresh compute resources become available; when vehicles leave, often unexpectedly, their resources depart with them, creating a highly dynamic environment. In turn, the dynamically changing availability of compute resources, due to vehicles joining and leaving the VC, unexpectedly leads to a volatile computing environment where reasoning about system performance becomes challenging [177,178].

The mobility attribute of VCs, combined with the fact that the presence of vehicles in close proximity to an event is very often unplanned, implies that the pooling of the resources of those vehicles into a VC in support of mitigating the event must occur spontaneously by the common recognition of a need for which there are no preassigned or dedicated resources available. This *agility of action* seldom exists in conventional clouds and turns out to be an important defining characteristic of VCs.

The early papers [7] and [163] hinted at an entire array of possible applications of VCs. In addition to numerous interesting applications of VCs, these early papers identified an entire series of research challenges. Given their large array of potential applications, it is reasonable to expect that, once adopted, the VCs will be the next paradigm shift, with a lasting technological and societal impact. In fact, the VCs may well turn out to be the “killer app” taking vehicular networks to the next level of relevance

and innovation [371]. Recent years have seen the emergence of VCs as an active topic of research. Recent surveys [69], [372], and [373] summarized trends and advances in VC research and pointed out new application domains and implementation challenges ahead.

VC resource management. Resource management is a fundamental task in conventional clouds and remains of great importance in VCs. Recall that one of the defining ways in which VCs differ from conventional clouds is resource volatility. Therefore, managing VC resources is an important, albeit very challenging, task. The main goal of this section is to discuss a few important instances of resource management in VCs.

Estimating the expected number of vehicles in the VC. One of the key tasks to be performed by the VC *resource manager* is to predict, with some degree of accuracy, the number of vehicles present in the VC, given generic information about the vehicle arrival and departure rates. Specifically, assume that at time $t = 0$, the VC contains $n_0 \geq 0$ vehicles. After that, vehicles arrive at a time-dependent rate $\lambda(t)$ and depart at a time-dependent rate $\mu(t)$. Of interest are:

- the probability that, at a given time $t > 0$, there are exactly k vehicles in the VC;
- the variance of the number vehicles in the VC at time $t > 0$;
- the limit as $t \rightarrow \infty$ of the expected number of vehicles in the VC (provided such a limit exists).

Consider the counting process $\{N(t) \mid t \geq 0\}$ of continuous parameter t , where for some arbitrary integer k the event $\{N(t) = k\}$ occurs if the VC contains k vehicles at time t . We let $P_k(t) = \Pr[\{N(t) = k\}]$ denote the probability that the event $\{N(t) = k\}$ occurs and assume that the VC has infinite capacity. We model our problem as a nonhomogeneous birth-and-death process.

Letting $G(z, t) = \sum_k P_k(t)z^k$ denote the probability generating function of $P_k(t)$, it is fairly routine to see that $G(z, t)$ satisfies the partial differential equation

$$\frac{\partial G(z, t)}{\partial t} + \mu(t)(z - 1)\frac{\partial G(z, t)}{\partial z} = \lambda(t)(z - 1)G(z, t), \quad (13.6)$$

with the boundary condition $G(z, 0) = z^{n_0}$. It is also straightforward to confirm that the solution to (13.6) is

$$\begin{aligned} G(z, t) &= \left[e^{-\int_0^t \mu(u)du} z + (1 - e^{-\int_0^t \mu(u)du}) \right]^{n_0} \\ &\times \exp \left[-(1 - z)e^{-\int_0^t \mu(u)du} \int_0^t \lambda(u)e^{\int_0^u \mu(s)ds} du \right]. \end{aligned} \quad (13.7)$$

In spite of its complexity, Eq. (13.7) reveals much of the structure of the process $\{N(t) \mid t \geq 0\}$. Indeed, observe that $G(z, t)$ is the product of two factors:

- $\left[e^{-\int_0^t \mu(u)du} z + (1 - e^{-\int_0^t \mu(u)du}) \right]^{n_0}$ that is the probability generating function of a binomial random variable with parameter n_0 and “success” probability $p(t) = e^{-\int_0^t \mu(u)du}$;

- $\exp \left[-(1-z) e^{-\int_0^t \mu(u) du} \int_0^t \lambda(u) e^{\int_0^u \mu(s) ds} du \right]$ that, upon examination, turns out to be the probability generating function of a (nonhomogeneous) Poisson random variable with parameter

$$\Lambda(t) = \frac{\int_0^t \lambda(u) e^{\int_0^u \mu(s) ds} du}{e^{\int_0^t \mu(u) du}}. \quad (13.8)$$

We define two additional counting processes:

- $\{R(t) \mid t \geq 0\}$ that keeps track of the number of the n_0 vehicles in the VC at time $t = 0$ that are still in the VC at time t ; it is clear that the “success” probability $p(t) = e^{-\int_0^t \mu(u) du}$ is precisely the probability that such a generic vehicle is still in the VC at time t ;
- $\{S(t) \mid t \geq 0\}$ that keeps track of the number of vehicles in the VC at time t that were not in the VC at time $t = 0$; as already mentioned, this is a nonhomogeneous Poisson process with parameter $\Lambda(t)$ defined in (13.8).

It is immediately clear that, for all $t \geq 0$, $R(t)$ and $S(t)$ are independent random variables. Furthermore, the expression of $G(z, t)$ as a product implies that for all $t \geq 0$, $N(t)$ is the convolution of $R(t)$ and $S(t)$, and so

$$N(t) = R(t) + S(t). \quad (13.9)$$

A direct consequence of (13.9) and of the linearity of expectation is that

$$\begin{aligned} E[N(t)] &= E[R(t)] + E[S(t)] \\ &= n_0 e^{-\int_0^t \mu(u) du} + e^{-\int_0^t \mu(u) du} \int_0^t \lambda(u) e^{\int_0^u \mu(s) ds} du \\ &= e^{-\int_0^t \mu(u) du} \left[n_0 + \int_0^t \lambda(u) e^{\int_0^u \mu(s) ds} du \right]. \end{aligned} \quad (13.10)$$

Even though the standard way of obtaining $E[N(t)]$ would be to compute the partial derivative of $G(z, t)$ with respect to z and then take $z = 1$, the merit of Eq. (13.10) is that it provides a more direct way of computing the expected number of vehicles in the VC. The same observation holds, *mutatis mutandis*, for computing the variance $Var[N(t)]$ of $N(t)$.

Virtual machine migration in VC. We assume that each vehicle has a *virtualizable* on-board computer preloaded with a Virtual Machine Monitor in charge of mapping between the guest VM and the host vehicle’s resources. VM migration is one instance of *resource management* in cloud computing. It remains a fundamental performance attribute in VCs, as well.

One of the fundamental resources in the VC are the various VMs hosted by vehicles. When a vehicle leaves the VC while its hosted VM is still active, a failure occurs and various strategies need to be implemented to minimize disruption and performance degradation. It is intuitively clear that, the longer and more predictable the vehicle residency times are in the VC, the easier it is to predict the optimal moment when VM migration has to be undertaken. By contrast, when vehicular residency times in the VC are short and/or unpredictable, VM migration becomes very challenging.

Baron et al. [49] used mobility traces of city buses in Dublin, Ireland, to investigate the feasibility of VM migration between buses at communication hotspots. Refaat et al. [409] have proposed a Vehicular

Virtual Machine Migration (VVMM) framework. They have studied strategies for selecting the target vehicle for VM migration. Unlike the VM migration strategy of [49], VVMM requires preinstalled roadside infrastructure. Finally, Florin et al. [178] have investigated VM migration in the context of VCs running compute-intensive applications in support of Big Data.

Estimating vehicular residency times. To identify the types of workloads supported by a given VC architecture, it is exceedingly important to model vehicular residency times as a function of its specific mobility parameters. Although few authors seem to pay much attention, it is of fundamental importance to distinguish between the characteristics of vehicle residency times in the VC and the VM lifetime. VMs are associated with job execution, while vehicle residency times reflect the amount of time vehicles spend in the VC.

Recently, three papers investigated *quantitatively* vehicular residency times. The first, [31], was concerned with estimating the number of vehicles present in a VC given stochastically changing arrival and departure parameters. The second, [540], assumed a VC built on top of moving vehicles and have studied the amount of residency times of vehicles in the VC assuming various characteristics of (urban) traffic flow. Finally, [524] also looked at the residency times of vehicles in the downtown area of a smart city. Their result can be extended with minor changes to assessing the vehicular residency times in a VC set up in the downtown area of a smart city, assuming that a sufficient level of supporting roadside infrastructure is available.

Identifying feasible VC workloads and applications. The VC literature is full of papers discussing VC architectures implemented on top of moving vehicles and supported by ubiquitous and pervasive roadside infrastructure. However, some of those authors do not seem to be concerned with the obvious fact that moving vehicles' residency times in the VC may, indeed, be very short and therefore so is their contribution to the amount of useful work performed. While such supporting infrastructure may well become available in the future, it does not exist today. Similarly, zero-infrastructure VCs built on top of fast-moving vehicle are utterly unrealistic under present-day technology since the DSRC-mandated V2V transmission range, combined with the limited bandwidth inherent to wireless communications, cannot guarantee sufficient time to permit VM migration and/or data offloading in case vehicles exit the VC.

Most conventional clouds are supported by tens of thousands of servers. By comparison, VCs involve a much smaller number of vehicles. As a result, while conventional clouds can handle arbitrary workloads and provide high reliability of the order of “four nines” (i.e., 99.99%) availability, the workloads supported by VCs depend, to a large extent, on the number and residency characteristics of participating vehicles.

Thus, various VC architectures, distinguished by the number of participating vehicles, existing communication bandwidth, and vehicle residency times, are, in all likelihood, apt to support multiple types of user applications and workloads.

Given the significant differences between various instances of VCs, one of the important tasks is to identify *feasible* applications for each of them. For example, VCs with a very short vehicle residency times are clearly not suited for long-lived applications involving a lot of data. This is so, among other reasons, because VM migration and data replication will be next to impossible to support efficiently.

The workloads specific to these application categories differ substantially and require different architectural support. Can VCs support these applications? Very little work along these lines has been reported in the literature. It is clear that a VC, built on top of a large number of participating vehicle known to spend a predictable amount of time in the VC, is suitable for demanding applications.

By contrast, VCs set up in the parking lot of a convenience store where vehicles spend a short amount of time will not be able to support effectively the same application. They, however, can support different, more lightweight applications.

We suggest using VCs in support of the following types of workloads: (i) transaction-oriented processing for decision support systems and business analytics; and (ii) mobile interactive applications that process large volumes of data from different types of sensors. As an illustration, imagine a VC set up in the parking lot of a shopping mall. During the normal business hours, the various businesses at the mall can leverage the compute power of the vehicles in the shopping-mall parking lot for business analytics, mostly involving transaction-oriented processing. Likewise, the mall security personnel can harness the compute power of the cars in the parking lot to process various sensor inputs, including face recognition and other similar software programs intended to provide security for the shoppers at the mall.

Assigning job to vehicles in VC. As mentioned before, the dynamically changing availability of compute resources due to vehicles joining and leaving the VC leads to a volatile computing environment where the task of assigning incoming user jobs to vehicles is quite challenging.

To get a feel for the problem, assume that a user job was assigned to a vehicle in the VC. If the vehicle remains in the VC until the job completes, all is well. Difficulties arise when the vehicle leaves the VC before job completion. In this case, unless special measures are taken, the work is lost, and the job must be restarted, taking chances on another vehicle, and so on, until eventually the job is completed. Clearly, losing the entire work of a vehicle, in case of a premature departure, must be mitigated. One possible method is to use checkpointing, a strategy originally proposed in databases. This strategy requires making regular backups of the job's state, and storing them in a separate location. In the event a vehicle leaves before job completion, the job can be restarted, on a new vehicle, using the last backup. While this seems simple and intuitive, checkpointing is not implemented efficiently in VCs.

Reliability and availability can be enhanced, as it is often done in other computer systems, by employing various forms of redundant job-assignment strategies. Recently, Ghazizadeh et al. [196] have studied redundancy-based job assignment strategies and have derived analytical expressions for the corresponding MTTF. Building on the results in [196], Florin et al. [177] have developed an analytical model for estimating job completion time in VCs assuming a redundant job-assignment strategy.

Research Challenges. Several important VC research challenges are discussed next.

Architectural Challenges. The architectural challenges faced by the VC community include issues related to the organization of the cyber structure of the VC and its interaction with the physical resources. There is a critical need to efficiently manage host mobility and heterogeneity (including compute, communication, and storage capabilities) and vehicle membership.

One of the obvious applications of VCs is to help the authorities mitigate the effects of traffic-related events whose timely resolution cannot be met by preassigned assets or in a proactive fashion. It is, therefore, clear that the VCs can only achieve their full potential if their basic architecture is engineered to offer seamless integration of the resources of the participating vehicles. In particular, the architecture must adapt its managed vehicular resources allocated to an application according to dynamically changing requirements and system conditions.

Operational and Policy Challenges. For the VC to operate seamlessly, issues related to authority establishment and management, decision support and control structure, the establishment of incentives,

accountability metrics, assessment and intervention strategies, rules and regulations, standardization, etc. must be addressed. By the same token, there is a need for economic models and metrics to determine reasonable pricing and billing for VC services. We anticipate that some forms of VCs will be contract-driven, where the owner of the vehicle or vehicular fleet consents to renting out some form of excess computational or storage capacity. At the same time, mobility concerns dictate that, in addition to the contract-based form of VC, it should be possible to form a VC in an ad hoc manner as necessitated by dynamically changing situations or demand.

VC server consolidation. In conventional clouds, server consolidation is one of the fundamental performance booster and cost-reduction strategies. The idea behind server consolidation is to migrate VMs from lightly loaded servers so that they can be powered off (to save energy) or to migrate VMs from overloaded servers to improve processing or response times. Many strategies for VM consolidation to optimize a given objective functions have been proposed in the cloud computing literature.

It is known that server consolidation can be reduced to *bin-packing*, a known NP-hard problem. Thus, general exact solutions for server consolidation are very unlikely to exist. This has motivated the search for heuristics of which some are quite efficient. Server consolidation in cloud computing is still an active area of investigation. In spite of its fundamental importance, server consolidation has not been addressed in the context of VCs.

Promoting VCs reliability and availability. Dependability is an integrative concept that encompasses a number of attributes. Of these, *reliability* and *availability* are fundamental to conventional clouds and also to VCs. Availability is a measure of the delivery of correct service with respect to the alternation of correct and incorrect service. Reliability is a measure of continuous delivery of correct service, equivalently of the time to the next failure. Depending on the application domain, different emphasis is placed on these attributes.

The high availability requirements of conventional clouds can only be met through redundancy of storage (in addition to execution redundancy, where each user job is be executed by two or more servers). Storage redundancy implies that virtually each data item must be stored at several locations in the network of data centers, often in data centers at various locations around the world.

The dynamically changing availability of compute resources due to vehicles joining and leaving the VC unexpectedly leads to a volatile computing environment where promoting system reliability and availability is very challenging. Yet, if VCs are to see an adoption rate and success similar to that of conventional clouds, reliability and availability issues must be addressed adequately.

Recently, Florin et al. [180] studied the reliability of VCs with *short vehicular residency times*, as is the case in shopping mall parking lots or the short-term parking lot of a major airport. They showed how to enhance system reliability and availability in these types of VCs through a family of redundancy-based job assignment strategies that attempt to mitigate the effect of compute-resource volatility. They offered a theoretical prediction of the corporate MTTF and availability offered by these strategies in terms of the MTTF of individual vehicles and of the MTTR. Extensive simulations using data from shopping mall statistics have confirmed the accuracy of their analytical predictions. To the best of our knowledge, [180] is the *only* paper in the literature that addresses evaluating the MTTF and availability in VCs built on top of vehicles with a short residency time.

Incentivizing participation in VCs. The early VC papers [7, 163] have suggested that, to set up credible VCs, the owners of the vehicles involved will have to be suitably incentivized. For example, we anticipate that in the near future air travelers will park and plug their cars in airport long-term parking lots.

In return for free parking and other perks, they will allow their cars to participate, during their absence, in the airport data center, as discussed in [31].

We now discuss a number of possible approaches to incentivize participation in VCs. The first possible approach will be a lightweight bidding/auction-based solution. Specifically, the data center will announce the list of resources needed by the application along with (1) their qualitative and quantitative attributes, and (2) a specified remuneration level for each resource. In response to the solicitation, vehicles will bid. To make the workforce recruiting operation fast, the number of bidding rounds will be kept to a minimum.

A second approach for incentivizing VC participation is an extension of the approach proposed for scaling cloud resources in [334]. The idea is to provide atomic pricing/compensation per resource used and provide incentives at the individual resource level. For simple tasks involving a small number of basic resources to be provided by edge devices, this approach seems to be very promising.

VC Security and Privacy. It goes without saying that, if VCs are to gain universal acceptance, security and privacy issues need to be suitably addressed. The establishment of trust relationships between the various players is a key component of secure computation and communication. Since some of the vehicles involved in a VC may have met before, the task of establishing proactively a basic trust relationship between vehicles is possible and may be even desirable. Research is needed on developing a trustworthy base, a negotiation and strategy formulation methodology, efficient communication protocols, data processing, etc. VS security and privacy in VC were first addressed in [523] under very general conditions; the authors pointed out novel security problems and challenges specific to VCs and proposed initial solutions to a number of the identified security and privacy problems.

Motivated by this early paper and the importance of security and privacy, a number of authors are investigating security and privacy problems in VCs. For example, Ahmad et al. [14] have identified a number of security threats at various levels of the VC architecture. Similarly, Huang et al. [247] have looked at privacy issues in VCs. It would be interesting to see if lightweight security solutions can be ported to VCs. In spite of these efforts, providing security and privacy in VC remains an open challenge.

VC support for compute-intensive applications. One of the significant research challenges in VCs is to identify conditions under which VCs can support Big Data applications. It is apparent that compute-intensive applications, with stringent data processing requirements, cannot be supported by ephemeral VCs, where the residency time of vehicles in the cloud is too short to support VM setup and migration.

Recently, Florin et al. [177] have identified sufficient conditions under which Big Data applications can be effectively supported by data centers built on top of vehicles in a parking lot. This is the first time researchers are looking at evaluating the feasibility of the VC concept and its suitability for supporting Big Data applications. The main findings of [177] are that, if the residency times of the vehicles are sufficiently long and the interconnection fabric has a sufficient amount of bandwidth, then Big Data applications can be supported effectively by VCs. In spite of this result, more work is needed to understand what it takes for VCs to gracefully support data- and processing-intensive applications.

VC support for Smart Cities. Visionaries depicted the smart cities of the future as fully connected entities supported by various forms of *road-side infrastructure*, as well as advanced in-vehicle resources, such as embedded powerful computing and storage devices, cognitive radios, and multi-modal programmable sensor nodes. As a result, in the near future, vehicles equipped with computing, communication, and sensing capabilities will be organized into ubiquitous and pervasive networks with a

significant Internet presence while on the move. This will revolutionize the driving experience, making it safer, more enjoyable, and more environmentally friendly.

One of the characteristics of Smart Cities is the interconnectivity of the city's infrastructure, which allows traffic data to be collected from various human-generated or machine-generated sources or, indeed, by using the vehicles participating in the traffic. Future vehicles with powerful on-board computers, communication capabilities, and vast sensor arrays are perfect candidates in this hyper-connected environment to utilize to create a fluid VC capable of performing large-scale computations.

The way we see it, the main challenge for VCs in the context of Smart Cities should be aligned with the 2015–2019 strategic priorities recently spelled out by the US-DOT [484]. To demonstrate the relevance of VCs to Smart Cities, the following objectives were proposed:

- 1. Enhance urban mobility by exploring methods and management strategies that increase system efficiency and improve individual mobility through information sharing:** VCs should combine detailed knowledge of real-time traffic flow data with stochastic predictions within a given time horizon to help: (1) the formation of urban platoons containing vehicles with a similar destination and trajectory; (2) adjust traffic-signal timing to reduce unnecessary platoon idling at traffic lights; and (3) present the driving public with high-quality information that will enable them to reduce their trip time and its variability, thus eliminating the conditions that lead to congestion or reduce its effects.
- 2. Avoid congestion of key transportation corridors through cooperative navigation systems:** Congestion-avoidance techniques that become possible in SC environments will be supplemented by route guidance strategies to reduce unnecessary idling and will limit environmental impact of urban transportation.
- 3. Handling nonrecurring congestion:** VCs will explore strategies to efficiently dissipate congestion by a combination of traffic-light retiming and route guidance to avoid more traffic buildup in congested areas.

Managing VC ecosystems. Recently, reference [334] made the point that, given its current trajectory, the complexity of cloud ecosystems will evolve to where traditional resource management strategies will struggle to remain fit for the purpose. Unlike conventional clouds where, for cost efficiency purposes, the compute resources are as uniform as possible, in VCs, the compute resources of individual vehicles are heterogeneous. This, in conjunction with the possible federation of VCs (suggested for the first time by [163]) and the emergence of new services that these VC will have to offer, will make it increasingly hard to manage VC resources efficiently. To the best of our knowledge, the topics of VC federation and creation of VC ecosystems have not been addressed in the VC literature. This promises to be an exciting area for future work.

Concluding remarks. It is well known that implementing cloud services, whether in conventional clouds or VCs, requires substantial architectural support ranging from virtualization, to server consolidation, and to file-system support. The amount of such support must be compared with the resulting performance; conversely, for a desired level of performance, it is of great interest to determine the level of architectural support required. This leads to interesting cost-performance tradeoffs that are fundamental to understanding the feasibility of cloud services and, therefore, cannot be ignored.

In the past decade, various VC architectures and services were outlined in terms of desirable *qualitative* characteristics without much regard to, or credible study of, their feasibility and quantitative

performance characteristics. All this is changing—more and more researchers are turning their attention to *quantitative* aspects of VCs and of the services they contemplate.

The success of conventional cloud computing was due, to a large extent, to its ability to provide quantifiable (i.e., quantitative) functional characteristics such as high scalability, reliability and availability. Cloud service providers have come to equate reliability and availability with customer satisfaction and, ultimately, with revenue. By the same token, if the VCs are to see a widespread adoption, the same quantitative aspects have to be addressed here, too. Feasibility issues in terms of sufficient compute power, communication bandwidth, reliability, availability, and job-duration time are all fundamental quantitative aspects of VCs that need to be studied and understood before one can claim with any degree of certainty that they can support the workload for which they are intended.

It is our belief that reliability, availability, user job duration [177,179,180] and other similar quantitative aspects will play a fundamental contributing role in the large-scale adoption of VCs. We therefore believe that it is time for the research community to devote effort to deepen our understanding of these attributes.

13.5 Final thoughts

A cursory look at the cloud computing literature reveals the extraordinary attention given to this emerging field of computer science. Areas, such as computer architecture, concurrency, data management and databases, resource management, scheduling, and mobile computing, have bloomed in response to the need to find efficient solutions to the challenges brought about by cloud computing. Even somewhat ossified areas such as operating systems have been brought back to life by problems posed by virtualization and containerization.

Several areas critical for the future of cloud computing, including communication and security, demand special attention. Increasing the bandwidth and lowering the communication latency will make cloud computing more attractive for real-time applications and for integrations of services related to IoT. Optimization of communication protocols could lower the latency up to the limits imposed by the laws of physics. Ultimately, communication latency depends on the distance between the producer and the consumer of data.

Computer clouds and mobile devices are in a symbiotic relationship with one another and effective communication to/from clouds and inside the cloud infrastructure has to keep pace with advances in processor and storage technology. Faster cloud interconnects are also necessary to accommodate data-intensive and communication-intensive applications in need of a large number of servers working in concert. Applications in computational sciences and engineering exhibiting fine-grained parallelism would greatly benefit from lower latency.

Data security and privacy are major concerns not properly addressed by existing SLAs. Though sensitive information has been leaked or stolen from large data centers, many cloud users are unaware of the potential dangers they are exposed to when entrusting their data to a third party and trusting the protection guaranteed by SLAs.

Strong encryption protects data in storage, but processing encrypted data is only feasible for some types of queries. Most applications only operate with plaintext data; thus encrypted data has to be decrypted before processing. This creates a window of vulnerability that can be exploited by insider attacks. Hybrid clouds offer an alternative to protect sensitive information. In this case, effective mech-

anisms to hide sensitive information stored on the public cloud and revealed only on the private side of the cloud must be conceived.

Virtualization, in spite of its benefits, creates significant complications for software maintenance. A checkpointed virtual machine containing an older version of an operating system without the current security patches may be activated at a later time, opening a window of vulnerability that could affect the entire cloud infrastructure.

Response times plagued by a heavy-tail distribution cannot be tolerated by most interactive or real-time applications, but eliminating the tail of the latency at the scale of clouds is an enduring challenge. Another enduring challenge is reducing energy consumption and, implicitly, increasing the average server utilization. Elasticity without over provisioning requires accurate knowledge of resource consumption.

Resource reservations can help, but reservations place an additional burden on cloud users expected to know well the needs of their applications. Moreover, accurately predicting resource consumption is possible if and only if the system enforces strict performance isolation, yet another major headache for systems based on multitenancy.

Even skeptics cautioning about the dangers inherent to systems “too big to fail” have to recognize that the cloud ecosystem plays an important role in the modern society and that it has democratized computing, in the same way the web has completely changed the manner we access and use information. The Internet will continue to morph, and the web will evolve to a semantic web or Web 3.0. Thus, it is fair to expect that computer clouds will continue to change under the pressure from consumers of cloud services and from new technologies.

It is hard to predict how the cloud ecosystem will look five or ten years from now, but it should be clear that the disruptive qualities of computer clouds ultimately demand a new thinking in system design. The design of large-scale systems requires an *ab initio* preparation for the unexpected because low probability events occur and can cause major disruptions.

We have seen that the separation of control and routing planes in the Internet is partially responsible for the rapid assimilation of new communication technologies. Only a holistic approach could lead to a similar separation of concerns for computer clouds and allow computing technology to evolve at a lightning speed.

There is a glimmer of hope that machine learning, data analytics, and market-based resource management will play a transformative role in cloud computing [42,332]. As more data are collected after the execution of all instances of an application, it may be possible to construct the application profile, optimize its execution, and, ultimately, optimize the overall system performance. It may also be possible to identify conditions leading to phase transitions and prevent their occurrence that often lead to data center shutdown.

In this maze of challenges and uncertainties, there is one prediction few could argue against: The interest in cloud computing, as well as the need for individuals well trained in this field, will continue to grow for the foreseeable future. The incredible pace of developments in cloud computing poses its own challenges and demands the grasp of fundamental concepts in many areas of computer science and computer engineering, as well as curiosity and the desire to continually learn.

Cloud projects

A

Several projects for students enrolled in a cloud computing class and a research project involving an extensive simulation are discussed in this section. These projects reflect the great appeal of cloud computing for various types of applications. Large-scale simulation of complex systems, such as the one discussed in Sections A.1 and A.3, can only be done using the large pool of resources provided by clouds. Cloud services, such as the one discussed in Sections A.2 and A.4, are another important class of applications.

Design projects in which multiple alternatives must be evaluated and compared, such as the one discussed in Section A.5, benefit from a cloud environment. Several, possibly many, design alternatives can be simulated concurrently; then, the selection of the best alternatives based on different performance objectives can be done by multiple instances running concurrently. High-performance computing applications such as the one discussed in Section A.6 require resources that can only be provided by supercomputers or computer clouds. Clouds are better alternatives than supercomputers due to easier access and lower cost. Applications of machine learning and healthcare applications take advantage of GPUs and TPUs available in some AWS instances. A simulation study of ML algorithm scalability, a cloud-based task management system, and a cloud-based health-monitoring application are discussed in Sections A.7, A.8, and A.9, respectively.

A.1 Cloud simulation of a distributed trust algorithm

Cloud-based simulation for trust evaluation in a Cognitive Radio Networks (CRN) [63] is discussed first. The available communication spectrum is a precious commodity, and the objective of a CRN is to use the communication bandwidth effectively, while attempting to avoid interference with licensed users. Two main functions necessary for the operation of a CRN are spectrum sensing and spectrum management; the former detects unused spectrum, and the latter determines the optimal use of the available spectrum. Spectrum sensing in CRNs is based on information provided by the nodes of the network. The nodes compete for the free channels, and some may supply deliberately distorted information to gain advantage over the other nodes; thus, trust determination is critical for the management of CRNs.

Cognitive radio networks. Research in the last decade reveals a significant temporal and spatial underutilization of the allocated spectrum, thus the motivation to opportunistically harness the vacancies of spectrum at a given time and place.

The original goal of cognitive radio, first proposed at Bell Labs [351,352], was to develop a software-based radio platform that allows a reconfigurable wireless transceiver to automatically adapt

its communication parameters to network availability and to user demands. Today, the focus of cognitive radio is on spectrum sensing [76].

We recognize two types of devices connected to a CRN, primary and secondary; *primary* devices have exclusive rights to specific regions of the spectrum, while *secondary* devices enjoy dynamic spectrum access and are able to use a channel, provided that the primary, licensed to use that channel, is not communicating. Once a primary starts its transmission, the secondary using the channel is required to relinquish it and identify another free channel to continue its operation; this mode of operation is called an *overlay mode*.

Cognitive radio networks are often based on *cooperative spectrum-sensing* strategy. In this mode of operation, each device determines the occupancy of the spectrum based on its own measurements combined with information from its neighbors and then shares its own spectrum occupancy assessment with its neighbors [186,463,464].

Information sharing is necessary because a device alone cannot determine true spectrum occupancy. A secondary device has a limited transmission and reception range; device mobility combined with wireless channel impairments, such as multipath fading, shadowing, and noise, add to the difficulties of gathering accurate information by a single device.

Individual devices of a centralized, or infrastructure-based, CRN send the results of their measurements regarding spectrum occupancy to a central entity, whether a base station, an access point, or a cluster head. This entity uses a set of *fusion rules* to generate the spectrum occupancy report and then distributes it to the devices in its jurisdiction. The area covered by such networks is usually small because global spectrum decisions are affected by the local geography. There is another mode of operation based on the idea that a secondary device operates at a much lower power level than a primary one. In this case, the secondary can share the channel with the primary as long as its transmission power is below a threshold, μ , that has to be determined periodically. In this scenario, the receivers wishing to listen to the primary are able to filter out the “noise” caused by the transmission initiated by secondaries if the signal-to-noise ratio (S/N) is large enough.

We are only concerned with the overlay mode whereby a secondary device maintains an *occupancy report* that gives a snapshot of the current status of the channels in the region of the spectrum it is able to access. The occupancy report is a list of all the channels and their states, e.g., 0 if the channel is free for use and 1 if the primary is active. Secondary devices continually sense the channels they can access to gather accurate information about available channels.

The secondary devices of an ad hoc CRN compete for free channels, and the information one device may provide to its neighbors could be deliberately distorted; malicious devices will send false information to the fusion center in a centralized CRN. Malicious devices could attempt to deny the service or to cause other secondary devices to violate spectrum allocation rules. To *deny the service*, a device will report that free channels are used by the primary. To entice the neighbors to commit FCC violations, the occupancy report will show that channels used by the primary are free. This attack strategy is called *secondary spectrum-data falsification (SSDF)* or a Byzantine attack.¹ Thus, trust determination is a critical issue for CR networks.

Trust. The actual meaning of *trust* is domain and context specific. Consider for example networking; at the MAC-layer (Medium Access Control) multiple-access protocols assume that all senders follow

¹ See Section 10.13 for a brief discussion of Byzantine attacks.

the channel access policy, e.g., in CSMA-CD, a sender senses the channel and then attempts to transmit if no one else does. In a store-and-forward network, trust assumes that all routers follow a best-effort policy to forward packets towards their destination. We shall use node instead of device throughout the remaining of this section.

In the context of cognitive radio, trust is based on the quality of information regarding channel activity provided by a node. The status of individual channels can be assessed by each node, based on the results of its own measurements combined with information provided by its neighbors, as is the case of several algorithms discussed in the literature [103, 463].

The alternative discussed in Section A.2 is to have a cloud-based service that collects information from individual nodes, evaluates the state of each channel based on the information received, and supplies this information on demand. Evaluation of the trust and identification of untrustworthy nodes are critical for both strategies [390].

A distributed algorithm for trust management in cognitive radio. The algorithm computes the trust of node $1 \leq i \leq n$ in each node in its vicinity, $j \in V_i$, and requires several preliminary steps. The basic steps executed by a node i at time t are:

1. Determine node i 's version of the occupancy report for each one of the K channels:

$$S_i(t) = \{s_{i,1}(t), s_{i,2}(t), \dots, s_{i,K}(t)\}. \quad (\text{A.1})$$

In this step node i measures the power received on each of the K channels.

2. Determine the set $V_i(t)$ of the nodes in the vicinity of node i . Node i broadcasts a message, and individual nodes in its vicinity respond with their *NodeId*.
3. Determine the distance to each device $j \in V_i(t)$ using the algorithm described in this section.
4. Infer the power as measured by each device $j \in V_i(t)$ on each channel $k \in K$.
5. Use the location and power information determined in the previous two steps to infer the status of each channel:

$$s_{i,k,j}^{infer}(t) \quad \text{with } 1 \leq k \leq K, j \in V_i(t). \quad (\text{A.2})$$

A secondary node j should have determined: 0 if the channel is free for use, 1 if the primary device is active, and X if it cannot be determined:

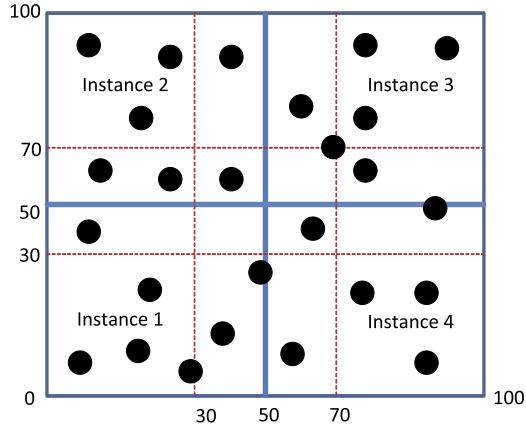
$$s_{i,k,j}^{infer}(t) = \begin{cases} 0 & \text{if secondary node } j \text{ decides that channel } k \text{ is free} \\ 1 & \text{if secondary node } j \text{ decides that channel } k \text{ is used by the primary} \\ X & \text{if no inference can be made.} \end{cases} \quad (\text{A.3})$$

6. Receive the information provided by neighbor $j \in V_i(t)$, $S_{i,k,j}^{recv}(t)$.
7. Compare the information provided by neighbor $j \in V_i(t)$

$$S_{i,k,j}^{recv}(t) = \{s_{i,1,j}^{recv}(t), s_{i,2,j}^{recv}(t), \dots, s_{i,K,j}^{recv}(t)\} \quad (\text{A.4})$$

with the information inferred by node i about node j

$$S_{i,k,j}^{infer}(t) = \{s_{i,1,j}^{infer}(t), s_{i,2,j}^{infer}(t), \dots, s_{i,K,j}^{infer}(t)\}. \quad (\text{A.5})$$

**FIGURE A.1**

Data partitioning for the simulation of a trust algorithm; the area covered is of size 100×100 units. The nodes in the four subareas of size 50×50 units are processed by an instance of the cloud application. The subareas allocated to an instance overlap to enable an instance to have all the information about a node in its coverage area.

8. Compute the number of matches, mismatches, and cases when no inference is possible, respectively:

$$\alpha_{i,j}(t) = \mathcal{M} \left[S_{i,k,j}^{infer}(t), S_{i,k,j}^{recv}(t) \right], \quad (\text{A.6})$$

with \mathcal{M} being the number of matches between the two vectors,

$$\beta_{i,j}(t) = \mathcal{N} \left[S_{i,k,j}^{infer}(t), S_{i,k,j}^{recv}(t) \right], \quad (\text{A.7})$$

with \mathcal{N} being the number of mismatches between the two vectors, and $X_{i,j}(t)$ being the number of cases where no inference could be made.

9. Use the quantities $\alpha_{i,j}(t)$, $\beta_{i,j}(t)$, and $X_{i,j}(t)$ to assess the trust in node j . For example, compute the trust of node i in node j at time t as:

$$\xi_{i,j}(t) = [1 + X_{i,j}(t)] \frac{\alpha_{i,j}(t)}{\alpha_{i,j}(t) + \beta_{i,j}(t)}. \quad (\text{A.8})$$

Simulation of the distributed-trust algorithm. The cloud application is a simulation of a CRN to assess the effectiveness of a particular trust-assessment algorithm. Multiple instances of the algorithm run concurrently on an AWS cloud. The area where the secondary nodes are located is partitioned in several overlapping subareas as in Fig. A.1. The secondary nodes are identified by an instance Id, iId , as well as a global Id, gId . The simulation assumes that the primary nodes cover the entire area; thus their position is immaterial.

The simulation involves a controller and several cloud instances; in its initial implementation, the controller runs on a local system under Linux Ubuntu 10.04 LTS. The controller supplies the data, the

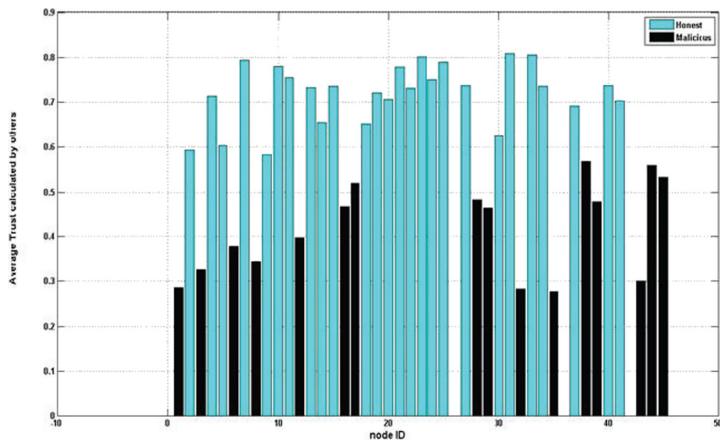


FIGURE A.2

The trust values computed using the distributed-trust algorithm. The secondary nodes programmed to act maliciously have a trust value less than 0.6 and many less than 0.5, lower than that of the honest nodes.

trust program, and the scripts to the cloud instances; the cloud instances run under the Basic 32-bit Linux image on AWS, the so-called *t1.micro*. The instances run the actual trust program and compute the instantaneous trust inferred by a neighbor; the results are then processed by an *awk*² script to compute the average trust associated with a node as seen by all its neighbors. On the next version of the application, the data is stored on the cloud using the S3 service, and the controller also runs on the cloud.

In the simulation discussed here, the nodes with

$$gId = \{1, 3, 6, 8, 12, 16, 17, 28, 29, 32, 35, 38, 39, 43, 44, 45\} \quad (\text{A.9})$$

were programmed to be dishonest. The results show that the nodes programmed to act maliciously have a trust value lower than that of the honest nodes; their trust value is always lower than 0.6 and, in many instances, lower than 0.5; see Fig. A.2. We also observe that the node density affects the accuracy of the algorithm; the algorithm predicts more accurately the trust in densely populated areas. As expected, nodes with no neighbors are unable to compute the trust. In practice, node density is likely to be nonuniform, high density in a crowded area such as a shopping mall and considerably lower in surrounding areas. This indicates that, when trust is computed using information provided by all secondary nodes, we can expect a higher accuracy of trust determination.

² The AWK utility is based on a scripting language and used for text processing; in this application, it is used to produce formatted reports.

A.2 A trust management service

The cloud service discussed in this section [63] is an alternative to the distributed-trust management scheme analyzed in Section A.1. Mobile devices are ubiquitous nowadays and their use will continue to increase. Clouds are emerging as the computing and the storage engines of the future for a wide range of applications. There is a symbiotic relationship between the two: Mobile devices can consume, as well as produce, very large amounts of data, while computer clouds have the capacity to store and deliver such data to the user of a mobile device. To exploit the potential of this symbiotic relationship, we propose a new cloud service for the management of wireless networks.

Mobile devices have limited resources. While new generations of smart phones and tablet computers are likely to use multicore processors and have a fair amount of memory, power consumption is still, and will continue to be in the near future, a major concern. Thus, it seems reasonable to delegate compute-and data-intensive tasks to the cloud.

Transferring computations related to CRN management to a cloud supports the development of new, possibly more accurate, resource management algorithms. For example, algorithms to discover communication channels currently in use by a primary transmitter could be based on past history but are not feasible when the trust is computed by the mobile device. Such algorithms require massive amounts of data and can also identify malicious nodes with high probability.

Mobile devices such as smartphones and tablets are able to communicate using two networks: (i) a cellular wireless network; and (ii) a WiFi network. The service we propose assumes that a mobile device uses the cellular wireless network to access the cloud, while the communication over the WiFi channel is based on cognitive radio (CR). The amount of data transferred using the cellular network is limited by the subscriber's data plan, but no such limitation exists for the WiFi network. The cloud service discussed next will allow mobile devices to use the WiFi communication channels in a cognitive radio network environment and will reduce the operating costs for the end-users.

While the focus of our discussion is on trust management for CRN networks, the cloud service we propose can be used for tasks other than the bandwidth management. For example, routing in a mobile ad hoc network, detection, and isolation of noncooperative nodes, and other network management and monitoring functions could benefit from the identification of malicious nodes.

Model assumptions. The cognitive radio literature typically analyzes networks with a relatively small number of nodes active in a limited geographic area; thus, all nodes in the network sense the same information on channel occupancy. Channel impairments such as signal fading, noise, and so on, cause errors and lead trustworthy nodes to report false information. We consider networks with a much larger number of nodes distributed over a large geographic area; because the signal strengths decays with the distance, we consider several rings around a primary tower. We assume a generic fading model given by the following expression:

$$\gamma_k^i = T_k \times \frac{A^2}{s_{ik}^\alpha}, \quad (\text{A.10})$$

where γ_k^i is the received signal strength on channel k at location of node i , A is the frequency constant, $2 \leq \alpha \leq 6$ is path loss factor, s_{ik}^α is the distance between primary tower P_k and node i , and T_k is the transition power of primary tower P_k transmitting on channel k .

In our discussion, we assume that there are K channels labeled $1, 2, \dots, K$ and that the primary transmitter P^k transmits on channel k . The algorithm is based on several assumptions regarding the

secondary nodes, the behavior of malicious nodes, and the geometry of the system. First, we assume the following:

- The secondary nodes are mobile devices; some are slow-moving, while others are fast-moving.
- They cannot report their position because they are not equipped with a GPS system.
- The clocks of the mobile devices are not synchronized.
- The transmission and reception range of a mobile device can be different.
- The transmission range depends on the residual power of each mobile device.

We assume that the malicious nodes in the network are a minority and their behavior is captured by the following assumptions:

- The misbehaving nodes are malicious, rather than selfish; their only objective is to hinder the activity of other nodes whenever possible, a behavior distinct from the one of selfish nodes motivated to gain some advantage.
- The malicious nodes are uniformly distributed in the area we investigate.
- The malicious nodes do not collaborate in their attack strategies.
- The malicious nodes change the intensity of their Byzantine attack in successive time slots; similar patterns of malicious behavior are easy to detect, and an intelligent attacker is motivated to avoid detection.

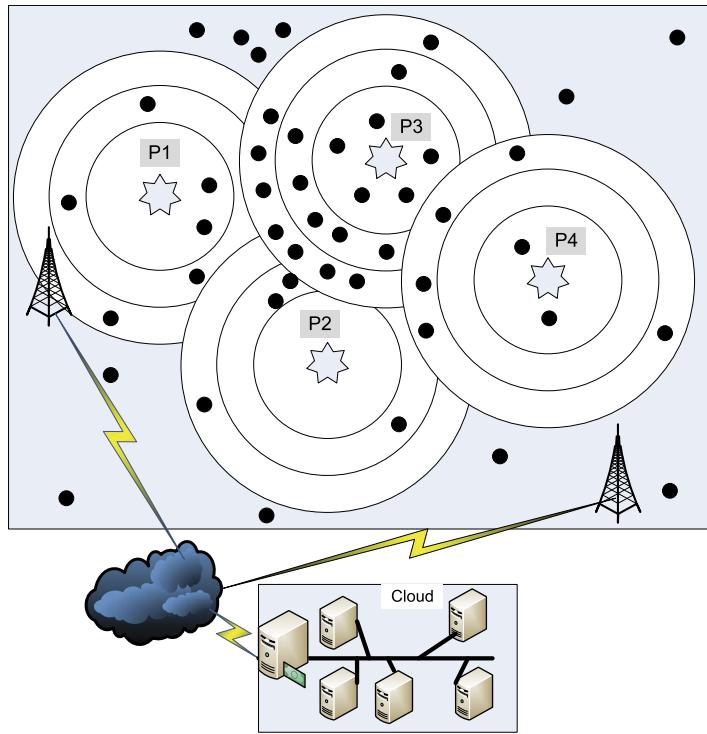
The geometry of the system is captured by Fig. A.3. We distinguish primary and secondary nodes and the cell towers used by the secondary nodes to communicate with the service running on the cloud.

We use a majority voting rule for a particular ring around a primary transmitter; the global decision regarding the occupancy of a channel requires a majority of the votes. Since the malicious nodes are a minority and they are uniformly distributed, the malicious nodes in any ring are also a minority; thus, a ring-based majority fusion is representative of accurate occupancy for the channel associated with the ring.

All secondary nodes are required to register first and then to transmit periodically their current power levels, as well as their occupancy report for each one of the K channels. As mentioned in the introductory discussion, the secondary nodes connect to the cloud using the cellular network. After a mobile device is registered, the cloud application requests the cellular network to detect its location; the towers of the cellular network detect the location of a mobile device by triangulation with an accuracy that is a function of the environment and is of the order of 10 m. The location of the mobile device is reported to the cloud application every time they provide an occupancy report.

The nodes that do not participate in the trust computation will not register in this cloud-based version of the resource management algorithm, thus they do not receive the occupancy report and cannot use it to identify free channels. Obviously, if a secondary node does not register, it cannot influence other nodes and prevent them from using free channels, or tempt them to use busy channels.

In the registration phase, a secondary node transmits its MAC address, and the cloud responds with the tuple (Δ, δ_s) . Here, Δ is the time interval between two consecutive reports, chosen to minimize the communication and the overhead for sensing the status of each channel. To reduce the communication overhead, secondary nodes should transmit only the changes from the previous status report. $\delta_s < \Delta$ is the time interval to the first report expected from the secondary node. This scheme provides a pseudo-synchronization, so that the data collected by the cloud and used to determine the trust is based on observations made by the secondary nodes at about the same time.

**FIGURE A.3**

Schematic representation of a CR layout; four primary nodes, P_1-P_4 , a number of mobile devices, two towers for a cellular network, and a cloud are shown. Not shown are the hotspots for the WiFi network.

An algorithm for trust evaluation based on historical information. The cloud computes the probable distance d_i^k of each secondary node i from the known location of a primary transmitter, P^k . Based on signal attenuation properties, we conceptualize N circular rings centered at the primary where each ring is denoted by \mathcal{R}_r^k , with $1 \leq r \leq N$ the ring number.

The radius of a ring is based on the distance d_r^k to the primary transmitter P^k . A node at a distance $d_i^k \leq d_1^k$ is included in the ring \mathcal{R}_1^k , nodes at distance $d_1^k < d_i^k \leq d_2^k$ are included in the ring \mathcal{R}_2^k , and so on. The closer to the primary, the more accurate the channel occupancy report of the nodes in the ring should be. Call n_r^k the number of nodes in ring \mathcal{R}_r^k .

At each report cycle at time t_q , the cloud computes the occupancy report for channel $1 \leq k \leq K$ used by primary transmitter P^k . The status of channel k reported by node $i \in \mathcal{R}_r^k$ is denoted as $s_i^k(t_q)$. Call $\sigma_{one}^k(t_q)$ the count of the nodes in the ring \mathcal{R}_r^k reporting that the channel k is not free (reporting $s_i^k(t_q) = 1$) and $\sigma_{zero}^k(t_q)$ the count of those reporting that the channel is free (reporting $s_i^k(t_q) = 0$):

$$\sum_{one}^k(t_q) = \sum_{i=1}^{n_r^k} s_i^k(t_q) \quad \text{and} \quad \sigma_{zero}^k(t_q) = n_r^k - \sigma_{one}^k(t_q). \quad (\text{A.11})$$

Then, the status of channel k reported by the nodes in the ring R_r^k is determined by majority voting as:

$$\sum_{R_r}^k(t_q) \begin{cases} = 1 & \text{when } \sum_{one}^k(t_q) \geq \sum_{zero}^k(t_q) \\ = 0 & \text{otherwise} \end{cases}. \quad (\text{A.12})$$

To determine the trust in node i , we compare $s_i^k(t_q)$ with $\sigma_{R_r}^k(t_q)$; call $\alpha_{i,r}^k(t_q)$ and $\beta_{i,r}^k(t_q)$ the number of matches and, respectively, mismatches in this comparison for each node in the ring R_r^k . We repeat this procedure for all rings around P^k and construct

$$\alpha_i^k(t_q) = \sum_{r=1}^{n_r^k} \alpha_{i,r}^k(t_q) \quad \text{and} \quad \beta_i^k(t_q) = \sum_{r=1}^{n_r^k} \beta_{i,r}^k(t_q). \quad (\text{A.13})$$

Node i will report the status of the channels in the set $C_i(t_q)$, the channels with index $k \in C_i(t_q)$; then, quantities $\alpha_i(t_q)$ and $\beta_i(t_q)$ with $\alpha_i(t_q) + \beta_i(t_q) = |C_i(t_q)|$ are

$$\alpha_i(t_q) = \sum_{k \in C_i} \alpha_i^k(t_q) \quad \text{and} \quad \beta_i(t_q) = \sum_{k \in C_i} \beta_i^k(t_q). \quad (\text{A.14})$$

Finally, the global trust in node i is a random variable

$$\zeta_i(t_q) = \frac{\alpha_i(t_q)}{\alpha_i(t_q) + \beta_i(t_q)}. \quad (\text{A.15})$$

The trust in each node at each iteration is determined using a strategy similar to the one discussed earlier; its status report, $S_j(t)$, contains only information about the channels it can report on and only if the information has changed from the previous reporting cycle.

Then, a statistical analysis of the random variables for a window of time W , $\zeta_j(t_q), t_q \in W$ allows us to compute the moments and a 95% confidence interval. Based on these results, we assess if node j is trustworthy and eliminate the untrustworthy ones when we evaluate the occupancy map at the next cycle. We continue to assess the trustworthiness of all nodes and may accept the information from node j when its behavior changes.

Let us now discuss the use of historical information to evaluate trust. We assume a sliding window $W(t_q)$ consists of n_w time slots of duration τ . Given two decay constants k_1 and k_2 , with $k_1 + k_2 = 1$, we use an exponential averaging giving decreasing weight to old observations. We choose $k_1 \ll k_2$ to give more weight to the past actions of a malicious node. Such nodes attack only intermittently and try to disguise their presence with occasional good reports; the misbehavior should affect the trust more than the good actions. The history-based trust requires the determination of the two quantities

$$\alpha_i^H(t_q) = \sum_{i=0}^{n_w-1} \alpha_i(t_q - i\tau) k_1^i \quad \text{and} \quad \beta_i^H(t_q) = \sum_{i=0}^{n_w-1} \beta_i(t_q - i\tau) k_2^i. \quad (\text{A.16})$$

Then, the history-based trust for node i valid only at times $t_q \geq n_w \tau$ is

$$\zeta_i^H(t_q) = \frac{\alpha_i^H(t_q)}{\alpha_i^H(t_q) + \beta_i^H(t_q)}. \quad (\text{A.17})$$

For times $t_q < n_w \tau$, the trust will be based only on a subset of observations rather than a full window based on n_w observations.

This algorithm can also be used in regions where the cellular infrastructure is missing. An ad hoc network could allow the nodes that cannot connect directly to the cellular network to forward their information to nodes closer to the towers and then to the cloud-based service.

Simulation of the history-based algorithm for trust management. The aim of the history-based trust evaluation is to distinguish between trustworthy and malicious nodes. We expect the ratio of malicious to trustworthy nodes and the node density to play an important role in this decision. The node density ρ is the number of nodes per unit of area. In our simulation experiments, the size of the area is constant, but the number of nodes increases from 500 to 2000; thus the node density increases by a factor of four. The ratio of the number of malicious to the total number of nodes varies between $\alpha = 0.2$ to a worst case of $\alpha = 0.6$.

The performance metrics we consider are: the average trust for all nodes, the average trust of individual nodes, and the error of honest/trustworthy nodes. We wish to see how the algorithm behaves when the density of the nodes increases; we consider four cases with 500, 1 000, 1 500, and 2 000 nodes on the same area; thus we allow the density to increase by a factor of four. We also investigate the average trust when α , the ratio of malicious nodes to the total number of nodes increases from $\alpha = 0.2$ to $\alpha = 0.4$ and, finally, to $\alpha = 0.6$.

This straightforward data partitioning strategy for the distributed trust management algorithm is not a reasonable one for the centralized algorithm because it would lead to excessive communication among the cloud instances. Individual nodes may contribute data regarding primary transmitters in a different subarea; to evaluate the trust of each node, the cloud instances would have to exchange a fair amount of information. This data partitioning would also complicate our algorithm that groups together secondary nodes based on the distance from the primary one.

Instead, we allocate to each instance a number of channels, and all instances share the information about the geographic position of each node; the distance of a secondary node to any primary one can then be easily computed. This data partitioning strategy scales well in the number of primaries. Thus, it is suitable for simulation in large metropolitan areas but may not be able to accommodate cases when the number of secondaries is of the order of 10^8 – 10^9 .

The objectives of our studies are to understand the limitations of the algorithm; the aim of the algorithm is to distinguish between trustworthy and malicious nodes. We expect that the ratio of malicious to trustworthy nodes, and the node density, should play an important role in this decision. The measures we examine are the average trust for all nodes and the average trust of individual nodes.

The effect of the malicious versus trustworthy node ratio on the average trust. We report the effect of the malicious versus trustworthy node ratio on the average trust when the number of nodes increases. The average trust is computed separately for the two classes of nodes and enables us to determine if the algorithm is able to clearly separate them.

Recall that the area is constant, thus, when the number of nodes increases, so does the node density. First, we consider two extreme case; the malicious nodes represent only 20% of the total number of nodes and an unrealistically high presence, of 60%. Then, we report on the average trust when the number of nodes is fixed and the malicious nodes represent an increasing fraction of the total number of nodes.

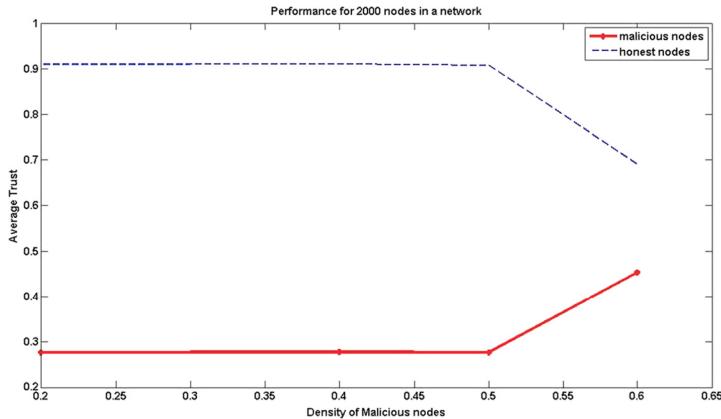


FIGURE A.4

The average trust function of α for a population size of 2000 nodes. As long as malicious nodes represent 50% or less of the total number of nodes, the average trust of malicious nodes is below 0.3, while the one of trustworthy nodes is above 0.9 in a scale of 0 to 1.0. As the number of nodes increases, the distance between the average trust of the two classes becomes larger, and even larger when $\alpha > 0.5$, i.e., the malicious nodes are in majority.

Results reported in [63] show that, when the malicious nodes represent only 20% of all nodes, there is a clear distinction between the two groups. The malicious nodes have an average trust of 0.28, and the trustworthy ones have an average trust index of 0.91, regardless of the number of nodes.

When the malicious nodes represent 60% of all the nodes, then the number of nodes plays a significant role; when the number of nodes is small, the two groups cannot be distinguished since their average trust index is almost equal, 0.55, although the honest nodes have a slightly above-average trust value. When the number of nodes increases to 2000 and the node density increases four folds, then the average trust of the first (malicious) group decreases to 0.45, and for the second (honest) group, it increases to about 0.68.

This result is not unexpected; it only shows that the history-based algorithm is able to classify the nodes properly even when the malicious nodes are a majority, a situation we do not expect to encounter in practice. This effect is somewhat surprising; we did not expect that, under these extreme conditions, the average of the trust of all nodes will be so different for the two groups. A possible explanation is that our strategy to reward constant good behavior, rather than occasional good behavior, and designed to mask the true intentions of a malicious node, works well.

Fig. A.4 shows the average trust function of α , the ratio of malicious versus total number of nodes. The results confirm the behavior discussed earlier; we see a clear separation of the two classes only when the malicious nodes are in the minority. When the density of malicious nodes approaches a high value so that they are in majority, the algorithm still performs as evident from the figure that the average trust for honest nodes even at high value of α is more than for malicious nodes. Thus, the trusts reflect the aim of isolating the malicious from the honest sets of nodes. We also observe that the separation is clearer when the number of nodes in the network increases.

The benefits of a cloud-based service for trust management. A cloud service for trust management in cognitive networks can have multiple technical and economical benefits [95]. The service is likely to have a broader impact than the one discussed here; it could be used to support a range of important policies in wireless network where many decisions require the cooperation of all nodes. A history-based algorithm to evaluate the trust and detect malicious nodes with high probability is at the center of the solution we have proposed [63].

A centralized, history-based algorithm for bandwidth management in CRNs has several advantages over the distributed algorithms discussed in the literature:

- It drastically reduces the computations a mobile device is required to carry out to identify free channels and avoid penalties associated with interference with primary transmitters.
- It enables a secondary node to get information about channel occupancy as soon as it joins the system and later on demand; this information is available even when a secondary node is unable to receive reports from its neighbors, or when it is isolated.
- It does not require the large number of assumptions critical to the distributed algorithms.
- The dishonest nodes can be detected with high probability, and their reports can be ignored; thus, over time, the accuracy of the results increases. Moreover, historic data could help detect a range of Byzantine attacks orchestrated by a group of malicious nodes.
- It is very likely to produce more accurate results than the distributed algorithm because the reports are based on information from all secondary nodes reporting on a communication channel used by a primary, not only those in its vicinity; a higher node density increases the accuracy of the predictions. The accuracy of the algorithm is a function of the frequency of the occupancy reports provided by the secondary nodes.

The centralized trust-management scheme has several other advantages. First, it can be used not only to identify malicious nodes and provide channel occupancy reports, but it also can manage the allocation of free channels. In the distributed case two nodes may attempt to use a free channel and collide; this situation is avoided in the centralized case. At the same time, malicious nodes can be identified with high probability and be denied access to the occupancy report.

The server could also collect historic data regarding the pattern of behavior of the primary nodes and use this information for the management of free channels. For example, when a secondary node requests access for a specific length of time, the service may attempt to identify a free channel likely to be available for that time.

The trust management may also be extended to other network operations, such as routing in a mobile ad hoc network; the strategy in this case would be to avoid routing through malicious nodes.

A.3 Simulation of traffic management in a smart city

The objective of the project is to study the traffic crossing the center of a city for variations in traffic intensities and traffic-light scheduling strategies.

The layout of the city center. Rectangular grid with n rows and m columns. There are NS (North–South) avenues and EW (East–West) streets. All avenues and streets are one way and have multiple lanes.

- The NS and SN avenues are $\mathcal{A}^{NS}(i)$ and \mathcal{A}_i^{SN} , $i \leq m$, respectively. The EW and WE streets are \mathcal{S}_j^{EW} and \mathcal{S}_j^{WE} , $j \leq n$, respectively.
- All distances are measured in c units; one unit equals the average car length plus an average required distance from the car ahead.
- The distance between rows i and $i + 1$ and between columns j and $j + 1$ are denoted by d_i , $1 \leq i \leq n$ and d_j , $1 \leq j \leq m$, respectively, and $\min(d_i, d_j) > kc$ with $k > 100$.
- The direction of one-way traffic alternates on both avenues and streets; \mathcal{A}_j^{NS} , $j \in \{1, 3, \dots\}$ and \mathcal{A}_j^{SN} , $j \in \{2, 4, \dots\}$; similarly, \mathcal{S}_i^{EW} , $i \in \{1, 3, \dots\}$ and \mathcal{S}_i^{WE} , $i \in \{2, 4, \dots\}$. Avenues and streets have either two or three lanes. Call L_j^{NS} the number of lanes of \mathcal{A}_j^{NS} .

The traffic lights are installed at all intersections $\mathcal{I}_{i,j}$, $1 \leq i \leq n$, $1 \leq j \leq m$. The traffic lights $\mathcal{I}_{i,j}$, $1 < i < n$, $1 < j < m$ allow left turns. Call $\tau_{i,j}^{gNS}(t)$, $\tau_{i,j}^{gNW}(t)$, $\tau_{i,j}^{gEW}(t)$ and $\tau_{i,j}^{gES}(t)$ the duration of the green light for the cycle starting at time t for directions NS, NW, EW, and ES, respectively, of traffic light $\mathcal{I}_{i,j}$.

Cars enter and exit the grid from all directions: NS, SN, EW, and WE.

- Each car has an associated path. For example, the path P^k of car C_i^k , $1 \leq i \leq MaxConv$ entering the grid is described by the pair entry point $\mathcal{I}_{i,j}^{k,entry}$ and exit point $\mathcal{I}_{i,j}^{k,exit}$, with $i = 1$ or $i = n$, and $j = 1$ or $j = m$, and at most two intermediate turning points $\mathcal{I}_{i,j}^{k,turn1}$ and $\mathcal{I}_{i,j}^{k,turn2}$ with $1 < i < n$ and $1 < j < m$. Call t_i^{in} the time when a car enters the grid.
- Cars entering the grid are grouped in “convos” of various sizes. Convoys can be split when cars leave it to turn left or right or when the light turns red. Convoys are merged when cars enter the convoy or when cars from another convoy join one convoy stopped at a traffic light.
- Cars whose paths require a right turn should be in the right lane of an avenue or street, and cars whose path requires a left turn should be in the left lane. Cars that do not turn should be in the center lane.
- Call $v_{i,j}^{in}(t, s)$ and $v_{p,q}^{out}(t, s)$ the number of cars in the convoys entering and, respectively, exiting the grid at intersections $\mathcal{I}_{i,j}$ and $\mathcal{I}_{p,q}$, in the interval $s - t$.
- Call $\Phi_{i,j}^{in}(t, s)$ the traffic intensity as the number of cars entering the grid in the interval $s - t$.

Hints for project implementation. Create a description file with separate sections describing: (i) the layout including n , m , c ; for example, $n = 100$, $m = 80$, and $c = 12$ m; (ii) the initial traffic lights setting; and (iii) the car arrival process and the convoys at each entry point. The simulation will require several data structures including:

- Car records. Each record includes: (i) static data such as: the CarId, the entry point, the exit point, the entry time, and the path; and (ii) dynamic data such as: a list of all traffic lights it had to wait for the green light and the waiting time, the speed on each block, and the time of exit.
- Traffic-light records. Each record should include:
 1. static data such as: TrafficLightId; the location as the intersection of avenue \mathcal{A}_i^{XX} with street \mathcal{S}_j^{YY} ;
 2. dynamic data such as: the number of cycles (red + yellow + green); for each cycle, it should include the times for green in each direction and the identity of convoys waiting.

To evaluate the performance, compute the *average transition time* of all cars in the interval (s, t) for several traffic management algorithms:

Dumb scheduling—traffic lights follow a static, deterministic schedule.

Individual self-managed scheduling—a traffic light switches to green in direction DD when the length of the queue of cars waiting to cross the intersection exceeds a threshold L_{sw} .

Coordinated scheduling—modify the previous algorithm to include communication among neighboring nodes; use a publish–subscribe algorithm in which each traffic light subscribes to messages sent by its four neighbors.

Convoy-aware algorithm—attempts to create a green wave for the largest convoys and anticipate the setting of green lights to next intersections at the time the convoy reaches the intersection.

Project implementation.³ The simulation includes a graphical user interface displaying the smart city center. The four algorithms are called: Dumb, Self-Managed, Coordinated, and Convoy. The implementation uses Java Swing API and defines several classes:

1. Road—manages the layout of the grid.
2. Traffic Point—handles intersections and entry and exit points.
3. PaintGrid—is responsible for continuously painting and repainting the grid for a given time interval.
4. Frame—sets the size of the Frame Window used by the Canvas to draw the grid.
5. Car—creates cars and feeds them to a Java hashmap dynamic list.
6. Schedule—manages changing of the traffic lights, and implements the logic of the four scheduling algorithms and controls the car arrival process.
7. Convoy—generates new convoy and adds cars to a convoy.
8. Statistics—gathers traffic statistics.
9. StatWindow—displays traffic statistics.

The car arrival stochastic process has an exponential distribution of the interarrival times

$$p(t) = \lambda e^{-\lambda t}. \quad (\text{A.18})$$

Lower values of parameter λ in Eq. (A.18) result in higher concentration of cars for the first few simulation cycles, while higher values of λ scatter cars on larger number of simulation cycles; see Fig. A.5. A record describing the path of a car is created at the time when a car enters the grid. The path of a car could have zero turns if the car enters and exits on the same street or avenue. The path will have one turn if the car enters on one street and exits from an avenue, or enters an avenue and exits a street. The path will have two turns if the car enters one street and exits from another street or enters an avenue and exits from a different avenue.

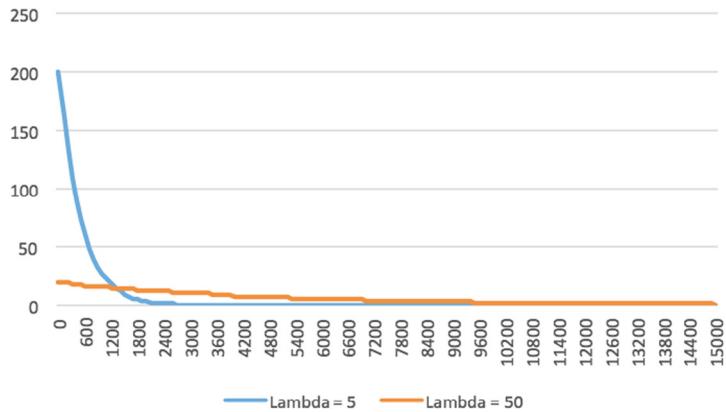
The next step involves the logic of car movement through an intersection. The traffic light checks if the car is moving straight or turning, using car path information. The state of the traffic-light status determines whether the car should decelerate to stop or move ahead and accelerate towards the next intersection. If the car turns, then a car speed-switch function is called.

³ The implementation of the project is due to the group of five students: Ahmed Alhazmi, Austin Jerome, Ahmad Qutbuddin, Sai Lalitha Renduchintala, and Wendelyn Sanabria.

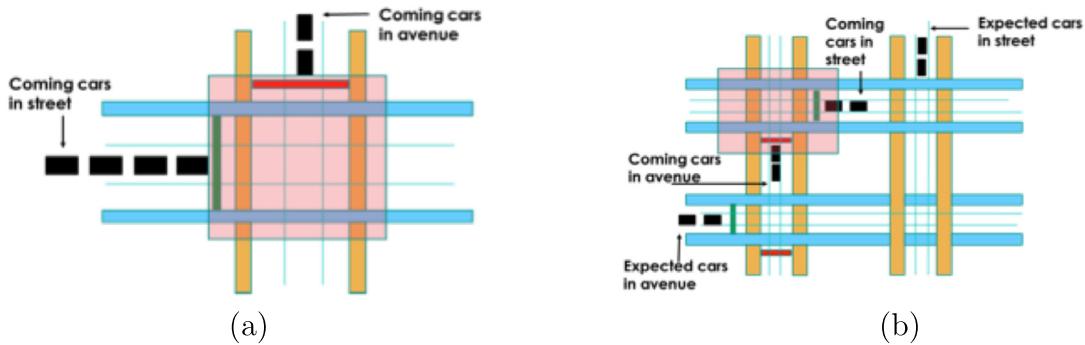
Shown next is a segment of the configuration file specifying the layout and some of parameters related to cars.

```
import java.io.BufferedReader;
import java.io.File;
import java.io.FileInputStream;
import java.io.FileNotFoundException;
import java.io.IOException;
import java.io.InputStreamReader;
import java.util.StringTokenizer;

public class Configuration {
    // Simulation class:
    //      Exponential Car Insertion Rate; Number of Cars
    protected int Lambda = 15;
    protected int NumberOfCars = 350;
    // Grid class:
    //      Number of Streets and Avenues
    //      Maximum and Minimum Block Side Length in c units
    protected int NumberOfStreets = 3;
    protected int NumberOfAvenues = 3;
    protected int MaximumBlockSide = 35;
    protected int MinimumBlockSide = 35;
    // Road class:
    //      Number of Forward and Turning Lanes
    protected int NumberOfForwardLanes = 2;
    protected int NumberOfTurningLanes = 1;
    // TrafficLight class
    //      Maximum Red and Green time in seconds
    //      Maximum Red Time could be Maximum Green Time + Yellow Time
    //      Yellow Time in seconds; Intersection light initial status (TBD)
    //      Scheduling Scheme D, S, C, V
    protected int MaxRedTime = 4000;
    protected int MaxGreenTime = 3000;
    protected int YellowTime = 1000;
    protected char SchedulingScheme = 'C';
    // Car class:
    //      Maximum Car Speed in c/second unit
    //      Car Acceleration in c/second2 unit; Car length and width in pixels
    protected int CarSpeed = 5;
    protected int CarAcceleration = 1;
    protected int CarLength = 6;
    protected int CarWidth = 3;
    protected int Clearance = 2;
```

**FIGURE A.5**

The number of cars entering the system function of time represented by simulation cycles. Higher values of the parameter λ in Eq. (A.18) scatter cars for a larger number of simulation cycles.

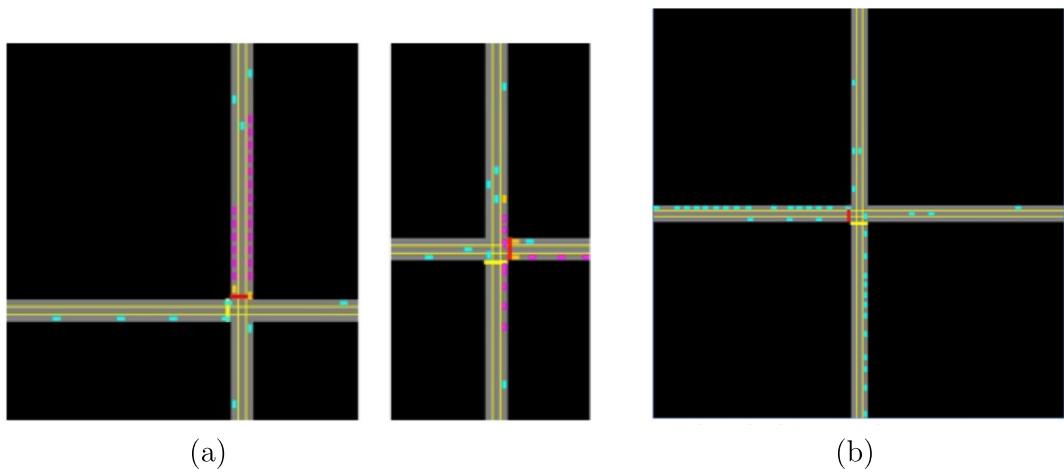
**FIGURE A.6**

(a) Self-managed scheduling uses the information about the number of cars to change the traffic light. (b) In the coordinated scheduling, self-managed intersections coordinate their traffic lights.

The self-managed scheduling uses two queues for the traffic lights at the intersection of streets and avenues. The length of the car queue at the traffic lights triggers changes of the traffic lights, as seen in Fig. A.6(a).

The coordinated scheduling algorithm was implemented as an extension of the self-managed scheduling algorithm. Now, an intersection communicates with its neighboring intersections to determine the number of cars that are passing through. This allows the next intersection along the path to turn green when a larger number of cars will reach the intersection; see Fig. A.6(b).

The convoy scheduling algorithm treats a number of n cars separated by a car length, c , to be assimilated with the larger car of length $n \times c$. Once a convoy is in motion, a traffic light extends its

**FIGURE A.7**

(a) Convoy of cars waiting at an intersection (left); traffic lights extend the yellow-light holds to allow a convoy to cross the intersection (right). (b) The simulation of maximum capacity scenario based on an ideal orchestration.

yellow state to allow the entire convoy to pass. Once a car turns, the whole convoy is broken up; see Fig. A.7(a).

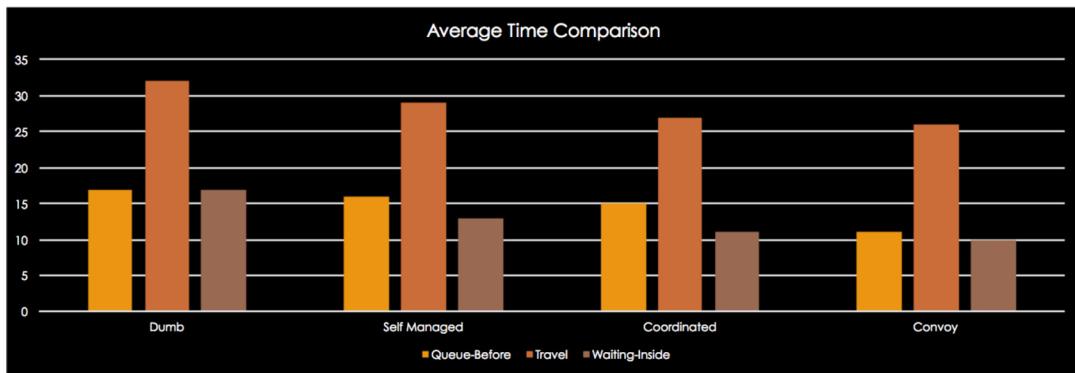
We have also investigated an ideal orchestrations scenario when the distance between cars allows the two continuous streams of cars entering every intersection from the two directions to pass alternatively through without stopping. This ideal traffic light-less scenario seen in Fig. A.7(b) allows us to determine the maximum capacity of the grid.

Some of the simulation results are discussed next. The average transit time decreases steadily from about 36.8 units of simulated time for the dumb traffic lights scheduling algorithm to 32.6 for the self-managed, 32 for the coordinated, and 31.4 for the convoy. The average waiting time is 19.8, 17.4, and 15.2 for the self-managed, coordinated, and convoys traffic lights scheduling algorithms, respectively.

Fig. A.8 shows also a comparison of four algorithms. For each algorithm, it shows the waiting time for entering the grid, the transit time through the grid, and the waiting time at traffic lights while transiting through grid. The convoy scheduling algorithm performs best. Fig. A.9 shows the confidence intervals from 50 simulation runs for the coordinated-scheduling algorithm.

A.4 A cloud service for adaptive data streaming

In this section, we discuss a cloud application related to data streaming [394]. Data streaming is the name given to the transfer of data at a high rate with real-time constraints. Multimedia applications, such as music and video streaming, high-definition television (HDTV), scientific applications that process a continuous stream of data collected by sensors, the continuous backup copying to a storage medium of the data flow within a computer, and many other applications require the transfer of real-time data

**FIGURE A.8**

The waiting time before entering the grid, the transit time through the grid, and the waiting time at traffic lights while transiting the grid. The results displayed are for the four traffic-lights scheduling algorithms, dumb, self-managed, correlated, and convoys.

	Total Number of Cars	Average Transient Time	Average Waiting Time
Maximum	495	36	23
Minimum	354	31	13
Average	431.32	32.92	17.94
Standard Deviation	34.83299585	1.11067547	2.395078287
95 % Confidence	9.655036433	0.307857876	0.663869631
Upper	441	33	19
Lower	421	32	17

FIGURE A.9

Confidence intervals for the car transit time.

at a high rate. For example, to support real-time human perception of the data, multimedia applications have to make sure that enough data is being continuously received without any noticeable time lag.

We are concerned with the case when the data streaming involves a multimedia application connected to a service running on a computer cloud. The stream could originate from the cloud, as is the case of the iCloud service provided by Apple, or could be directed toward the cloud, as in the case of a real-time data collection and analysis system.

Data streaming involves three entities: the sender, a communication network, and a receiver. The resources necessary to guarantee the timing constraints include CPU cycles and buffer space at the sender and the receiver and network bandwidth. Adaptive data streaming determines the data rate based on the available resources. Lower data rates imply lower quality, but they reduce the demands for system resources.

Adaptive data streaming is possible only if the application permits trade-offs between quantity and quality. Such trade-offs are feasible for audio and video streaming that allow lossy compression, but are not acceptable for many applications that process a continuous stream of data collected by sensors.

Data streaming requires accurate information about all resources involved, and this implies that the network bandwidth has to be constantly monitored; at the same time, the scheduling algorithms should be coordinated with memory management to guarantee the timing constraints. Adaptive data streaming poses additional constraints because the data flow is dynamic. Indeed, once we detect that the network cannot accommodate the data rate required by an audio or video stream we have to reduce the data rate to convert to a lower-quality audio or video. Data conversion can be done on the fly, and, in this case, the data flow on the cloud has to be changed.

Accommodating dynamic data flows with timing constraints is nontrivial; only about 18% of the top 100 global video web sites use ABR (Adaptive Bit Rate) technologies for streaming [460].

This application stores the music files in S3 buckets, and the audio service runs on the EC2 platform. In EC2, each virtual machine functions as a virtual private server and is called an *instance*; an instance specifies the maximum amount of resources available to an application and the interface for that instance, as well as the cost per hour.

EC2 allows the import of VM images from the user environment to an instance through a facility called *VM import*. It also distributes automatically the incoming application traffic among multiple instances using the *elastic load balancing* facility. EC2 associates an *elastic IP address* with an account. This mechanism allows a user to mask the failure of an instance and remap a public IP address to any instance of the account, without the need to interact with the software support team.

The adaptive audio streaming involves a multiobjective optimization problem. We wish to convert the highest quality audio file stored on the cloud to a resolution corresponding to the rate that can be sustained by the available bandwidth; at the same time, we wish to minimize the cost on the cloud site and also minimize the buffer requirements for the mobile device to accommodate the transmission jitter. Finally, we wish to reduce to a minimum the start-up time for the content delivery.

A first design decision is whether or not data streaming should only begin after the conversion from the WAV to MP3 format has been completed, or if it should proceed concurrently with conversion, in other words, start as soon as several MP3 frames have been generated; another question is whether or not the converted music file should be saved for later use or discarded.

To answer these questions, we experimented with conversion from the highest quality audio files, which require a 320 Kbps data rate to lower quality files corresponding to 192, 128, 64, 32, and finally 16 Kbps. If the conversion time is small and constant, there is no justification for pipelining data conversion and streaming, a strategy that complicates the processing flow on the cloud. It makes sense to cache the converted copy for a limited period of time with the hope that it will be reused in the next future.

Another design decision is how the two services should interact to optimize the performance; two alternatives come to mind:

1. The audio service running on the EC2 platform requests the data file from the S3, converts it, and, eventually, sends it back. The solution involves multiple delays and it is far from optimal.
2. The S3 bucket is mounted as an EC2 drive. This solution reduces considerably the start-up time for audio streaming.

Table A.1 Conversion time in s on a EC2 *t1.micro* server platform; the source file is of the highest audio quality, 320 Kbps. The individual conversions are labeled C1 to C10; \bar{T}_c is the mean conversion time.

Bit-rate (Kbps)	Audio file size (MB)	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	\bar{T}_c
192	6.701974	73	43	19	13	80	42	33	62	66	36	46.7
128	4.467982	42	46	64	48	19	52	52	48	48	13	43.2
64	2.234304	9	9	9	9	10	26	43	9	10	10	14.4
32	1.117152	7	6	14	6	6	7	7	6	6	6	7.1
16	0.558720	4	4	4	4	4	4	4	4	4	4	4

Table A.2 Conversion time in s on a EC2 *t1.micro* server platform; the source file is of high audio quality, 192 Kbps. The individual conversions are labeled C1 to C10; \bar{T}_c is the mean conversion time.

Bit-rate (Kbps)	Audio file size (MB)	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	\bar{T}_c
128	4.467982	14	15	13	13	73	75	56	59	72	14	40.4
64	2.234304	9	9	9	32	44	9	23	9	45	10	19.9
32	1.117152	6	6	6	6	6	6	20	6	6	6	7.4
16	0.558720	6	6	6	6	6	6	20	6	6	6	5.1

The conversion from a high-quality audio file to a lower quality, thus a lower bit rate, is performed using the LAME library.

The conversion time depends on the desired bit-rate and the size of the original file. Tables A.1, A.2, A.3, and A.4 show the conversion time in s when the source MP3 file are of 320 Kbps and 192 Kbps, respectively; the size of the input files is also shown.

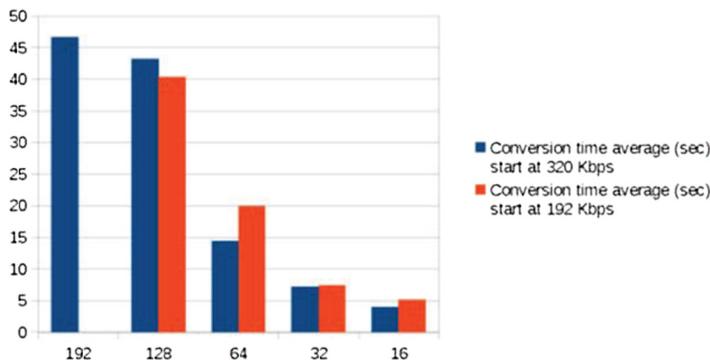
The platforms used for conversion are: (a) the EC2 *t1.micro* server for the measurements reported in Tables A.1 and A.2 and (b) the EC2 *c1.medium* for the measurements reported in Tables A.3 and A.4. The instances run the Ubuntu Linux operating system.

The results of our measurements when the instance is the *t1.micro* server exhibit a wide range of the conversion times, 13–80 s, for the large audio file of about 6.7 MB when we convert from 320 to 192 Kbps. A wide range, 13–64 s, is also observed for an audio file of about 4.5 MB when we convert from 320 to 128 Kbps. For low-quality audio, the file size is considerably smaller, about 0.56 MB, and the conversion time is constant and small, 4 s.

Fig. A.10 shows the average conversion time for the experiments summarized in Tables A.1 and A.2. It is somewhat surprising that the average conversion time is larger when the source file is smaller as it is in the case when the target bit rates are 64, 32 and 16 Kbps. Fig. A.11 shows the average conversion time for the experiments summarized in Tables A.3 and A.4.

The results of our measurements when the instance runs on the EC2 *c1.medium* platform show consistent and considerably lower conversion times; Fig. A.11 presents the average conversion time.

To understand the reasons for our results, we took a closer look at the two types of EC2 instances, “micro” and “medium,” and their suitability for the adaptive data-streaming service. The t1.micro sup-

**FIGURE A.10**

The average conversion time on a EC2 *t1.micro* platform. The left bars and the right bars correspond to the original file at: the highest resolution (320 Kbps data rate) and next highest resolution (192 Kbps data rate), respectively.

Table A.3 Conversion time T_c in seconds on a EC2 *c1.medium* platform; the source file is of the highest audio quality, 320 Kbps. The individual conversions are labeled C1 to C10; \bar{T}_c is the mean conversion time.

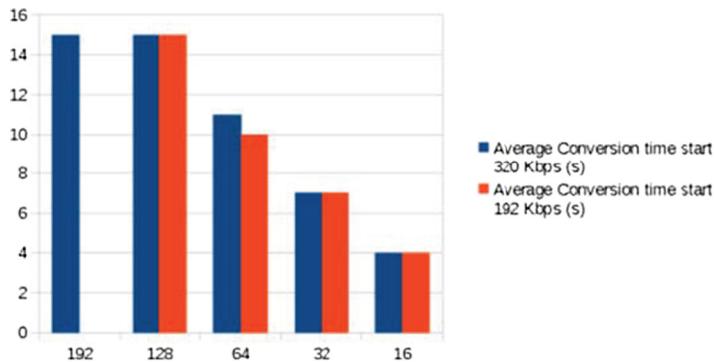
Bit-rate (Kbps)	Audio file size (MB)	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	\bar{T}_c
192	6.701974	15	15	15	15	15	15	15	15	15	15	15
128	4.467982	15	15	15	15	15	15	15	15	15	15	15
64	2.234304	11	11	11	11	11	11	11	11	11	11	11
32	1.117152	7	7	7	7	7	7	7	7	7	7	7
16	0.558720	4	4	4	4	4	4	4	4	4	4	4

Table A.4 Conversion time in s on a EC2 *c1.medium* platform; the source file is of high audio quality, 192 Kbps. The individual conversions are labeled C1 to C10; \bar{T}_c is the mean conversion time.

Bit-rate (Kbps)	Audio file size (MB)	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	\bar{T}_c
128	4.467982	15	15	15	15	15	15	15	15	15	15	15
64	2.234304	10	10	10	10	10	10	10	10	10	10	10
32	1.117152	7	7	7	7	7	7	7	7	7	7	7
16	0.558720	4	4	4	4	4	4	4	4	4	4	4

ports bursty applications, with a high average-to-peak ratio for CPU cycles, e.g., transaction processing systems. EBS provides block level storage volumes; the “micro” instances are only EBS-backed.

The “medium” instances support compute-intensive applications with a steady and relatively high demand for CPU cycles. Our application is compute-intensive, thus there should be no surprise that our measurements for the EC2 *c1.medium* platform show consistent and considerably lower conversion times.

**FIGURE A.11**

The average conversion time on a EC2 *c1.medium* platform. The left bars and the right bars correspond to the original file at: the highest resolution (320 Kbps data rate) and the next highest resolution (192 Kbps data rate), respectively.

A.5 Optimal FPGA synthesis

We now discuss another class of applications that could benefit from cloud computing. The benchmarks presented in Section 11.10 compared the performance of several codes running on a cloud with runs on supercomputers. As expected, the results showed that a cloud is not an optimal environment for applications exhibiting fine- or medium-grained parallelism. Indeed, communication latency is considerably larger on a cloud than on a supercomputer with a more expensive custom interconnect. This simply means that we have to identify applications that do not involve extensive communication, or applications exhibiting coarse-grained parallelism.

Computer clouds provide an ideal running environment for scientific applications involving model development when multiple cloud instances could concurrently run slightly different models of the system. When the model is described by a set of parameters, the application can be based on the SPMD paradigm combined with an analysis phase when the results from the multiple instances are ranked based on a well-defined metric.

In this case, there is no communication during the first phase of the application when partial results are produced and then written to storage server. The individual instances signal the completion, and a new instance to carry out the analysis and display the results is started. A similar strategy can be used by engineering applications of mechanical, civil, electrical, or electronic engineering, or any other system design area. In this case, the multiple instances run a concurrent design for different sets of parameters of the system.

A cloud application for optimal design of field-programmable gate arrays (FPGAs) is discussed next. As the name suggests, an FPGA is an integrated circuit designed to be configured/adapted/programmed in the field to perform a well-defined function [428]. Such a circuit consists of *logic blocks* and *interconnects* that can be “programmed” to carry out logical and/or combinatorial functions; see Fig. A.12.

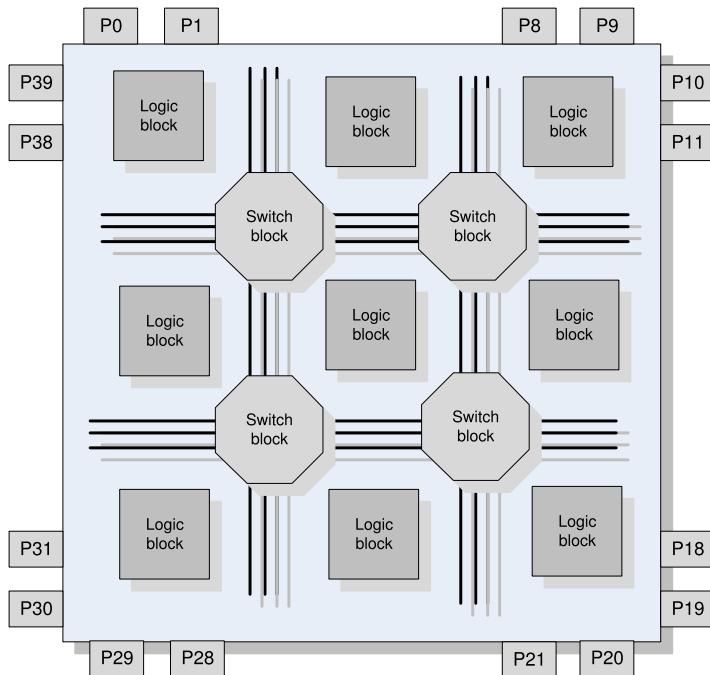


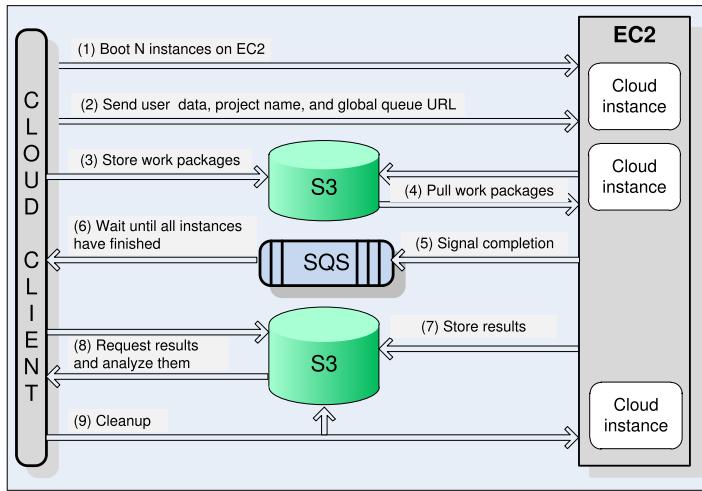
FIGURE A.12

The structure of a Field Programmable Gate Array (FPGA) with 30 pins, $P1-P29$, 9 logic blocks, and 4 switch-blocks.

The first commercially viable FPGA, XC2064, was produced in 1985 by Xilinx. Today, FPGAs are used in many areas, including digital signal processing, CRNs, aerospace, medical imaging, computer vision, speech recognition, cryptography, and computer hardware emulation. FPGAs are less energy efficient and slower than application-specific integrated circuits (ASICs). The widespread use of FPGAs is due to their flexibility and the ability to reprogram them.

Hardware description languages (HDLs) such as VHDL and Verilog are used to program FPGAs; HDLs are used to specify a register-transfer level (RTL) description of the circuit. Multiple stages are used to synthesize FPGA.

A cloud-based system was designed to optimize the routing and placement of components. The basic structure of the tool is shown in Fig. A.13. The system uses the PlanAhead tool from Xilinx—see <http://www.xilinx.com/>—to place system components and route chips on the FPGA logical fabric. The computations involved are fairly complex and take a considerable amount of time; for example, a fairly simple system consisting of a software core processor (Microblaze), a block random access memory (BRAM), and a couple of peripherals can take up to 40 min to synthesize on a powerful workstation. Running N design options in parallel on a cloud speeds up the optimization process by a factor close to N .

**FIGURE A.13**

The architecture of a cloud-based system to optimize the routing and placement of components on an FPGA.

A.6 Tensor network contraction on AWS

A numerical simulation project related to research in condensed matter physics is discussed in [332] and overviewed in this section. To illustrate the problems posed by Big Data, we analyze various options offered by 2016-vintage Amazon Web Services for running the application. *M4* and *C4* seem to be the best choices for applications such as Tensor Network Contraction (TNC).

Tensor contraction. In linear algebra, the *rank* \mathcal{R} of an object is given by the number of indices necessary to describe its elements. A scalar has rank 0, a vector $a = (a_1, a_2, \dots, a_n)$ has rank 1 and n elements, and a matrix $\mathcal{A} = [a_{ij}]$, $1 \leq i \leq n$, $1 \leq j \leq m$ has rank 2 and $n \times m$ elements

$$\mathcal{A} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \vdots & & & \\ a_{n1} & a_{n2} & \dots & a_{nm} \end{bmatrix}. \quad (\text{A.19})$$

Tensors have rank $\mathcal{R} \geq 3$; the description of tensor elements is harder. For example, consider a rank 3 tensor $\mathcal{B} = [b_{jkl}]$ with elements $b_{jkl} \in \mathbb{R}^{2 \times 2 \times 2}$. The eight elements of this tensor are: $\{b_{111}, b_{112}, b_{121}, b_{122}\}$ and $\{b_{211}, b_{212}, b_{221}, b_{222}\}$. We can visualize the tensor elements as the vertices of a cube where the first group of elements are in the plane $j = 1$ and $j = 2$, respectively. Similarly, the tensor elements $\{b_{111}, b_{211}, b_{112}, b_{212}\}$ and $\{b_{121}, b_{221}, b_{122}, b_{222}\}$ are in the planes $k = 1$ and $k = 2$, respectively, while $\{b_{111}, b_{121}, b_{211}, b_{221}\}$ and $\{b_{112}, b_{122}, b_{212}, b_{222}\}$ are in the planes $l = 1$ and $l = 2$, respectively.

Tensor contraction is the summation over repeated indices of the two tensors or of a vector and a tensor. Let \mathcal{C} be the contraction of two arbitrary tensors \mathcal{A} and \mathcal{B} . The rank of the tensor resulting after contraction is

$$\mathcal{R}(\mathcal{C}) = \mathcal{R}(\mathcal{A}) + \mathcal{R}(\mathcal{B}) - 2. \quad (\text{A.20})$$

For example, when $\mathcal{A} = [a_{ij}]$, $\mathcal{B} = [b_{jkl}]$ and we contract over j , we obtain $\mathcal{C} = [c_{ikl}]$ with

$$c_{ikl} = \sum_j a_{ij} b_{jkl}. \quad (\text{A.21})$$

The rank of \mathcal{C} is $\mathcal{R}(\mathcal{C}) = 2 + 3 - 2 = 3$. Tensor \mathcal{C} has eight elements:

$$\begin{array}{ll} c_{111} = \sum_{j=1}^2 a_{1j} b_{j11} = a_{11} b_{111} + a_{12} b_{211} & c_{121} = \sum_{j=1}^2 a_{1j} b_{j21} = a_{11} b_{121} + a_{12} b_{221} \\ c_{212} = \sum_{j=1}^2 a_{2j} b_{j12} = a_{21} b_{112} + a_{22} b_{212} & c_{222} = \sum_{j=1}^2 a_{2j} b_{j22} = a_{21} b_{122} + a_{22} b_{222} \\ c_{112} = \sum_{j=1}^2 a_{1j} b_{j12} = a_{11} b_{112} + a_{12} b_{212} & c_{122} = \sum_{j=1}^2 a_{1j} b_{j21} = a_{11} b_{122} + a_{12} b_{222} \\ c_{211} = \sum_{j=1}^2 a_{2j} b_{j11} = a_{21} b_{111} + a_{22} b_{211} & c_{221} = \sum_{j=1}^2 a_{2j} b_{j21} = a_{21} b_{121} + a_{22} b_{221} \end{array} \quad (\text{A.22})$$

Tensor networks and tensor network contraction (TNC). A *tensor network* is defined as follows: Let $[A_1], \dots, [A_n]$ be n tensors with index sets $x^{(1)}, \dots, x^{(n)}$ where each $\{x^{(i)}\}$ is a subset of $\{x_1, \dots, x_N\}$ with N very large. We assume that the “big” tensor $[A]_{\{x_1, \dots, x_N\}}$ can be expressed as the product of the “smaller” tensors $[A_1], \dots, [A_n]$:

$$[A]_{\{x_1, \dots, x_N\}} = [A_1]_{\{x^{(1)}\}} \dots [A_n]_{\{x^{(n)}\}}. \quad (\text{A.23})$$

We wish to compute the scalar

$$Z_A = \sum_{\{x_1, \dots, x_N\}} [A_i]_{\{x_1, \dots, x_N\}}. \quad (\text{A.24})$$

For example, $N = 7$ and $n = 4$ and the index set is $\{x_1, x_2, \dots, x_7\}$ for the TNC in Fig. A.14. The four “small” tensors and their respective subsets of the index set are

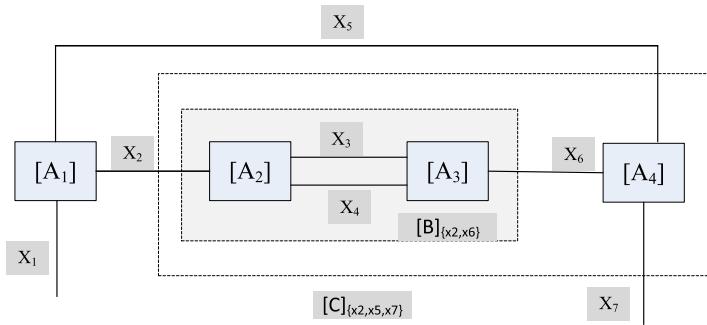
$$[A_1]_{\{x_1, x_2, x_5\}}, [A_2]_{\{x_2, x_3, x_4\}}, [A_3]_{\{x_3, x_4, x_6\}} \text{ and } [A_4]_{\{x_5, x_6, x_7\}}. \quad (\text{A.25})$$

The “big” tensor $[A]$ is the product of the four “small” tensors

$$[A]_{\{x_1, x_2, x_3, x_4, x_5, x_6, x_7\}} = [A_1]_{\{x_1, x_2, x_5\}} \otimes [A_2]_{\{x_2, x_3, x_4\}} \otimes [A_3]_{\{x_3, x_4, x_6\}} \otimes [A_4]_{\{x_5, x_6, x_7\}}. \quad (\text{A.26})$$

To calculate Z_A , we first contract $[A_2]$ and $[A_3]$, and the result is tensor $[B]$

$$[B]_{\{x_2, x_6\}} = \sum_{x_3} \sum_{x_4} [A_2]_{\{x_2, x_3, x_4\}} \otimes [A_3]_{\{x_3, x_4, x_6\}}. \quad (\text{A.27})$$

**FIGURE A.14**

The ordering of tensor contraction when the index set is $\{x_1, x_2, \dots, x_7\}$ and the four “small” tensors are $[A_1]_{\{x_1, x_2, x_5\}}$, $[A_2]_{\{x_2, x_3, x_4\}}$, $[A_3]_{\{x_3, x_4, x_6\}}$ and $[A_4]_{\{x_5, x_6, x_7\}}$. Tensors $[B]$ and $[C]$ are the results of contraction of $[A_2]$, $[A_3]$, and $[A_4]$, $[B]$, respectively.

Next, we contract $[B]$ and $[A_4]$ to produce $[C]$

$$[C]_{\{x_2, x_5, x_7\}} = \sum_{x_6} [B]_{\{x_2, x_6\}} \otimes [A_4]_{\{x_5, x_6, x_7\}}. \quad (\text{A.28})$$

Finally, we compute

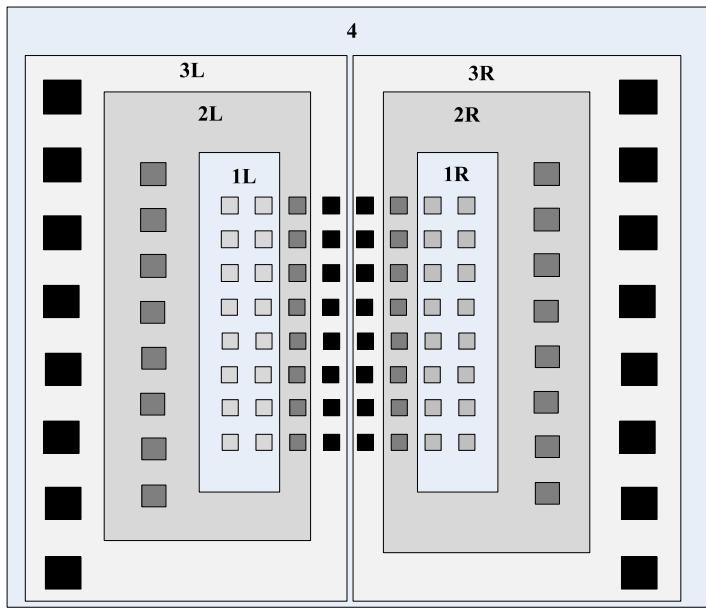
$$Z_{\{x_1, x_7\}} = \sum_{x_2} \sum_{x_5} [A_1]_{\{x_1, x_2, x_5\}} \otimes [C]_{\{x_2, x_5, x_7\}}. \quad (\text{A.29})$$

Tensor network contraction is CPU- and memory-intensive. If the tensor network has an arbitrary topology, TNC is considerably more intensive than in the case of a regular topology, e.g., a 2-D lattice.

An example of TNC. We now discuss the case of an application where the tensors form a 2-D, $L \times L$ rectangular lattice. Each tensor in the interior of the lattice has four indices, each one running from 1 to D^2 , while outer tensors have only three indices, and the ones at the corners have only two. The resulting tensors form a product of vectors (top and bottom tensors) and matrices (interior tensors) with vertical orbitals running from 1 to D^{2L} . The space required for tensor network contraction can be very large: We expect parameter values as large as $D = 20$ and $L = 100$. This is a Big Data application; 20^{200} is a very large number indeed!!

Fig. A.15 illustrates the *generic TLC algorithm* for $L = 8$. The first iteration of the computation contracts the left-most (1L) and the right-most columns (1R) of the tensor network. The process continues until we end up with the “big” vector after $L/2 = 4$ iterations. The left and right contractions, (1L, 2L, and 3L) and (1R, 2R, and 3R), are mirror images of one another and are carried out concurrently.

TNC algorithm for condensed matter physics. In quantum mechanics, vectors in an n -dimensional Hilbert space describe quantum states, and tensors describe transformations of quantum states. Tensor network contraction has applications in condensed matter physics, and our discussion is focused only on the algorithmic aspects of the problem.

**FIGURE A.15**

Contraction when $L = 8$ and we have an 8×8 tensor network. The first iteration contracts columns 1 and 2 and columns 7 and 8; see 1L and 1R boxes. During the second iteration, the two resulting tensors are contracted with columns 3 and 6, respectively, as shown by 2L and 2R boxes. During the 3rd iteration, the new tensors are contracted with columns 4 and 5, respectively, as shown by the 3L and 3R boxes. Finally, during the 4th iteration, the “big” vector is obtained by contracting the two tensors produced by the 3rd iteration.

The algorithm for tensor network contraction should be flexible, efficient, and cost effective. Flexibility means the ability to run problems of various sizes, with a range of values for D and L parameters. An efficient algorithm should support effective multithreading and optimal use of available system resources.

The notations used to describe the contraction algorithms for tensor network T are:

- $N^{(i)}$ —number of vCPUs for iteration i ; $N = 2$ for all iterations of Stage 1, while the number of vCPUs for Stage 2 may increase with the number of iterations;
- m —the amount of memory available on the vCPU of the current instance;
- $\mathcal{T}^{(i)} = [\mathcal{T}_{j,k}^{(i)}]$ —version of the $[T]$ after the i th iteration;
- $L^{(i)}$ —the number of columns of $[T]$ at iteration i ;
- $\mathcal{T}_{j,k}^{(i)}$ —tensor in row j and column k of $\mathcal{T}^{(i)}$; $T_{j,k}^1 = T_{j,k}$;
- $\mathcal{T}_k^{(i)}$ —column k of $\mathcal{T}^{(i)}$;
- $\mathbb{C}(\mathcal{T}_k^{(i)}, T_j)$, $i > 1$ —contraction operator applied to columns $\mathcal{T}_k^{(i)}$ and T_j in Stage 1;
- $\mathbb{V}(\mathcal{T}^{(L_{col})})$ —vertical contraction operator applied to the “big tensor” obtained after L_{col} column contractions;

- μ —amount of memory for $T_{j,k}$, a tensor of the original T ;
- $\mu^{(i)}$ —storage for a tensor $\mathcal{T}_{j,k}^{(i)}$ created at iteration i ;
- I_{max} —maximum number of iterations for Stage 1 of the TNC algorithm.

The *generic contraction algorithm* for a 2-D tensor network, with L_{row} rows and L_{col} columns, $T = [T_{j,k}]$, $1 \leq j \leq 2L_{row}$, $1 \leq k \leq 2L_{col}$, is an extension of the one in Fig. A.15. TNC is an iterative process: At each iteration, two pairs of columns are contracted concurrently. During the first iteration the two pairs of columns of T , $(1, 2)$, and $(2L, 2L - 1)$ are contracted. At iterations $2 \leq i \leq L$, the new tensor network has $L^{(i)} = L - 2i$ columns, and the contraction is applied to column pairs: the column resulting from the contractions at iteration $(i - 1)$, now columns 1 and $L^{(i)}$ with columns 2 and $L^{(i)} - 1$, respectively.

The TNC algorithm is organized in multiple stages with different AWS instances for different stages. Small- to medium-sized problems need only the first stage to produce the results and use low end instances with 2 vCPUs, while large problems must feed the results of the first stage to a second one running on more powerful AWS instances. A third stage may be required for extreme cases when the size of one tensor exceeds the amount of vCPU memory, some 4 GB at this time. The three stages of the algorithm are discussed next.

- *Stage 1.* An entire column of $\mathcal{T}^{(i)}$ can be stored in the vCPU memory and successive contraction iterations can proceed seamlessly when

$$L_{row} \left(2\mu + \mu^{(i-1)} + \mu^{(i)} \right) < m. \quad (\text{A.30})$$

This is feasible for the first iterations of the algorithm and for relatively small values of L_{row} . Call I_{max} the largest value of i which satisfies the Eq. (A.30).

Use a low-end M4 or C4 instance with 2 vCPUs, $N = 2$. The computation runs very fast with optimal use of the secondary storage and network bandwidth. Each vCPU is multithreaded: Multiple threads carry out the operations required by the contraction operator \mathbb{C} , while one thread reads the next column of the original tensor network, \mathcal{T} in preparation for the next iteration.

- *Stage 2.* After a number of iterations, the condition in Eq. (A.30) is no longer satisfied, and the second phase should start. Now an iteration consists of partial contractions when subsets of column tensors are contracted independently. In this case the number of vCPUs is $N > 2$.
- *Stage 3.* As the amount of space needed for a single tensor increases and the vCPU memory cannot store a single tensor,

$$\mu_i > m. \quad (\text{A.31})$$

In this extreme case we use several instances with the largest number of vCPUs, e.g., either M4.10xlarge or C4.10xlarge.

Stage 1 TNC algorithm. The algorithm is a straightforward implementation of the generic TNC algorithm:

1. Start an instance with $N = 2$, e.g., C4.large;
2. Read input parameters e.g., L_{row} , L_{col} ;
3. Compute I_{max} ;

4. First iteration
 - a. vCP1—read T_1 and T_2 , apply $\mathbb{C}(T_1, T_2)$; start reading T_3 ;
 - b. vCP2—read $T_{L_{col}}$ and $T_{L_{col}-1}$, apply $\mathbb{C}(T_{L_{col}}, T_{L_{col}-1})$, start reading $T_{L_{col}-2}$;
5. Iterations $2 \leq i \leq \min[I_{max}, L_{col}]$. The column numbers correspond to the contracted tensor network with $L_{col}^{(i)} = L_{col} - 2(i - 1)$ columns
 - a. vCP1—apply $\mathbb{C}(\mathcal{T}_1^{(i)}, T_2)$; start reading T_3 ;
 - b. vCP2—apply $\mathbb{C}(\mathcal{T}_{L_{col}^{(i)}}^{(i)}, T_{L_{col}^{(i)}-1})$; start reading $T_{L_{col}^{(i)}-2}$;
6. If $L_{col} \leq I_{max}$, carry out vertical compression of the “big tensor” and finish;
 - a. Apply $\mathbb{V}(T^{(L_{col})})$;
 - b. Write result;
 - c. Kill the instance;
7. Else, prepare the data for the Stage 2 algorithm;
 - a. vCPU1—save $\mathcal{T}_i^{(i)}$;
 - b. vCPU2—save $\mathcal{T}_{L_{col}^{(i)}}^{(i)}$;
 - c. Kill the instance.

Stage 2 TNC algorithm. This stage starts with a tensor network $\mathcal{T}^{(I_{max})}$ with $2(L_{col} - I_{max})$ columns and L_{row} rows. Multiple partial contractions will be done for each column of $\mathcal{T}^{(I_{max})}$ during this stage.

The number of vCPUs for the instance used for successive iterations may increase. Results of a partial iteration have to be saved at the end of the partial iteration. The parameters for this phase are:

- $\mu_{pc}^{(i)}$ —the space per tensor required for partial contraction at iteration i

$$\mu_{pc}^{(i)} = \mu + \mu^{(i-1)} + \mu^{(i)}; \quad (\text{A.32})$$

partial contraction increases the space required by each tensor;

- $\mathbb{C}_{pc}(\mathcal{T}_k^{(i)}, T_j, s)$ —partial contraction operator applied to segment s of columns $\mathcal{T}_k^{(i)}$ and T_j in Stage 2;
- $n_r^{(i)}$ —number of rows of a column segment for each partial contraction at iteration i given by

$$n_r^{(i)} = \left\lceil \frac{m}{\mu_{pc}^{(i)}} \right\rceil. \quad (\text{A.33})$$

- $p^{(i)}$ —number of partial contractions per column at iteration i ;

$$p^{(i)} = \left\lceil \frac{L_{row}}{n_r^{(i)}} \right\rceil. \quad (\text{A.34})$$

The total number of partial contractions at iteration i is $2p^{(i)}$;

- The number of vCPUs for iteration i is

$$N^{(i)} = 2p^{(i)}. \quad (\text{A.35})$$

- $L_{col}^{(i)}$ —the number of columns at iteration i of Stage 2;
- $I_{Max} = L_{col} - I_{max}$ —the number of iterations of Stage 2 assuming that Stage 3 is not necessary;
- $\mathbb{A}_{pc}(\mathcal{T}_{k,p^{(i)}}^i)$ —assembly operator for the $p^{(i)}$ segments resulting from partial contraction of column k at iteration i .

Stage 2 TNC consists of the following steps:

1. For $i = 1, I_{Max}$
 - a. Compute $\mu_{pc}, n_r^{(i)}, p^{(i)}, N^{(i)}$;
 - b. If $N \leq 40$, start an instance with $N = N^{(i)}$; else, start multiple C4.10xlarge instances to run concurrently all partial contractions;
 - c. For $j = 1, p^{(i)}$
 - i. $vCPU_j$
 - Read $\mathcal{T}_{1,j}^{(i)}$ and $T_{2,j}$ and apply $\mathbb{C}_{pc}(\mathcal{T}_{1,j}^{(i)}, T_{2,j})$,
 - Store the result $\mathcal{T}_{1,j}^{(i+1)}$;
 - ii. $vCPU_{j+p^{(i)}}$
 - Read $\mathcal{T}_{L_{col}^{(i)},j}^{(i)}$ and $T_{L_{col}^{(i)}-1,j}$ and apply $\mathbb{C}_{pc}(\mathcal{T}_{L_{col}^{(i)},j}^{(i)}, T_{L_{col}^{(i)}-1,j})$,
 - Store the result, $\mathcal{T}_{L_{col}^{(i)},j}^{(i+1)}$;
 - d. Assemble partial contractions
 - i. $vCPU_1$
 - Apply $\mathbb{A}_{pc}(\mathcal{T}_1^i, p^{(i)})$,
 - Store $\mathcal{T}_1^{(i+1)}$;
 - ii. $vCPU_2$
 - Apply $\mathbb{A}_{pc}(\mathcal{T}_{L_{col}^{(i)}}^i, p^{(i)})$,
 - Store $\mathcal{T}_{L_{col}^{(i+1)}}^{(i+1)}$.
2. If $i < I_{Max}$, proceed to the next iteration, $i = i + 1$; else,
 - a. Apply $\mathbb{V}\mathcal{T}^{(I_{Max})}$;
 - b. Write TNC result;
 - c. Kill the instance.

Stage 3 TNC algorithm. The algorithm is similar with the one for Stage 2, but now a single tensor is distributed to multiple vCPUs.

An analysis of the memory requirements for TNC. Let us assume that we have L tensors per column, and each tensor has dimension D . Consider the leftmost, or equivalently the rightmost column, and note that the number of bonds differs for different tensors; the top and the bottom tensors have *two* bonds,

and the other $L - 1$ have *three* bonds, so the total number of elements in this column is

$$\mathcal{N}_1^{(0)} = 2D^2 + (L - 2)D^3. \quad (\text{A.36})$$

The top and bottom tensors of the next column have three bonds, and the remaining $L - 2$ have four bonds; thus, the total number of elements in the second column is

$$\mathcal{N}_2^{(0)} = 2D^3 + (L - 2)D^4. \quad (\text{A.37})$$

After contraction, the number of elements becomes

$$\mathcal{N}_1^{(1)} = 2D^3 + (L - 2)D^5. \quad (\text{A.38})$$

Each tensor element requires two double-precision floating-point numbers, thus the amount of memory needed for the first iteration is

$$\begin{aligned} \mathcal{M}^{(1)} &= 2 \times 8 \times [D^2 + (L - 2)D^3 + 2D^3 + (L - 2)D^4 + 2D^3 + (L - 2)D^5] \\ &= 16 \times [2D^3 + (L - 2)D^4 + 2D^2(1 + D) + (L - 2)D^3(1 + D^2)]. \end{aligned} \quad (\text{A.39})$$

The amount of memory needed for iterations 2 and 3 are

$$\begin{aligned} \mathcal{M}^{(2)} &= 16 \times [2D^3 + (L - 2)D^5 + 2D^3 + (L - 2)D^4 + 2D^4 + (L - 2)D^7] \\ &= 16 \times [2D^3 + (L - 2)D^4 + 2D^3(1 + D) + (L - 2)D^5(1 + D^2)] \end{aligned} \quad (\text{A.40})$$

and

$$\begin{aligned} \mathcal{M}^{(3)} &= 16 \times [2D^4 + (L - 2)D^7 + 2D^3 + (L - 2)D^4 + 2D^5 + (L - 2)D^9] \\ &= 16 \times [2D^3 + (L - 2)D^4] + 2D^4(1 + D) + (L - 2)D^7(1 + D^2). \end{aligned} \quad (\text{A.41})$$

It follows that the amount of memory for iteration i is

$$\mathcal{M}^{(i)} = 16 \times [2D^3 + (L - 2)D^4 + 2D^{i+1}(1 + D) + (L - 2)D^{2i+1}(1 + D^2)]. \quad (\text{A.42})$$

When $D = 20$ and $L = 100$, the amount of memory for the first iteration is

$$16 \times [2 \times 20^3 + 98 \times 20^4 + 2 \times 20^2 \times (1 + 20) + 98 \times (20^3 + 20^5)] = 5\,281\,548\,800 \text{ bytes}. \quad (\text{A.43})$$

This example shows why only the most powerful systems with ample resources can be used for TNC. It also shows that an application has to adapt, the best it can, to the packages of resources provided by the CSP, while in a better world, an application-centric view should prevail, and the system should assemble and offer precisely the resources needed by an application, neither more nor less.

A.7 A simulation study of machine-learning scalability

Training Convolutional Neural Networks (CNNs) and Deep Neural Networks (DNNs) is computationally intensive. In Section 13.2, we discussed results reported by Google researchers showing that

**FIGURE A.16**

StarCraft is a real-time strategy game played against one opponent. Each player has hundreds of game pieces.

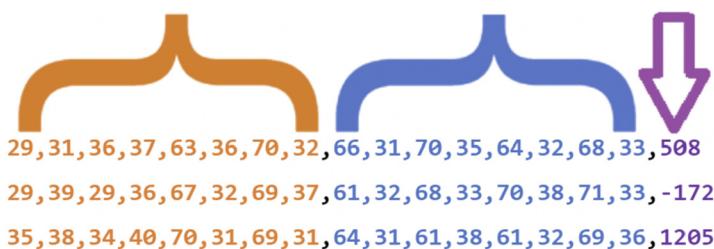
training a CNN with 89 layers, 100×10^6 weights, and 1750 operations/weight required 17.5×10^9 operations. The training set size for activity recognition for an eight-layer DNN required 10^6 videos, and the training time was reported to take some 30 days.

We posed to students in a graduate cloud computing class the question if it is feasible to use machine learning for predicting the computational effort required to train a DNN. A CNN used to predict the computational effort for DNN training should have as input data describing the DNN and data describing the computational effort for training a range of DNNs. CNN output should be the training time, and the number of arithmetic operations required for DNN training.

DNN data should include: hidden layer count; number of neurons in each hidden layer; number of neurons in input and output layers; size of the training set; learning rate; number of epochs; and the loss function determining how far a given output is from correct value. System monitoring data includes training time; number of different arithmetic operations; CPU utilization; number of messages exchanged; number of I/O operations; etc.

A preliminary study conducted in Spring 2020 by Thomas C. investigated the scalability of ML training function of DNN characteristics. In each experiment only the computation time was recorded. The project expanded on work done for a machine learning class to control game pieces in a *StarCraft* game; see Fig. A.16. The ML project had three phases: (i) build a dataset by running a specific game scenario 20 000 times; (ii) train the model using to predict a “best” action; and (iii) rerun game scenarios using the new predicted *best* action.

The new project was to simulate scaling up the workload to 10, 100, or 1 000 machines. The project used several AWS services including EC2, S3, SimpleDB, Elastic Beanstalk, SNS, Lambda, and API gateway to create a system of multiple AWS instances for running ML workloads and then to test scaling using a combination of real and simulated workloads. Each training dataset was created by running *StarCraft* with game pieces in a certain positions and included the input location of game pieces for the player and the opponent and an integer reflecting the battle score—a high battle score

**FIGURE A.17**

Each record included: (i) player units; (ii) opponent units; and (iii) battle score.

meaning that the battle went well for the player and a low or negative score if the battle did not go well for the player. The data set size reflected the number of records in the training data set; see Fig. A.17; the project used datasets with 4k, 6k, 8k, and 10k records.

Other training parameters are: (i) batch size, i.e., how many records to train on before updating the weights (64, 128, 256); (ii) the number of epochs, i.e., how many times the dataset is used to train the neural network (125, 250, 500, 1 000); and (iii) the loss function.

System organization is shown in Fig. A.18. *Model workload clients* perform several functions: (i) train models using seven training parameters: data set size, learning rate, hidden layer count, hidden layer size, batch size, number of epochs, and loss function; (ii) get training parameters using AWS HTTP Endpoint; and (iii) upload results including training time, model loss, and model accuracy. *Data set workload clients* generate and upload datasets to S3 buckets and upload the results.

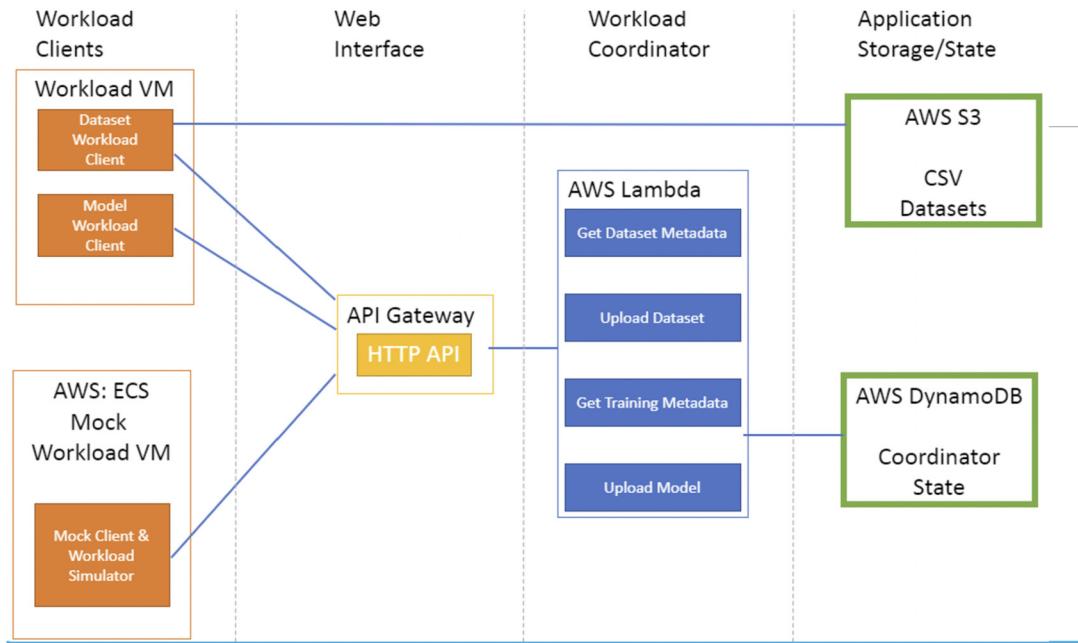
The *workload coordinator* uses the AWS Lambda service to get training parameters and data set metadata and then upload the model and the datasets. It also stores three DynamoDB tables, the coordinator state, models, and datasets.

Simulation experiments created four datasets with a total of 40 000 records and trained 5 000 models. The average training time was 143.7 s with a standard deviation of 250.2 s. Fig. A.19 shows a histogram of learning models function of training time and learning rate, respectively. As expected, the training time increases nearly linearly with the number of hidden layers and with the size of hidden layers; see Fig. A.20.

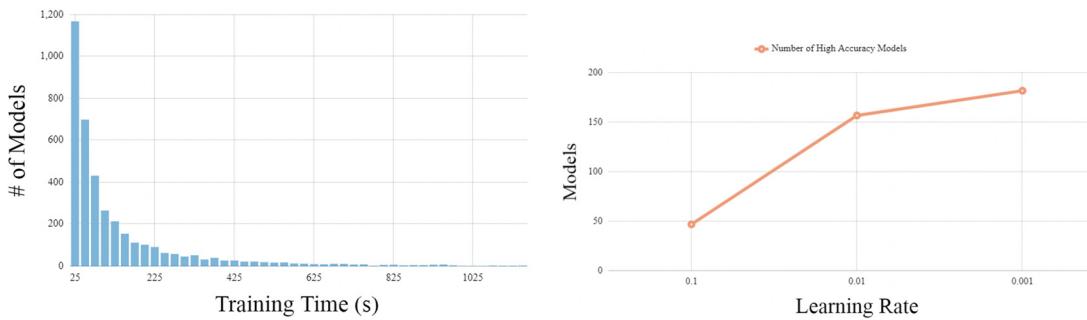
Training time decreases as the batch size increases; see Fig. A.21 (left); accuracy increases as the number of epochs increase; see Fig. A.21 (right). The results of a scaling experiment using a mock workload are summarized in Fig. A.22. The number of agents, i.e., servers involved increases linearly, and the number of requests increases faster than linearly as the time progresses until the experiment was stopped when reaching the Lambda free-tier limit.

A.8 Cloud-based task alert application

TaskAlert application developed by Brandon G. is designed to manage a variety of user tasks, some involving computer-based activities and tasks related to any other types of professional or personal activities. The application allows users on client systems access to a set of microservices performing

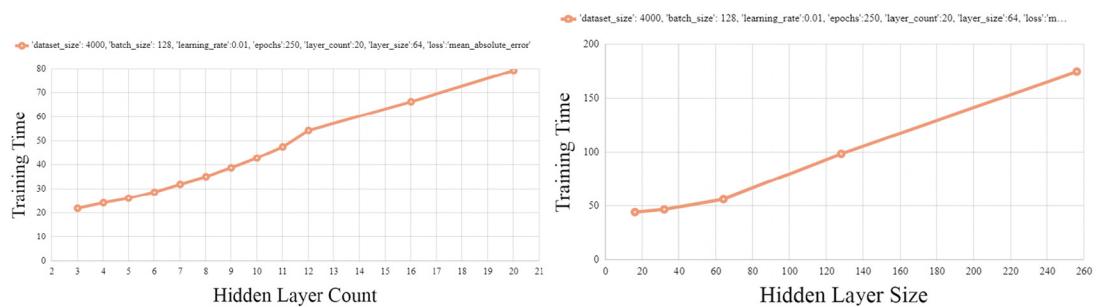
**FIGURE A.18**

System organization. There are two types of clients, *data set workload* to create datasets and *model workload* to train models. A *workload coordinator* assigns workload metadata and stores results. The *mock workload client* simulates workload clients.

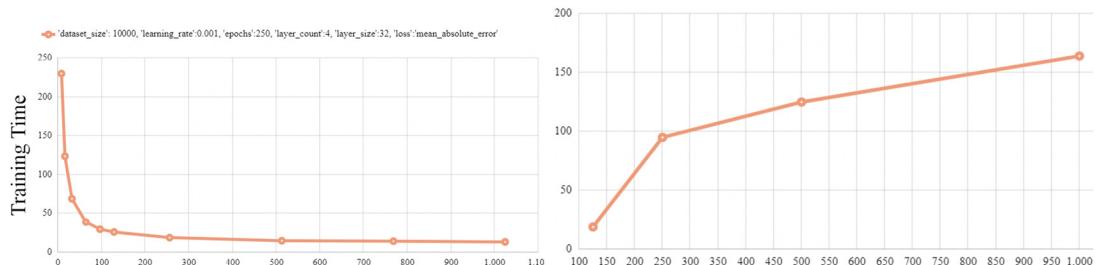
**FIGURE A.19**

Number of models function of training time and learning rate.

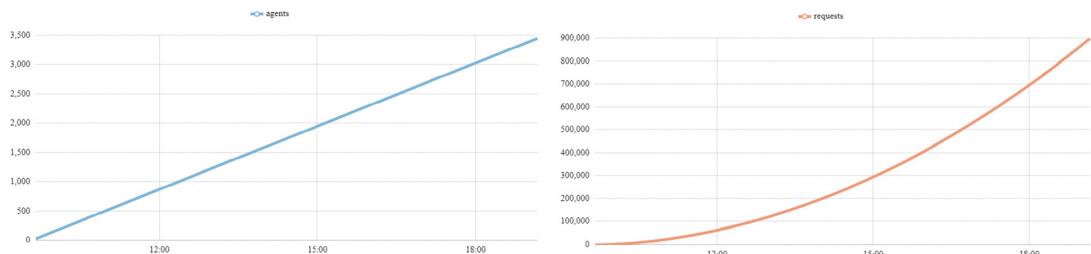
functions such as: (i) access a list of tasks to be completed; (ii) manage account information and register for the service; (iii) supply task information, to automatically create a task in response to an email sent to the application.

**FIGURE A.20**

Training time function of hidden layer (a) count; (b) size.

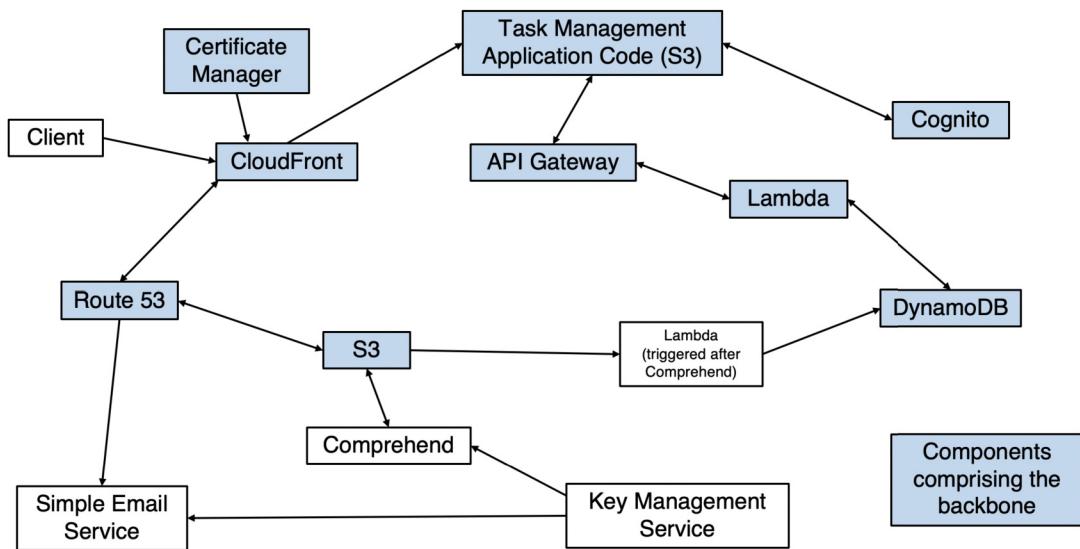
**FIGURE A.21**

Training time function of batch size and accuracy function of number of epochs.

**FIGURE A.22**

Scaling experiment using a mock workload. The experiment started at 9 am and was stopped after some 10 h after reaching AWS Lambda free-tier limit when 3 500 agents and 900 000 requests had been processed without any errors.

Task description identifies the task type, resources needed for task completion, deadlines, milestones, task priority, task dependencies, events and/or conditions that could affect task completions, interactions with other individuals required for task completion, and possible conflicts with other tasks.

**FIGURE A.23**

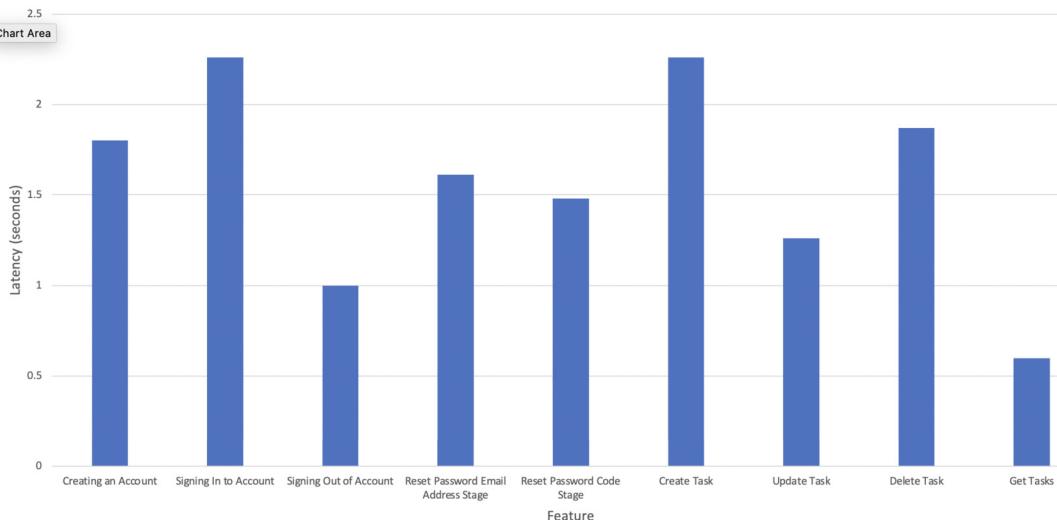
TaskAlert organization and AWS services used by the application.

Devices used to communicate with *TaskAlert* and the communication means; email, messaging, and phone alerts are also part of the task description.

Application front end hosts a user interface (UI) allowing authorized and authenticated users to interact with the back end; see Fig. A.23. The front end is implemented in Javascript with Node.js and React, a JavaScript library used to create interactive UIs. Several other libraries are used by the application front end: (i) *React-router-dom* library used to route users to the proper pages; and (ii) *Axios* library used to access custom REST API to interact with other AWS services. The *AWS Amplify* framework is used to access *Cognito* API to support user accounts. The application uses extensively the *Lambda* service, which runs user code without the need to provision or manage servers. Moreover, the user pays only for the compute time consumed by the application; it charges for every 100 ms when the code is triggered.

TaskAlert's back end supports task creation, deletion, modification, task retrieval, and maintaining user accounts using *Cognito*. Other AWS services used by *TaskAlert* are:

1. *Route 53*—allows users to access the service through *mytaskalert.com* domain.
2. *CloudFront*—converts user HTTP requests to HTTPS requests to *mytaskalert.com*.
3. *Certificate Manager*—issues a TLS/SSL certificate for *mytaskalert.com*.
4. *S3*—hosts the application code and dependences for *mytaskalert.com*
5. *API Gateway*—used by the application code to trigger the appropriate *Lambda* functions.
6. *Lambda*—performs functions specified by API Gateway to create, modify, delete, and retrieve user tasks.
7. *DynamoDB*—stores user tasks and allows task modification, deletion, and retrieval.

**FIGURE A.24**

TaskAlert execution profile. The time in seconds for various task management activities.

8. *Cognito*—allows users to sign in through social identity providers.

Fig. A.24 shows the time for various task management activities. The response time for task creation and signing into an existing account is about two s, while the response time for other task management activities is around one second.

A.9 Cloud-based health-monitoring application

The cloud-based health-monitoring application, developed by Mojtaba T., is designed to monitor outpatients who need constant supervision, individuals with chronic health problems, high-risk individuals, and frail elderly people. The system uses smartphones and wearable sensors to collect monitoring data from patients, as shown in Fig. A.25.

Existing health-monitoring systems seldom provide extended physiological data and a fast response time in case of emergencies. The cloud-based system takes advantage of a wide range of sensors and high-bandwidth wireless technology to connect patients with healthcare providers, emergency services, pharmacies, rehabilitation centers, and other relevant facilities. The system enables healthcare providers to exchange information, medical records and test results and to consult with each other for the benefit of a patient; it also provides access to healthcare databases.

Physiological and behavioral data is collected from a range of sources including: personal electrocardiogram (ECG) devices, electroencephalograph (EEG) headsets, photoplethysmographs (PPG), accelerometers, gyroscopes, and gait-cycle monitors. Photoplethysmography (PPG) is a optical technique used to detect volumetric changes in blood in peripheral circulation. A gait cycle is the time

**FIGURE A.25**

Physiological and behavioral data is collected from a range of sources.

period or sequence of events or movements from the instant one foot contacts the ground until the foot contacts the ground again and involves propulsion of the center of gravity in the direction of motion.

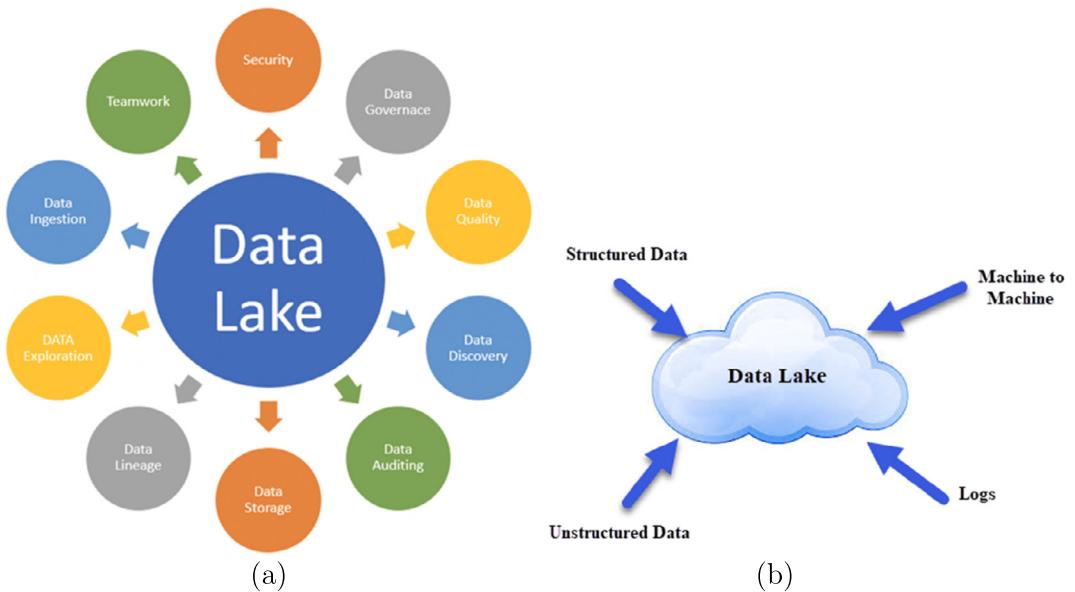
Some data is collected periodically, e.g., on hourly or daily bases, and others only when sensors signal out-of-range measurements. The vital signs of outpatients and other data such as gait cycles are streamed continually. It follows that the system must be nimble and in some cases provide real-time response.

The main functions of the system are: (i) user authentication; (ii) access control; (iii) data synchronization across multiple devices; (iv) patient behavior analysis; (v) data storage management; (vi) data collection from various sources; (vii) real-time data collection and response.

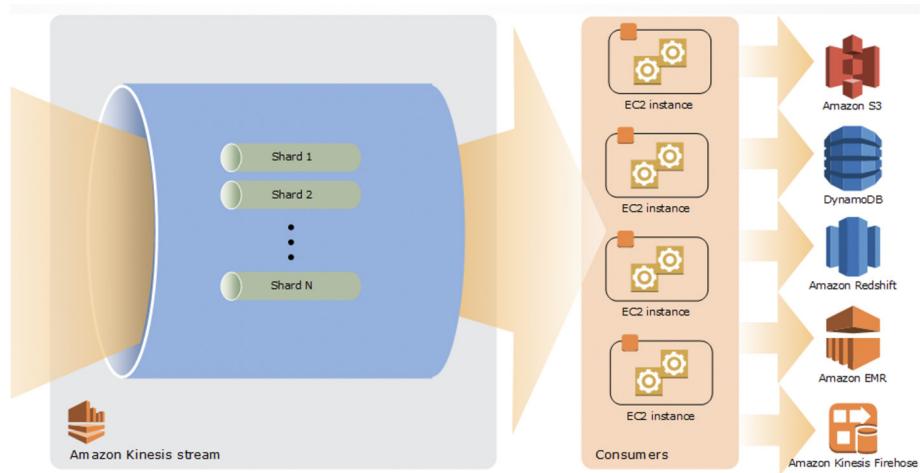
The application needs a more agile and flexible storage and uses AWS Lake Formation to set up a secure data lake. A data lake is a centralized, curated, and secured repository storing structured and unstructured data at any scale; see Fig. A.26(a). Data can be stored as is, without having first to be structured, and can be used for multiple types of analytics from dashboards and visualizations to big data processing, real-time analytics, and machine learning; see Fig. A.26(b).

The applications uses AWS Kinesis system for data streaming from multiple sources to multiple EC2 instances and then stores the results of processing kinesis streams to data lakes, as shown in Fig. A.27. Amazon Kinesis service supports process streaming data at any scale, along with the flexibility to choose the optimal processing for every data stream, including real-time data provided by various health-monitoring sensors.

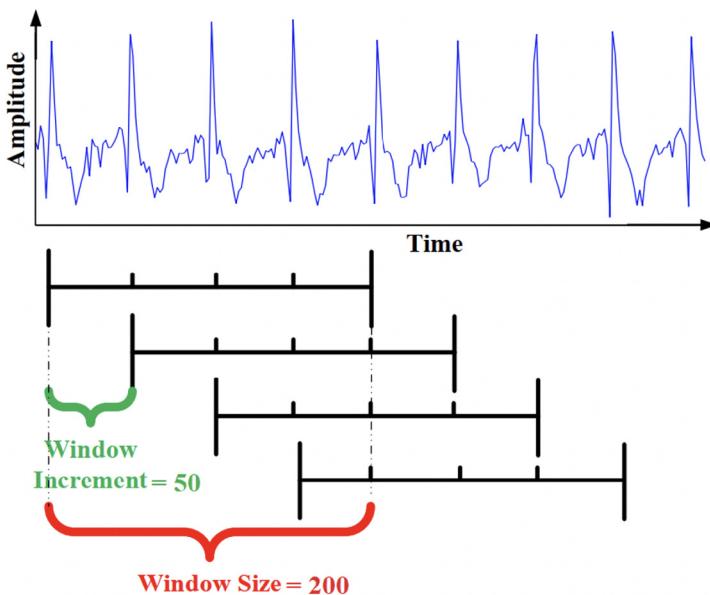
The application is designed to collect, preprocess, analyze, and store Big Data produced and consumed by a large user population. To accomplish its mission, the application must rely on a powerful framework, and it was decided to use Spark. Recall from Section 4.12 that Apache Spark is an open-source Big Data processing framework used to manage Big Data processing for various data types.

**FIGURE A.26**

(a) Multiple functionality supported by data lakes. (b) Various types of data can be stored in data lakes.

**FIGURE A.27**

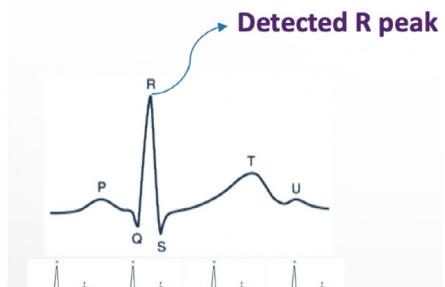
The flow of data from multiple sources through the AWS Kinesis to EC2 instances for processing and to data lakes for storing measurement results and raw data.

**FIGURE A.28**

EEG preprocessing. A band-pass filter for 0.5–4 Hz removes low-frequency power-line noise and high-frequency noise and segments the EEG signal to create 200 data-point windows overlapping each other on 50 data-point window.

	Feature	Unit	Description
Time – Domain	HRV		Heart Rate Variability
	AvgHR	bpm	Average Heart Rate
	MeanRR	ms	Mean of selected R-R series
	NN50	count	Number of consecutive R-R intervals that differs more than 50 millisecond
	SD_HR	1/min	Standard Deviation of Heart Rate
	SD_RR	1/min	Standard Deviation of R-R interval
	RMSSD	ms	Root Mean Square of the differences of selected R-R interval series
Frequency- Domain	SE	-	Sample Entropy
	PSE	-	Power Spectral Entropy

(a)



b)

FIGURE A.29

(a) EEG features in time-domain and in frequency-domain. (b) R-peak detection.

Along with MapReduce, it also supports SQL queries, streaming data, graph data processing, and machine learning techniques.

To illustrate the complexity of data analysis, we show the preprocessing and feature extraction for EEG in Figs. A.28 and A.29. ECG peak detection is used to detect the combination of Q, R, and S waves or the so-called QRS complex. The QRS complex is a peak model for ECG signal including Q-valley point, R-peak point, and S-valley point. Other important peak points in the ECG signal are the P-peak point and the T-peak point. The detection of the QRS complex is critical for ECG signal processing.

A band-pass filter for 0.5–4 Hz removes low-frequency powerline noise and high-frequency noise and segments the EEG signal to create 200 data-point windows overlapping each other on 50 data-point window. Fig. A.29(a) shows the EEG features. For each segment of the data, various features are extracted and results are saved in a file. Fig. A.29(b) illustrates the R-peak detection.

Cloud application development

B

In the previous chapters, our discussion was focused on research issues in cloud computing; now, we examine computer clouds from the perspective of an application developer. This chapter presents a few recipes useful to assemble a cloud computing environment on a local system and to use basic cloud functions.

It is fair to assume that the population of application developers and cloud users is, and will continue to be, very diverse. Some cloud users have developed and run parallel applications on clusters or other types of systems for many years and expect an easy transition to the cloud. Others are less experienced, but willing to learn and expect a smooth learning curve. Many view cloud computing as an opportunity to develop new businesses with minimum investment in computing equipment and human resources.

The questions we address are: How easy is it to use the cloud? How knowledgeable should an application developer be about networking and security? How easy is it to port an existing application to the cloud? How easy is it to develop a new cloud application?

The answers to these questions are different for the three cloud delivery models, SaaS, PaaS, and IaaS; the level of difficulty increases as we move towards the base of the cloud service pyramid, as shown in Fig. B.1. Recall that SaaS applications are designed for the end-users and are accessed over the web; in this case, the user must be familiar with the API of a particular application. PaaS provides a set of tools and services designed to facilitate application coding and deploying, while IaaS provides the hardware and the software for servers, storage, networks, including operating systems and storage management software. We restrict our discussion to the IaaS cloud computing model, and we concentrate on popular services offered by AWS.

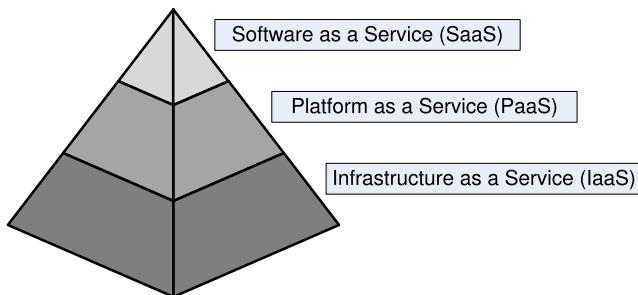
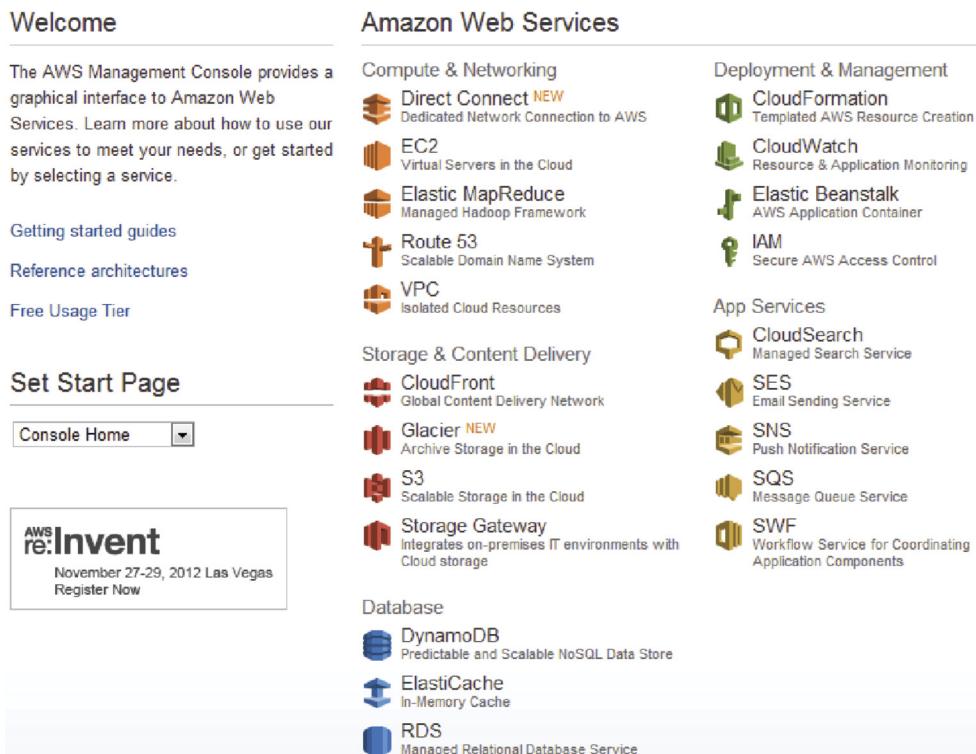


FIGURE B.1

A pyramid model of cloud computing paradigms; the infrastructure provides the basic resources; the platform adds an environment to facilitate the use of these resources, while software allows direct access to services.

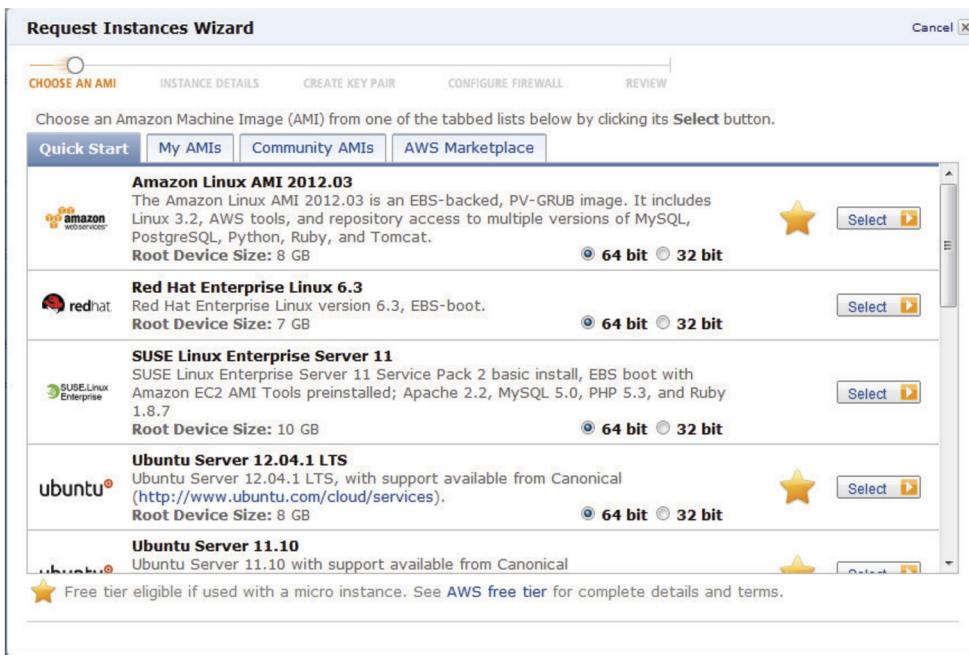
**FIGURE B.2**

Amazon Web Services accessible from the Amazon Management Console.

Though the AWS services are well documented, the environment they provide for cloud computing requires some effort to benefit from the full spectrum of services offered. In this section, we report on lessons learned from the experience of a group of students with a strong background in programming, networking, and operating systems; each one of them was asked to develop a cloud application for a problem of interest in their own research area. First, we discuss several issues related to cloud security, a major stumbling block for many cloud users; then, we present a few recipes for the development of cloud applications, and finally, we analyze several cloud applications developed by individuals in this group over a period of less than three months.

B.1 AWS EC2 instances

Fig. B.2 displays the Amazon Management Console window listing the Amazon Web Services offered at the time of this writing; the services are grouped in several categories: computing and networking, storage and content delivery, deployment and management, databases, and application services.

**FIGURE B.3**

The Instance menu allows the user to select from existing AMIs.

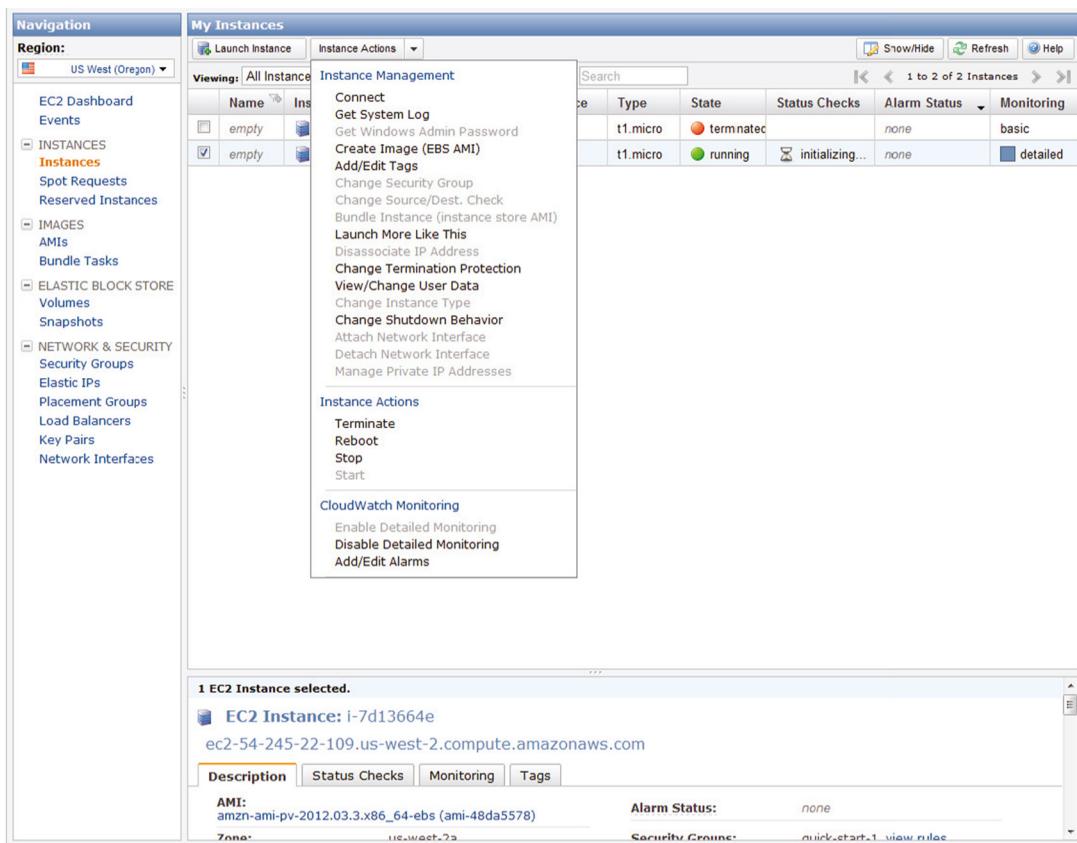
In spite of the wealth of information available from the providers of cloud services, the learning curve of an application developer is still relatively steep. The examples discussed in this chapter are designed to help overcome some of the hurdles faced when someone first attempts to use the AWS. Due to space limitations, we have chosen to cover only a few of the very large number of combinations of services, operating systems, and programming environments supported by AWS.

In Section 2.2, we mentioned that new services are continually added to AWS; the look and feel of the web pages change over time. The screen shots reflect the state of the system at the time of the writing of the first edition of the book, i.e., the second half of 2012.

To access AWS, one must first create an account at <http://aws.amazon.com/>; once the account has been created, the Amazon Management Console (AMC) allows the user to select one of the service, e.g., EC2, and then start an instance.

Recall that an EC2 instance is a virtual server started in a region and availability zone selected by the user. Instances are grouped into a few classes, and each class has a specific amount of resources, such as CPU cycles, main memory, secondary storage, communication, and I/O bandwidth available to it. Several operating systems are supported by AWS including: Amazon Linux, Red Hat Enterprise Linux 6.3, SUSE Linux Enterprise Server 11, Ubuntu Server 12.04.1, and several versions of Microsoft Windows, see Fig. B.3.

The next step is to create an AMI (Amazon Machine Image) on one of the platforms supported by AWS and start an instance using the *RunInstance* API. An AMI is a unit of deployment, an environment

**FIGURE B.4**

The *Instance Action* pull-down menu of the *Instances* panel of the *AWS Management Console* enables the user to interact with an instance, e.g., *Connect* and *Create an EBS AMI Image*, among others.

including all information necessary to set up and boot an instance. If an application needs more than 20 instances, then a special form must be filled out. The local instance persists in storage only for the duration of an instance. The data persist when an instance is started using the Amazon EBS (Elastic Block Storage), and then the instance can be restarted at a later time.

Once an instance has been created, the user can perform several actions, for example, connect to the instance, launch more instances identical to the current one, or create an EBS AMI. The user can also terminate, reboot, or stop the instance; see Fig. B.4. The *Network & Security* panel allows the creation of *Security Groups*, *Elastic IP addresses*, *Placement Groups*, *Load Balancers* and *Key Pairs* (see the discussion in Section B.3), while the *EBS* panel allows the specification of volumes and the creation of snapshots.

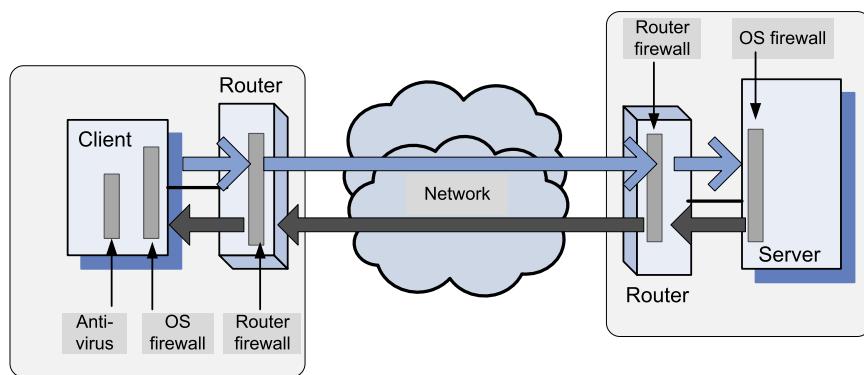


FIGURE B.5

Firewalls screen incoming and sometimes outgoing traffic. The first obstacle encountered by the inbound or outbound traffic is a router firewall, and the next one is the firewall provided by the host operating system; sometimes, the antivirus software provides a third line of defense.

B.2 Connecting clients to cloud instances through firewalls

A firewall is a software system based on a set of rules for filtering network traffic; its function is to protect a computer in a local area network from unauthorized access. The first generation of firewalls, deployed in the late 1980s, carried out *packet filtering*; they discarded individual packets that did not match a set of acceptance rules. Such firewalls operated below the transport layer and discarded packets based on the information in the headers of physical, data link, and transport-layer protocols.

The second generation of firewalls operate at the transport layer and maintain the state of all connections passing through them. Unfortunately, this traffic filtering solution opened the possibility of *denial-of-service attacks*; a denial-of-service (DOS) attack targets a widely used network service and forces the operating system of the host to fill the connection tables with illegitimate entries. DOS attacks prevent legitimate access to the service.

The third generation of firewalls “understand” widely-used application layer protocols, such as FTP, HTTP, TELNET, SSH, and DNS. These firewalls examine the header of application layer protocols and support *intrusion detection systems*.

Firewalls screen incoming traffic and, sometimes, filter outgoing traffic as well. A first filter encountered by the incoming traffic in a typical network is a firewall provided by the operating system of the router; the second filter is a firewall provided by the operating system running on the local computer; see Fig. B.5.

Typically, the local area network (LAN) of an organization is connected to the Internet via a router; a router firewall often hides the true address of hosts in the local network using the network address translation (NAT) mechanism. The hosts behind a firewall are assigned addresses in a “private address range,” and the router uses the NAT tables to filter the incoming traffic and translate external IP addresses to private ones. The mapping between the pair (*external address, external port*) and the (*internal address, internal port*) tuple carried by the network address translation function of the router firewall is called a *pinhole*.

Table B.1 Firewall rule setting. The columns indicate if a feature is supported or not by an operating system: the second column—a single rule can be issued to accept/reject a default policy; the third and fourth columns—filtering based on IP destination and source address, respectively; the fifth and sixth columns—filtering based on TCP/UDP destination and source ports, respectively; the seventh and eighth columns—filtering based on Ethernet MAC destination and source address, respectively; the ninth and tenth columns—inbound (ingress) and outbound (egress) firewalls, respectively.

Operating system	Def rule	IP dest addr	IP src addr	TCP/UDP dest port	TCP/UDP src port	Ether MAC dest	Ether MAC src	In-bound fwall	Out-bound fwall
Linux iptables	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
OpenBSD	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Windows XP	No	No	Yes	Partial	No	No	No	Yes	No
Cisco Access List	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Juniper Networks	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes

If one tests a client-server application with the client and the server in the same local area network, the packets do not cross a router; once a client from a different LAN attempts to use the service, the packets may be discarded by the firewall of the router. The application may no longer work if the router is not properly configured.

Table B.1 summarizes the options supported by various operating systems running on a host or on a router.

A *rule* specifies a filtering option at: (i) the network layer, when filtering is based on the destination/source IP address; (ii) the transport layer, when filtering is based on destination/ source port number; (iii) the MAC layer, when filtering is based on the destination/source MAC address.

In Linux or Unix systems the firewall can be configured only by someone with a *root* access using the *sudo* command. The firewall is controlled by a kernel data structure, the *iptables*. The *iptables* command is used to set up, maintain, and inspect the tables of the *IPv4* packet filter rules in the Linux kernel. Several tables may be defined; each table contains a number of built-in chains and may also contain user-defined chains.

A *chain* is a list of rules that can match a set of packets: the *INPUT* rule controls all incoming connections; the *FORWARD* rule controls all packets passing through this host; and the *OUTPUT* rule controls all outgoing connections from the host. A *rule* specifies what to do with a packet that matches: *Accept*—let the packet pass; *Drop*—discharge the packet; *Queue*—pass the packet to the user space; *Return*—stop traversing this chain and resume processing at the head of the next chain. For complete information on the *iptables*, see <http://linux.die.net/man/8/iptables>.

To get the status of the firewall, specify the L (List) action of the *iptables* command

```
sudo iptables -L
```

As a result of this command, the status of the *INPUT*, *FORWARD*, and *OUTPUT* chains will be displayed.

To change the default behavior for the entire chain, specify the action P (Policy), the chain name, and target name; e.g., to allow all outgoing traffic to pass unfiltered use

```
sudo iptables -P OUTPUT ACCEPT
```

To add a new security rule specify: the action, A (add), the chain, the transport protocol, *tcp* or *udp*, and the target ports, as in:

```
sudo iptables -A INPUT -p -tcp -dport ssh -j ACCEPT
sudo iptables -A OUTPUT -p -udp -dport 4321 -j ACCEPT
sudo iptables -A FORWARD -p -tcp -dport 80 -j DROP
```

To delete a specific security rule from a chain, set the action D (Delete) and specify the chain name and the rule number for that chain; the top rule in a chain has number 1:

```
sudo iptables -D INPUT 1
sudo iptables -D OUTPUT 1
sudo iptables -D FORWARD 1
```

By default, the Linux virtual machines on Amazon's EC2 accept all incoming connections.

The ability to access that virtual machine will be permanently lost when a user accesses an EC2 virtual machine using ssh and then issues the following command

```
sudo iptables -P INPUT DROP.
```

The access to a Windows firewall is provided by a GUI accessed as follows:

Control Panel -> System & Security -> Windows Firewall -> Advanced Settings

The default behavior for incoming and/or outgoing connections can be displayed and changed from the window *Windows Firewall with Advanced Security on Local Computer*.

The access to the Windows XP firewall is provided by a GUI accessed by selecting *Windows Firewall* in the *Control Panel*. If the status is *ON*, incoming traffic is blocked by default, and a list of Exceptions (as noted on the *Exceptions* tab) define the connections allowed. The user can only define exceptions for: *tcp* on a given port, *udp* on a given port, and a specific program. *Windows XP* does not provide any control over outgoing connections.

Antivirus software running on a local host may provide an additional line of defense. For example, the Avast antivirus software (see www.avast.com) supports several real-time shields. The *Avast network shield* monitors all incoming traffic; it also blocks access to known malicious websites. The *Avast web shield* scans the HTTP traffic and monitors all web-browsing activities. The antivirus also provides statistics related to its monitoring activities.

B.3 Security rules for application- and transport-layer protocols in EC2

A client must know the IP address of a virtual machine in the cloud to be able to connect to it. Domain Name Service (DNS) is used to map human-friendly names of computer systems to IP addresses on the Internet or in private networks. DNS is a hierarchical distributed database and plays a role reminiscent of an Internet phone book.

In 2010, Amazon announced a DNS service called Route 53 to route users to AWS services and to infrastructure outside of AWS. A network of DNS servers scattered across the globe enables customers to gain reliable access to AWS and places strict controls over who can manage their DNS system by allowing integration with AWS Identity and Access Management (IAM).

For several reasons, including security and the ability of the infrastructure to scale up, the IP addresses of instances visible to the outside world are mapped internally to private IP addresses. A virtual machine running under Amazon's EC2 has several IP addresses:

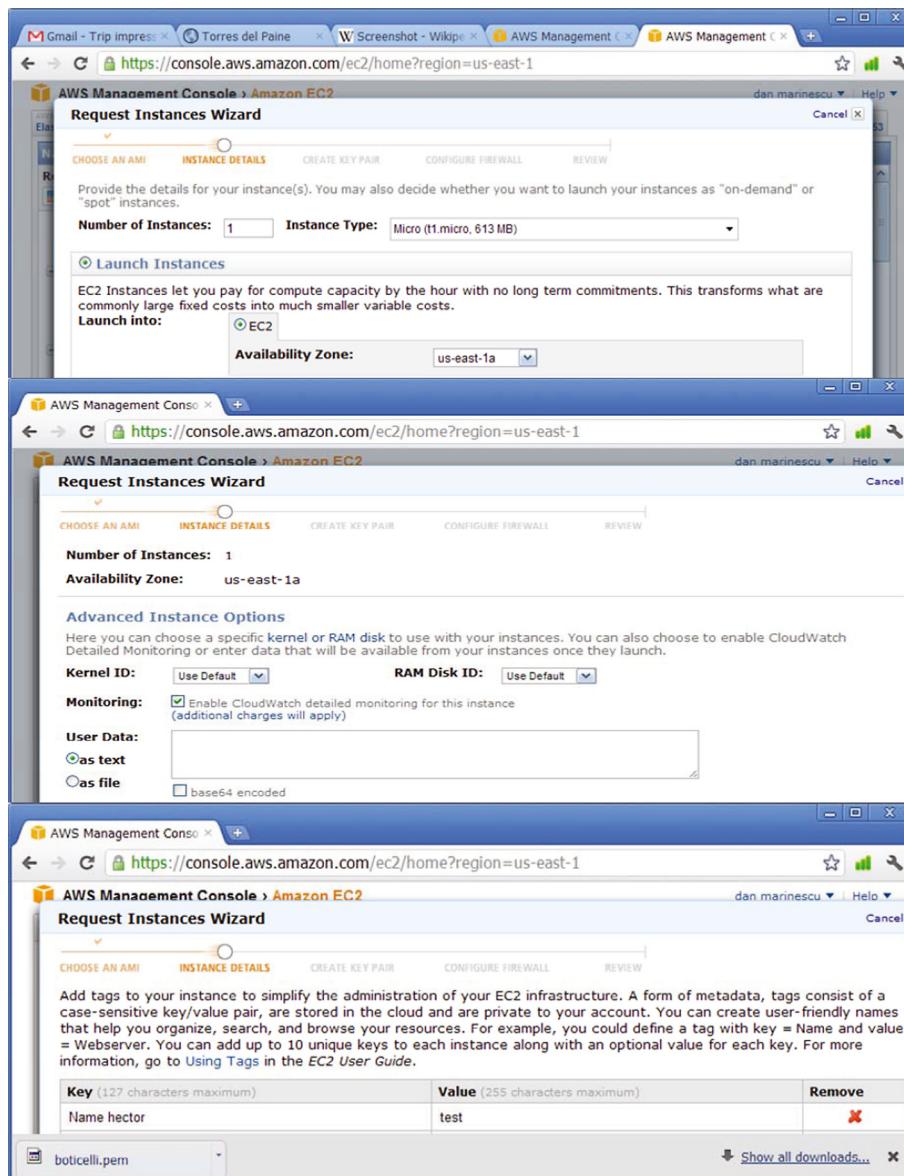
1. *EC2 Private IP Address*: The internal address of an instance; it is only used for routing within the EC2 cloud.
2. *EC2 Public IP Address*: Network traffic originating outside the AWS network must use either the public IP address or the elastic IP address of the instance. The public IP address is translated using the Network Address Translation (NAT) to the private IP address when an instance is launched and it is valid until the instance is terminated. Traffic to the public address is forwarded to the private IP address of the instance.
3. *EC2 Elastic IP Address*: The IP address allocated to an AWS account and used by traffic originated from outside AWS. NAT is used to map an elastic IP address to the private IP address. Elastic IP addresses allow the cloud user to mask instance or availability zone failures by programmatically remapping a public IP addresses to any instance associated with the user's account. This allows fast recovery after a system failure; for example, rather than waiting for a cloud maintenance team to reconfigure or replace the failing host, or waiting for DNS to propagate the new public IP to all of the customers of a web service hosted by EC2, the web-service provider can remap the elastic IP address to a replacement instance. Amazon charges a fee for unallocated Elastic IP addresses.

To control access to a user's VMs, AWS uses *security groups*. A VM instance belongs to one, and only one, security group that can only be defined before the instance is launched. Once an instance is running, the security group the instance belongs to cannot be changed. However, more than one instance can belong to a single security group.

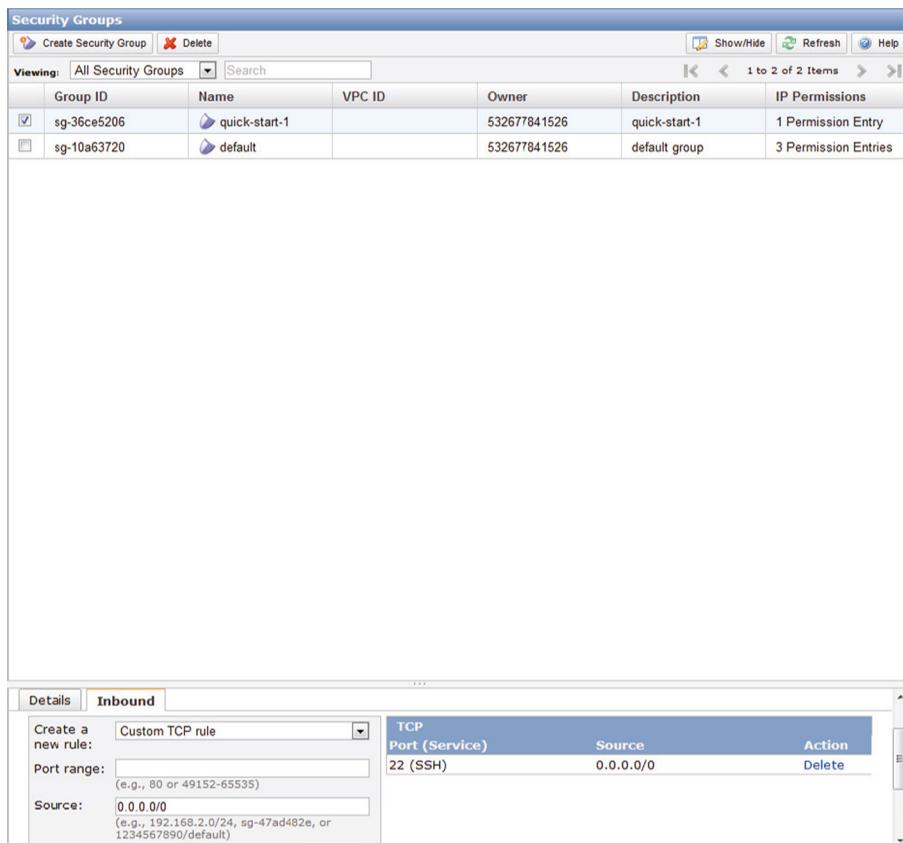
Security group rules control inbound traffic to the instance and have no effect on outbound traffic from the instance. The inbound traffic to an instance, either from outside the cloud or from other instances running on the cloud, is blocked, unless a rule stating otherwise is added to the security group of the instance. For example, assume a client running on instance A in the security group Σ_A is to connect to a server on instance B listening on TCP port P, where B is in security group Σ_B . A new rule must be added to security group Σ_B to allow connections to port P; to accept responses from server B, a new rule must be added to security group Σ_A .

The following steps allow the user to add a security rule:

1. Sign in to the AWS Management Console at <http://aws.amazon.com> using your email address and password and select EC2 service.
2. Use the *EC2 Request Instance Wizard* to specify the instance type, whether it should be monitored, and specify a key/value pair for the instance to help organize and search; see Fig. B.6.
3. Provide a name for the key pair; then, on the left hand side panel, choose *Security Groups* under *Network & Security*, select the desired security group, and click on the *Inbound* tab to enter the desired rule; see Fig. B.7.

**FIGURE B.6**

EC2 Request Instances Wizard is used to: (a) specify the number and type of instances and the zone; (b) specify the kernelId and the RAM diskId and enable the *CloudWatch* service to monitor the EC2 instance; (c) add tags to the instance; a tag is stored in the cloud and consists of a case-sensitive key/value pair private to the account.

**FIGURE B.7**

AWS security. Choose *Security Groups* under *Network & Security*, select the desired security group, and click on the *Inbound* tab to enter the desired rule.

To allocate an elastic IP address, use the *Elastic IPs* tab of the *Network & Security* left-hand side panel.

On Linux or Unix systems the port numbers below 1024 can only be assigned by the *root*. The plain ASCII file called *services* maps friendly textual names for internet services to their assigned port numbers and protocol types as in the following example:

```
netstat 15/tcp
ftp      21/udp
ssh      22/tcp
telnet   23/tcp
http     80/tcp
```

B.4 How to launch an EC2 Linux instance and connect to it

This section summarizes the step-by-step process to launch an EC2 Linux instance from a Linux platform.

A. Launch an instance

1. From the *AWS management console*, select EC2 and, once signed in, go to *Launch Instance Tab*.
2. To determine the processor architecture when you want to match the instance with the hardware, enter the command

```
uname -m
```

and choose an appropriate Amazon Linux AMI by pressing *Select*.

3. Choose *Instance Details* to control the number, size, and other settings for instances.
4. To learn how the system works, press *Continue* to select the default settings.
5. Define the instances security, as discussed in Section B.3: in the *Create Key Pair* page enter a name for the pair, and then press *Create and Download Key Pair*.
6. The key pair file downloaded in the previous step is a *.pem* file, and it must be hidden to prevent unauthorized access; if the file is in the directory *awmdir/dada.pem*, enter the commands

```
cd awmdir  
chmod 400 dada.pem
```

7. Configure the firewall; go to the page *Configure Firewall*, select the option *Create a New Security Group*, and provide a *Group Name*. Normally, one uses *ssh* to communicate with the instance; the default port for communication is port 8080, and one can change the port and other rules by creating a new rule.
8. Press *Continue* and examine the review page that gives a summary of the instance.
9. Press *Launch* and examine the confirmation page, and then press *Close* to end the examination of the confirmation page.
10. Press the *Instances* tab on the navigation pane to view the instance.
11. Look for your *Public DNS* name. As by default, some details of the instance are hidden; click on the *Show/Hide* tab on the top of the console and select *Public DNS*.
12. Record the *Public DNS* as *PublicDNSname*; it is needed to connect to the instance from the Linux terminal.
13. Use the *ElasticIP* panel to assign an elastic IP address if a permanent IP address is required.

B. Connect to the instance using ssh and the tcp transport protocol

1. Add a rule to the *iptables* to allow ssh traffic using the tcp protocol. Without this step, either an *access denied* or a *permission denied* error message appears when trying to connect to the instance:

```
sudo iptables -A iptables -p -tcp -dport ssh -j ACCEPT
```

2. Enter the Linux command

```
ssh -i abc.pem ec2-user@PublicDNSname
```

If you get the prompt *You want to continue connecting?*, respond *Yes*; a warning that the DNS name was added to the list of known hosts will appear.

3. An icon of the *Amazon Linux AMI* will be displayed.

C. Gain root access to the instance

By default the user does not have *root* access to the instance and thus cannot install any software. Once connected to the EC2 instance, use the following command to gain *root* privileges

```
sudo -i
```

Then, use *yum* install commands to install software, e.g., *gcc* to compile C programs on the cloud.

D. Run the service *ServiceName*

If the instance runs under *Linux* or *Unix*, the service is terminated when the *ssh* connection is closed; to avoid the early termination, use the command

```
nohup ServiceName
```

To run the service in the background and redirect *stdout* and *stderr* to files *p.out* and *p.err*, respectively, execute the command

```
nohup ServiceName > p.out 2 > p.err &
```

B.5 How to use S3 in Java

The Java API for Amazon Web Services is provided by the AWS SDK. A software development kit (SDK) is a set of software tools for the creation of applications in a specific software environment. Java Development Kit (JDK) is an SDK for Java developers available from Oracle.

JDK includes a set of programming tools, such as: *javac*, the Java compiler that converts Java source code into Java bytecode; *java*, the loader for Java applications—it can interpret the class files generated by the Java compiler; *javadoc* the documentation generator; *jar*, the archiver for class libraries; *jdb*, the debugger; *JConsole*, the monitoring and management console; *jstat*, JVM statistics monitoring; *jps*, JVM process status tool; *jinfo*, the utility to get configuration information from a running Java process; *jrunscript*, the command-line script shell for Java; *appletviewer* tool to debug Java applets without a web browser; and *idlj*, the IDL-to-Java compiler. The *Java Runtime Environment* is also a component of the JDK consisting of a Java Virtual Machine (JVM) and libraries.

Create an S3 client. S3 access is handled by the class *AmazonS3Client* instantiated with the account credentials of the AWS user

```
AmazonS3Client s3 = new AmazonS3Client(  
    new BasicAWSCredentials("your_access_key", "your_secret_key"));
```

The access and the secret keys can be found on the user's AWS account home page, as mentioned in Section B.3.

Buckets. An *S3 bucket* is analogous to a file folder or directory and it is used to store *S3 Objects*. Bucket names must be *globally unique*; hence it is advisable to check first if the name exists

```
s3.doesBucketExist("bucket_name");
```

This function returns “true” if the name exists and “false” otherwise. Buckets can be created and deleted either directly from the AWS Management Console or programmatically as follows:

```
s3.createBucket("bucket_name");
s3.deleteBucket("bucket_name");
```

S3 objects. An *S3 object* stores the actual data, and it is indexed by a key string. A single key points to only one S3 object in one bucket. Key names do not have to be globally unique, but if an existing key is assigned to a new object, then the original object indexed by the key is lost. To upload an object into a bucket, one can use the *AWS Management Console*, or programmatically a file *local_file_name* can be uploaded from the local machine to the bucket *bucket_name* under the key *key* using

```
File f = new File("local_file_name");
s3.putObject("bucket_name", "key", f);
```

A versioning feature for the objects in S3 was made available recently; it allows to preserve, retrieve, and restore every version of an S3 object. To avoid problems when uploading large files, e.g., the drop of the connection, use the *.initiateMultipartUpload()* with an API described at the *AmazonS3Client*. To access this object with key *key* from the bucket *bucket_name*, use:

```
S3Object myFile = s3.getObject("bucket_name", "key");
```

To read this file, you must use the S3Object's *InputStream*:

```
InputStream in = myFile.getObjectContent();
```

The *InputStream* can be accessed using *Scanner*, *BufferedReader*, or any other supported method. Amazon recommends closing the stream as early as possible because the content is not buffered, and it is streamed directly from the S3; an open *InputStream* means an open connection to S3. For example, the following code will read an entire object and print the contents to the screen:

```
AmazonS3Client s3 = new AmazonS3Client(
    new BasicAWSCredentials("access_key", "secret_key"));
InputStream input = s3.getObject("bucket_name", "key")
    .getObjectContent();
Scanner in = new Scanner(input);
while (in.hasNextLine())
{
    System.out.println(in.nextLine());
}
in.close();
input.close();
```

Batch Upload/Download. Batch upload requires repeated calls of *s3.putObject()* while iterating over local files.

To view the keys of all objects in a specific bucket, use

```
ObjectListing listing = s3.listObjects("bucket_name");
```

Object Listing supports several useful methods including `getObjectSummaries()`. `S3ObjectSummary` encapsulates most of an S3 object properties (excluding the actual data), including the key to access the object directly,

```
List<S3ObjectSummary> summaries = listing.getObjectSummaries();
```

For example, the following code will create a list of all keys used in a particular bucket, and all of the keys will be available in string form in `List < String > allKeys`:

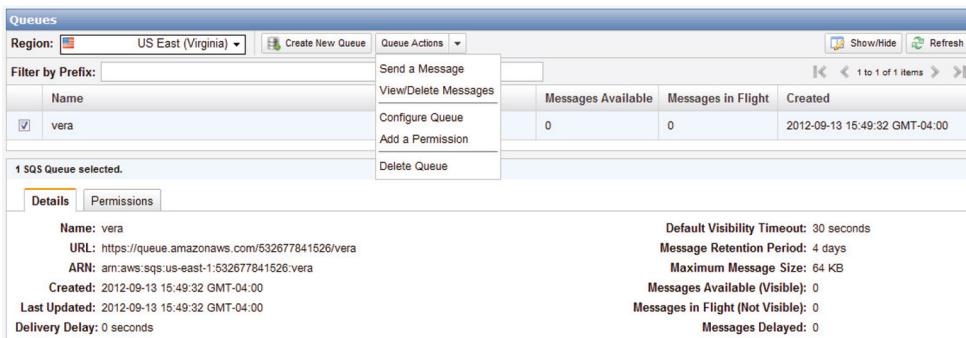
```
AmazonS3Client s3 = new AmazonS3Client(
    new BasicAWSCredentials("access_key", "secret_key"));
List<String> allKeys = new ArrayList<String>();
ObjectListing listing = s3.listObjects("bucket_name");
for (S3ObjectSummary summary:listing.getObjectSummaries())
{
    allKeys.add(summary.getKey());
}
```

Note that, if the bucket contains a very large number of objects, then `s3.listObjects()` will return a truncated list. Use the following command to test if the list is truncated one could use `listing.isTruncated()`; to get the next batch of objects

```
s3.listNextBatchOfObjects(listing);
```

To account for a large number of objects in the bucket, the previous example becomes

```
AmazonS3Client s3 = new AmazonS3Client(
    new BasicAWSCredentials("access_key", "secret_key"));
List<String> allKeys = new ArrayList<String>();
ObjectListing listing = s3.listObjects("bucket_name");
while (true)
{
    for (S3ObjectSummary summary :
        listing.getObjectSummaries())
    {
        allKeys.add(summary.getKey());
    }
    if (!listing.isTruncated())
    {
        break;
    }
    listing = s3.listNextBatchOfObjects(listing);
}
```

**FIGURE B.8**

Queue actions in SQS.

B.6 How to manage AWS SQS services in C#

Recall from Section 2.2 that SQS is a system for supporting automated workflows. Multiple components can communicate with messages sent and received via SQS. An example showing the use of message queues is presented in Section 11.6. Fig. B.8 shows the actions available for a given queue in *SQS*.

The following steps can be used to create a queue, send a message, receive a message, delete a message, and delete the queue in C#:

1. Authenticate an SQS connection

```
NameValueCollection appConfig =
    ConfigurationManager.AppSettings;
AmazonSQS sqs = AWSClientFactory.CreateAmazonSQSClient
    (appConfig["AWSAccessKey"], appConfig["AWSSecretKey"]);
```

2. Create a queue

```
CreateQueueRequest sqsRequest = new CreateQueueRequest();
sqsRequest.QueueName = "MyQueue";
CreateQueueResponse createQueueResponse =
    sqs.CreateQueue(sqsRequest);
String myQueueUrl;
myQueueUrl = createQueueResponse.CreateQueueResult.QueueUrl;
```

3. Send a message

```
SendMessageRequest sendMessageRequest =
    new SendMessageRequest();
sendMessageRequest.QueueUrl =
    myQueueUrl; //URL from initial queue
sendMessageRequest.MessageBody = "This is my message text.";
sqs.SendMessage(sendMessageRequest);
```

4. Receive a message

```
ReceiveMessageRequest receiveMessageRequest =
    new ReceiveMessageRequest();
receiveMessageRequest.QueueUrl = myQueueUrl;
ReceiveMessageResponse receiveMessageResponse =
    sqs.ReceiveMessage(receiveMessageRequest);
```

5. Delete a message

```
DeleteMessageRequest deleteRequest =
    new DeleteMessageRequest();
deleteRequest.QueueUrl = myQueueUrl;
deleteRequest.ReceiptHandle = messageReceiptHandle;
DeleteMessageResponse DelMsgResponse =
    sqs.DeleteMessage(deleteRequest);
```

6. Delete a queue

```
DeleteQueueRequest sqsDelRequest = new DeleteQueueRequest();
sqsdelRequest.QueueUrl =
    createQueueResponse.CreateQueueResult.QueueUrl;
DeleteQueueResponse delQueueResponse =
    sqs.DeleteQueue(sqsDelRequest);
```

B.7 How to install SNS on Ubuntu 10.04

SNS, the Simple Notification Service, is a web service for: monitoring applications, workflow systems, time-sensitive information updates, mobile applications, and other event-driven applications that require a simple and efficient mechanism for message delivery. SNS “pushes” messages to clients, rather than requiring a user to periodically poll a mailbox or another site for messages.

SNS is based on the publish–subscribe paradigm; it allows a user to define the topics, the transport protocol used (HTTP/HTTPS, Email, SMS, SQS), and the end-point (URL, email address, phone number, SQS queue) for notifications to be delivered.

Ubuntu is an open-source operating system for personal computers based on Debian Linux distribution. The desktop version of Ubuntu¹ supports Intel x86 32-bit and 64-bit architectures.

SNS supports the following actions:

- Add/Remove Permission
- Confirm Subscription
- Create/Delete Topic

¹ Ubuntu is an African humanist philosophy; “ubuntu” is a word in the Bantu language of South Africa meaning “humanity towards others.”

- Get/Set Topic Attributes
- List Subscriptions/Topics/Subscriptions By Topic
- Publish/Subscribe/Unsubscribe

The site <http://awsdocs.s3.amazonaws.com/SNS/latest/sns-qrc.pdf> provides detailed information about each one of these actions.

The following steps must be taken to install an SNS client:

1. Install Java in the *root* directory and then execute the commands

```
deb http://archive.canonical.com/lucid partner
update
install sun-java6-jdk
```

Then, change the default Java settings

```
update-alternatives -config java
```

2. Download the *SNS* client, unzip the file, and change permissions

```
wget http://sns-public-resources.s3.amazonaws.com/
      SimpleNotificationServiceCli-2010-03-31.zip
chmod 775 /root/ SimpleNotificationServiceCli-1.0.2.3/bin
```

3. Start the AWS management console and go to *Security Credentials*. Check the *Access Key ID* and the *Secret Access Key* and create a text file */root/credential.txt* with the following content:

```
AWSAccessKeyId= your_Access_Key_ID
AWSecretKey= your_Secret_Access_Key
```

4. Edit the *.bashrc* file and add

```
export AWS_SNS_HOME=~/SimpleNotificationServiceCli-1.0.2.3/
export AWS_CREDENTIAL_FILE=$HOME/credential.txt
export PATH=$AWS_SNS_HOME/bin
export JAVA_HOME=/usr/lib/jvm/java-6-sun/
```

5. Reboot the system

6. Enter on the command line

```
sns.cmd
```

If the installation was successful, the list of *SNS* commands will be displayed.

B.8 How to create an EC2 placement group and use MPI

An *EC2 Placement Group* is a logical grouping of instances that allows the creation of a virtual cluster. When several instances are launched as an *EC2 Placement Group*, the virtual cluster has a high bandwidth interconnect system suitable for network-bound applications. The cluster computing instances require an HVM (Hardware Virtual Machine) ECB-based machine image, while other instances use a PVM (Paravirtual Machine) image. Such clusters are particularly useful for high performance computing when most applications are communication intensive.

Once a placement group has been created, MPI can be used for communication among the instances in the placement group. MPI is a de facto standard for parallel applications using message passing; it is designed to ensure high performance, scalability, and portability; it is a language-independent “message-passing application programmer interface, together with a protocol and the semantic specifications for how its features must behave in any implementation” [211]. MPI supports point-to-point and collective communications; it is widely used by parallel programs based on the SPMD (Same Program Multiple Data) paradigm.

The following *C* code [211] illustrates the startup of MPI communication for a process group, *MPI_COM_PROCESS_GROUP* consisting of *nprocesses*; each process is identified by its *rank*. The runtime environment *mpirun* or *mpiexec* spawns multiple copies of the program, with the total number of copies determining the number of process ranks in *MPI_COM_PROCESS_GROUP*.

```
#include <mpi.h>
#include <stdio.h>
#include <string.h>
#define TAG 0
#define BUFSIZE 128

int main(int argc, char *argv[])
{
    char idstr[32];
    char buff[BUFSIZE];
    int nprocesses;
    int my_processId;
    int i;
    MPI_Status stat;
    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&nprocesses);
    MPI_Comm_rank(MPI_COMM_WORLD,&my_processId);
```

MPI_SEND and *MPI_RECEIVE* are blocking send and blocking receive, respectively; their syntax is:

```
int MPI_Send(void *buf, int count, MPI_Datatype datatype,
            int dest, int tag,MPI_Comm comm)
int MPI_Recv(void *buf, int count, MPI_Datatype datatype,
            int source, int tag, MPI_Comm comm, MPI_Status *status)
```

with

<i>buf</i>	— initial address of send buffer (choice)
<i>count</i>	— number of elements in send buffer (nonnegative integer)
<i>datatype</i>	— data type of each send buffer element (handle)
<i>dest</i>	— rank of destination (integer)
<i>tag</i>	— message tag (integer)
<i>comm</i>	— communicator (handle).

Once started, every process other than the coordinator, the process with $rank = 0$, sends a message to the entire group and then receives a message from each of the other members of the process group.

```

if(my_processId == 0)
{
    printf("%d: We have %d processes\n", my_processId, nprocesses);
    for(i=1;i<nprocesses;i++)
    {
        sprintf(buff, "Hello %d! ", i);
        MPI_Send(buff, BUFSIZE, MPI_CHAR, i, TAG, MPI_COMM_PROCESS_GROUP);
    }
    for(i=1;i<nprocesses;i++)
    {
        MPI_Recv(buff, BUFSIZE, MPI_CHAR, i, TAG, MPI_COMM_PROCESS_GROUP, &stat);
        printf("%d: %s\n", my_processId, buff);
    }
}
else
{
    /* receive from rank 0: */
    MPI_Recv(buff, BUFSIZE, MPI_CHAR, 0, TAG, MPI_COMM_PROCESS_GROUP, &stat);
    sprintf(idstr, "Processor %d ", my_processId);
    strncat(buff, idstr, BUFSIZE-1);
    strncat(buff, "reporting for duty\n", BUFSIZE-1);
    /* send to rank 0: */
    MPI_Send(buff, BUFSIZE, MPI_CHAR, 0, TAG, MPI_COMM_PROCESS_GROUP);
}

MPI_Finalize();
return 0;
}

```

An example of cloud computing using the MPI is described in [167]. An example of MPI use on EC2 is at <http://rc.fas.harvard.edu/faq/amazonec2>.

B.9 StarCluster—a cluster computing toolkit for EC2

StarCluster, <http://star.mit.edu/cluster/>, is an open-source cluster-computing toolkit for EC2. The system assigns user-friendly names to the nodes of the virtual cluster, and the cluster is so configured to allow *ssh* from any node of the cluster to any other nodes. It allows attaching EBS volumes to the cluster for persistent storage and provides an API for executing OS commands, such as copying files. StarCluster supports dynamic cluster reconfiguration and launching spot instances to reduce the service costs.

StarCluster AMIs consist of several scientific libraries:

1. OpenMPI—for writing parallel applications.
2. Automatically Tuned Linear Algebra Software (ATLAS)—optimized for Amazon EC2 larger instances; see <http://math-atlas.sourceforge.net/>.
3. NumPy/SciPy compiled against the optimized ATLAS install. SciPy is a Python-based ecosystem of open-source software for mathematics, science, and engineering. NumPy is a set of tools for integrating C/C++ and Fortran code for useful linear algebra, Fourier transform, and random number capabilities; see <https://www.scipy.org/scipylib/>.
4. IPython—interactive parallel computing in Python; see <https://ipython.org/>.

Several other important features of the StarCluster are: (1) Support for starting/stopping EBS-backed clusters on EC2. (2) Elastic Load Balancing—using Sun Grid Engine queue statistics.(3) Support for specifying instance types on a per-node basis. (4) A number of commands for EC2 and S3 operations, including the ones in Table B.2.

B.10 An alternative setting of an MPI virtual cluster

An alternative setting up for an MPI cluster is described in [313]. The instances used are *cc2.8xlarge* with 2x Intel Xeon E5-2670 processors and \approx 60 GB of RAM per node and spot instances rather than reserved ones were used, thus saving about 90% of the cost. The VM virtualization layer used by the *cc2.8xlarge* instances is thinner than the one of other instances.

Setting user's instances involves several steps:

1. Select a VM image and an instance type.
2. Define a placement group.
3. Configure the storage.
4. Define the VM tags to manage the instances.
5. Create of a key pair.
6. Set up the security group.
7. Launch the instances.

The nest phase is the configuration of the virtual cluster (VC). All instances of the virtual cluster are booted with the same operating system and share a 10 Gbps Ethernet subnet. The public IP address of the booted instances is available from the EC2 Dashboard by selecting their entry under the “Instances” page. The preliminary steps for the VC configurations are:

Table B.2 StarCluster commands for EC2 and S3 operations.

Command	Function
listinstances	List all running EC2 instances
listspots	List all EC2 spot instance requests
listimages	List all registered EC2 images (AMIs)
listpublic	List all public StarCluster images on EC2
listkeypairs	List all EC2 keypairs
createkey	Create a new Amazon EC2 keypair
removekey	Remove a keypair from Amazon EC2
s3image	Create a new instance-store (S3) AMI from a running EC2 instance
ebsimage	Create a new EBS image (AMI) from a running EC2 instance
removeimage	Deregister an EC2 image (AMI)
createvolume	Create a new EBS volume for use with StarCluster
listvolumes	List all EBS volumes
resizevolume	Resize an existing EBS volume
removevolume	Delete one or more EBS volumes
spothistory	Show spot instance pricing history stats
showconsole	Show console output for an EC2 instance
listregions	List all EC2 regions
listzones	List all EC2 availability zones in the current region
listbuckets	List all S3 buckets
showbucket	Show all files in an S3 bucket

- Create nodes aliases—add their internal IP addresses to */etc/hosts* as *node1, node2, node3, node4*.
- Create aliases for files to be transferred from the master (*node1*) to workers *node2, node3, node4*.
- Enable password-less ssh between instances.

There are two options to install and lunch MPI in every node, the first for *OpenMPI* and the second for *mpich*, both use the *yum* package.

```
sudo yum install openmpi-devel
sudo yum install mpich-devel
```

Once the GCC compiler, the MPI runtime, and the mpicc wrapper are available in every node, the following commands must be added to the *.bashrc* file on all compute nodes

```
export PATH=/usr/lib64/openmpi/bin:$PATH
export LD_LIBRARY_PATH=/usr/lib64/openmpi/lib
```

The following command runs a program “application.xx” available on every node

```
mpirun -np 32 -hostfile ~/nodefile ~/application.xx
```

with the hostfile called “nodefile” compatible with OpenMPI created as follows:

```
$ cat ~/nodefile
node1 slots=16
node2 slots=16
node3 slots=16
node4 slots=16
```

A 64-way MPI job is created. Each node has 16 cores and 16 hyperthreads.

B.11 How to install hadoop on eclipse on a windows system

Eclipse (<http://www.eclipse.org/>) is a software development environment. Eclipse consists of an integrated development environment (IDE) and an extensible plug-in system. It is written mostly in Java and can be used to develop applications in Java and, by means of various plug-ins, in C, C++, Perl, PHP, Python, R, Ruby, and several other languages. The IDE is often called Eclipse CDT for C/C++, Eclipse JDT for Java, and Eclipse PDT for PHP.

The software packages used are:

- Apache Hadoop is a software framework that supports data-intensive distributed applications under a free license. *Hadoop* was inspired by Google’s MapReduce; see Section 11.5 for a discussion of MapReduce and Section 11.6 for an application using *Hadoop*.
- *Cygwin* is a Unix-like environment for Microsoft Windows. It is open-source software, released under the GNU General Public License version 2. *Cygwin* consists of: (1) a dynamic-link library (DLL) as an API compatibility layer providing a substantial part of the POSIX API functionality; and (2) an extensive collection of software tools and applications that provide a Unix-like look and feel.

A. Prerequisites

- *Java 1.6*; set JAVA_Home = path where *JDK* is installed.
- *Eclipse Europa 3.3.2*

Note: the Hadoop plugin was specially designed for Europa, and newer releases of *Eclipse* might have some issues with *Hadoop* plugin.

B. SSH Installation

1. Install *cygwin* using the installer downloaded from <http://www.cygwin.com>. From the *Select Packages* window, select the *openssh* and *openssl* under Net.
Note: Create a desktop icon when asked during installation.
2. Display the “Environment Variables” panel

Computer -> System Properties -> Advanced System Settings
-> Environment Variables

Click on the variable named *Path* and press *Edit*; append the following value to the path variable

`;c:\cygwin\bin;c:\cygwin\usr\bin`

3. Configure the *ssh daemon* using *cygwin*. Left click on the *cygwin* icon on desktop and click “Run as Administrator.” Type in the command window of *cygwin*

```
ssh-host-config.
```

4. Answer “Yes” when prompted; *sshd* should be installed as a service; answer “No” to all other questions.
5. Start the *cygwin* service by navigating to

Control Panel -> Administrative Tools -> Services

Look for *cygwin sshd* and start the service.

6. Open *cygwin* command prompt and execute the following command to generate keys

```
ssh-keygen
```

7. When prompted for filenames and pass phrases press ENTER to accept default values. After the command has finished generating keys, enter the following command to change into your *.ssh* directory:

```
cd ~/.ssh
```

8. Check if the keys were indeed generated

```
ls -l
```

9. The two files *id_rsa.pub* and *id_rsa* with recent creation dates contain authorization keys.
10. To register the new authorization keys, enter the following command (Note: the sharply-angled double brackets are very important)

```
cat id_rsa.pub >> authorized_keys
```

11. Check if the keys were set up correctly

```
ssh localhost
```

12. Since it is a new *ssh* installation, you will be warned that authenticity of the host could not be established and will be asked whether you really want to connect. Answer YES and press ENTER. You should see the *cygwin* prompt again, which means that you have successfully connected.

13. Now, execute again the command:

```
ssh localhost
```

This time no prompt should appear.

C. Download *hadoop*

1. Download *hadoop 0.20.1* and place in a directory such as

```
C:\Java
```

2. Open the *cygwin* command prompt and execute

```
cd
```

3. Enable the home directory folder to be shown in the Windows Explorer window

```
explorer
```

```

C: ~/hadoop-0.20.2
$ cd hadoop-0.20.2
Rt@Rt-PC ~~/hadoop-0.20.2
$ ls -l
total 4885
-rw-r--r-- 1 Rt None 348624 Feb 19 2010 CHANGES.txt
-rw-r--r-- 1 Rt None 13366 Feb 19 2010 LICENSE.txt
-rw-r--r-- 1 Rt None 101 Feb 19 2010 NOTICE.txt
-rw-r--r-- 1 Rt None 1366 Feb 19 2010 README.txt
drwxr-xr-x+ 1 Rt None 0 Sep 22 19:36 bin
-rw-r--r-- 1 Rt None 74035 Feb 19 2010 build.xml
drwxr-xr-x+ 1 Rt None 0 Feb 19 2010 c++
drwxr-xr-x+ 1 Rt None 0 Sep 22 20:08 conf
drwxr-xr-x+ 1 Rt None 0 Feb 19 2010 contrib
drwxr-xr-x+ 1 Rt None 0 Sep 22 19:35 docs
-rw-r--r-- 1 Rt None 6839 Feb 19 2010 hadoop-0.20.2-ant.jar
-rw-r--r-- 1 Rt None 2689741 Feb 19 2010 hadoop-0.20.2-core.jar
-rw-r--r-- 1 Rt None 142466 Feb 19 2010 hadoop-0.20.2-examples.jar
-rw-r--r-- 1 Rt None 1563859 Feb 19 2010 hadoop-0.20.2-test.jar
-rw-r--r-- 1 Rt None 69940 Feb 19 2010 hadoop-0.20.2-tools.jar
drwxr-xr-x+ 1 Rt None 0 Sep 22 19:35 ivy
-rw-r--r-- 1 Rt None 8852 Feb 19 2010 ivy.xml
drwxr-xr-x+ 1 Rt None 0 Sep 22 19:35 lib
drwxr-xr-x+ 1 Rt None 0 Sep 22 19:35 librecordio
drwxr-xr-x+ 1 Rt None 0 Sep 22 20:08 logs
drwxr-xr-x+ 1 Rt None 0 Sep 22 19:36 src
drwxr-xr-x+ 1 Rt None 0 Feb 19 2010 webapps

```

FIGURE B.9

The result of unpacking *hadoop*.

4. Open another Windows Explorer window and navigate to the folder that contains the downloaded *hadoop* archive.
 5. Copy the *hadoop* archive into the home directory folder.
- D. Unpack *hadoop***
1. Open a new cygwin window and execute

```
tar -xzf hadoop-0.20.1.tar.gz
```

2. List the contents of the home directory

```
ls -l
```

A newly created directory called *hadoop-0.20.1* should be seen. Execute

```
cd hadoop-0.20.1
ls -l
```

The files listed in Fig. B.9 should be seen.

- E. Set properties in configuration file**
1. Open a new cygwin window and execute the following commands

```
cd hadoop-0.20.1
cd conf
explorer
```

```

11/09/26 13:40:41 INFO namenode.FSNamesystem: supergroup=supergroup
11/09/26 13:40:41 INFO namenode.FSNamesystem: isPermissionEnabled=true
11/09/26 13:40:42 INFO common.Storage: Image file of size 98 saved in 0 second

11/09/26 13:40:42 INFO common.Storage: Storage directory \tmp\hadoop-Rt\dfs\na
has been successfully formatted.
11/09/26 13:40:42 INFO namenode.NameNode: SHUTDOWN_MSG:
*****SHUTDOWN_MSG: Shutting down NameNode at Rt-PC/192.168.1.231
*****/
Rt@Rt-PC ~ /hadoop-0.20.2

```

FIGURE B.10

The creation of HDFS.

2. The last command will cause the Explorer window for the *conf* directory to pop up. Minimize it for now or move it to the side.
3. Launch *Eclipse* or a text editor such as *Notepad ++* and navigate to the *conf* directory and open the file *hadoop-site* to insert the following lines between *<configuration>* and *</configuration>* tags.

```

<property>
<name>fs.default.name</name>
<value>hdfs://localhost:9100</value>
</property>
<property>
<name>mapred.job.tracker</name>
<value>localhost:9101</value>
</property>
<property>
<name>dfs.replication</name>
<value>1</value>
</property>

```

F. Format the Namenode

Format the *namenode* to create a Hadoop Distributed File System (HDFS). Open a new *cygwin* window and execute the following commands:

```

cd hadoop-0.20.1
mkdir logs
bin/hadoop namenode -format

```

When the formatting of the *namenode* is finished, the message in Fig. B.10 appears.

B.12 Exercises and problems

Problem 1. Establish an AWS account. Use the AWS management console to launch an EC2 instance and connect to it.

- Problem 2.** Launch three EC2 instances; the computations carried out by the three instances should consist of two phases, and the second phase should be started only after all instances have finished the first stage. Design a protocol and use *Simple Queue Service (SQS)* to implement the barrier synchronization after the first phase.
- Problem 3.** Use the *Zookeeper* to implement the coordination model in Problem 2.
- Problem 4.** Use the *Simple Workflow Service (SWF)* to implement the coordination model in Problem 2. Compare the three methods.
- Problem 5.** Upload several (10–20) large image files to an S3 bucket. Start an instance that retrieves the images from the S3 bucket and compute the retrieval time. Use the *ElastiCache* service and compare the retrieval time for the two cases.
- Problem 6.** Numerical simulations are ideal applications for cloud computing. Output data analysis of a simulation experiment requires the computation of confidence intervals for the mean for the quantity of interest [297]. This implies that one must run multiple batches of simulation, compute the average value of the quantity of interest for each batch, and then calculate say 95% confidence intervals for the mean. Use the *CloudFormation* service to carry out a simulation using multiple cloud instances that store partial results in S3 and then another instance computes the confidence interval for the mean.
- Problem 7.** Run an application that takes advantage of the *Autoscaling* service.
- Problem 8.** Use the *Elastic Beanstalk* service to run an application and compare it with the case when the *Autoscaling* service was used.
- Problem 9.** Design a cloud service and a testing environment; use the *Elastic Beanstalk* service to support automatic scaling up and down and use the *Elastic Load Balancer* to distribute the incoming service request to multiple instances of the application.

Literature

- [1] W.M.P. van der Aalst, A.H. ter Hofstede, B. Kiepuszewski, A.P. Barros, Workflow patterns, Technical Report, Eindhoven University of Technology, 2000.
- [2] D. Abadi, Consistency tradeoffs in modern distributed database system design, Computer 45 (2) (2012) 37–42.
- [3] D. Abadi, et al., The Beckman report on database research, Communications of the ACM 59 (2) (2016) 92–99.
- [4] T.F. Abdelzaher, K.G. Shin, N. Bhatti, Performance guarantees for web server end-system: a control theoretic approach, IEEE Transactions on Parallel and Distributed Systems 13 (1) (2002) 80–96.
- [5] A. Abouzeid, K. Bajda-Pawlikowski, D.J. Abadi, A. Rasin, A. Silberschatz, HadoopDB: an architectural hybrid of MapReduce and DBMS technologies for analytical workloads, Proceedings of the VLDB Endowment 2 (1) (2009) 922–933.
- [6] B. Abrahao, V. Almeida, J. Almeida, A. Zhang, D. Beyer, F. Safai, Self-adaptive SLA-driven capacity management for Internet services, in: Proc. IEEE/IFIP Network Operations and Management Symp., 2006, pp. 557–568.
- [7] M. Abuelela, S. Olariu, Taking VANET to the clouds, in: Proc. 8-th ACM Int. Conf. on Advanced in Mobile Computing (MoMM’2010), Dec. 2010.
- [8] H. Abu-Libdeh, L. Princehouse, H. Weatherspoon, RACS: a case for cloud storage diversity, in: Proc. ACM Symp. Cloud Computing, CD Proc., ISBN 978-1-4503-0036-0, 2010.
- [9] D. Abts, The Cray XT4 and Seastar 3-D torus interconnect, in: D. Padua (Ed.), Encyclopedia of Parallel Computing, Part 3, Springer-Verlag, Heidelberg, 2011, pp. 470–477.
- [10] D. Abts, M.R. Marty, P.M. Wells, P. Klausler, H. Liu, Energy proportional datacenter networks, in: Proc. ACM IEEE Int. Symp. Comp. Arch., 2010, pp. 338–347.
- [11] B. Addis, D. Ardagna, B. Panicucci, L. Zhang, Autonomic management of cloud service centers with availability guarantees, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 220–227.
- [12] S. Agarwal, H. Milner, A. Kleiner, A. Talwarkar, M. Jordan, S. Madden, B. Mozafari, I. Stoica, Knowing when you’re wrong: building fast and reliable approximate query processing systems, in: Proc. ACM SIGMOD Int. Conf. on Management of Data, ACM Press, 2014, pp. 481–492.
- [13] M. Ahmadian, F. Plohan, Z. Roessler, D.C. Marinescu, SecureNoSQL: an approach for secure search of encrypted NoSQL databases in the public cloud, International Journal of Information Management 37 (2017) 63–74.
- [14] F. Ahmad, M. Kazim, A. Adnane, A. Awad, Vehicular cloud networks: architecture, applications and security, in: Proc. 8-th IEEE/ACM Int. Conf. on Utility and Cloud Computing, UUC’2015, 2015, pp. 571–576.
- [15] T. Akidau, A. Balikov, K. Bekiröglu, S. Chernyak, J. Haberman, R. Lax, S. McVeety, D. Mills, P. Nordstrom, S. Whittle, MillWheel: fault-tolerant stream processing at Internet scale, Proceedings of the VLDB Endowment 6 (11) (2013) 1033–1044.
- [16] T. Akidau, R. Bradshaw, C. Chambers, S. Chernyak, R.J. Fernandez-Moctezuma, R. Lax, S. McVeety, D. Mills, F. Perry, E. Schmidt, S. Whittle, The data-flow model: a practical approach to balancing correctness, latency, and cost in massive-scale, unbounded, out-of-order data processing, Proceedings of the VLDB Endowment 8 (12) (2015) 1792–1803.
- [17] R. Albert, H. Jeong, A.-L. Barabási, The diameter of the world wide web, Nature 401 (1999) 130–131.
- [18] R. Albert, H. Jeong, A.-L. Barabási, Error and attack tolerance of complex networks, Nature 406 (2000) 378–382.

- [19] R. Albert, A-L. Barabási, Statistical mechanics of complex networks, *Reviews of Modern Physics* 72 (1) (2002) 48–97.
- [20] M. Al-Fares, A. Loukissas, A. Vahdat, A scalable, commodity data center network architecture, in: Proc. ACM SIGCOM Conf. on Data Communication, 2008, pp. 63–74.
- [21] G.M. Amdahl, Validity of the single-processor approach to achieving large-scale computing capabilities, in: Proc. Conf. American Federation of Inf. Proc. Soc. Conf., AFIPS Press, 1967, pp. 483–485.
- [22] Y. Amir, B. Awerbuch, A. Barak, R.S. Borgstrom, A. Keren, An opportunity cost approach for job assignment in a scalable computing cluster, *IEEE Transactions on Parallel and Distributed Systems* 11 (7) (2000) 760–768.
- [23] G. Ananthanarayanan, A. Ghodsi, S. Shenker, I. Stoica, Disk-locality in datacenter computing considered irrelevant, in: Procs. 13th USENIX Conf. on Hot Topics in Operating Systems, 2011, pp. 12–17.
- [24] R. Ananthanarayanan, V. Basker, S. Das, A. Gupta, H. Jiang, T. Qiu, A. Reznichenko, D. Ryabkov, M. Singh, S. Venkataraman, Photon: fault-tolerant and scalable joining of continuous data streams, in: Proc. ACM SIGMOD Int. Conf. on Management of Data, 2013, pp. 577–588.
- [25] T. Andrei, Cloud computing challenges and related security issues, <http://www1.cse.wustl.edu/jain/cse571-09/ftp/cloud/index.html>. (Accessed August 2015).
- [26] D.P. Anderson, BOINC: a system for public-resource computing and storage, in: Proc. 5th IEEE/ACM Int. Workshop Grid Computing, 2004, pp. 4–10.
- [27] R. Aoun, E.A. Doumith, M. Gagnaire, Resource provisioning for enriched services in cloud environment, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 296–303.
- [28] Apache, Apache capacity scheduler, https://svn.apache.org/repos/asf/hadoop/common/tags/release-0.19.1/docs/capacity_scheduler.pdf. (Accessed August 2016).
- [29] D. Ardagna, M. Trubian, L. Zhang, SLA based resource allocation policies in autonomic environments, *Journal of Parallel and Distributed Computing* 67 (3) (2007) 259–270.
- [30] D. Ardagna, B. Panicucci, M. Trubian, L. Zhang, Energy-aware autonomic resource allocation in multi-tier virtualized environments, *IEEE Transactions on Services Computing* 5 (1) (2012) 2–19.
- [31] S. Arif, S. Olariu, J. Wang, G. Yan, W. Yang, I. Khalil, Datacenter at the airport: reasoning about time-dependent parking lot occupancy, *IEEE Transactions on Parallel and Distributed Systems* 23 (11) (2012) 2067–2080.
- [32] M. Armbrust, A. Fox, R. Griffith, A.D. Joseph, R. Katz, A. Konwinski, G. Lee, D. Paterson, A. Rabkin, I. Stoica, M. Zaharia, Above the clouds: a Berkeley view of cloud computing, Technical Report UCB/EECS-2009-28, 2009.
- [33] D. Artz, Y. Gil, A survey of trust in computer science and the Semantic Web, *Journal of Web Semantics. Science, Services, and Agents on World Wide Web* 1 (2) (2007) 58–71.
- [34] M.J. Atallah, C. Lock Black, D.C. Marinescu, H.J. Siegel, T.L. Casavant, Models and algorithms for co-scheduling compute-intensive tasks on a network of workstations, *Journal of Parallel and Distributed Computing* 16 (1992) 319–327.
- [35] M. Auty, S. Creese, M. Goldsmith, P. Hopkins, Inadequacies of current risk controls for the cloud, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 659–666.
- [36] L. Ausubel, P. Cramton, Auctioning many divisible goods, *Journal of the European Economic Association* 2 (2–3) (2004) 480–493.
- [37] L. Ausubel, P. Cramton, P. Milgrom, The clock-proxy auction: a practical combinatorial auction design, in: P. Cramton, Y. Shoham, R. Steinberg (Eds.), *Combinatorial Auctions*, MIT Press, Cambridge, Mass, 2006, Chapter 5.
- [38] A. Avisienis, J.C. Laprie, B. Randell, C. Landwehr, Basic concepts and taxonomy of dependable and secure computing, *IEEE Transactions on Dependable and Secure Computing* 1 (1) (2004) 11–33.
- [39] Y. Azar, A. Broder, A. Karlin, E. Upfal, Balanced allocations, in: Proc. 26th ACM Symp. on the Theory of Computing, 1994, pp. 593–602.

- [40] Ö. Babaoglu, K. Marzullo, Consistent global states, in: Sape Mullender (Ed.), *Distributed Systems*, 2nd edition, Addison-Wesley, Reading, Ma, 1993, pp. 55–96.
- [41] M.J. Bach, M.W. Luppi, A.S. Melamed, K. Yueh, A remote-file cache for RFS, in: Proc. Summer 1987 USENIX Conf., 1987, pp. 275–280.
- [42] X. Bai, H. Yu, G.Q. Wang, Y. Ji, G.M. Marinescu, D.C. Marinescu, L. Böloni, Coordination in intelligent grid environments, *Proceedings of the IEEE* 93 (3) (2005) 613–630.
- [43] J. Baker, C. Bond, J.C. Corbett, J.J. Furman, A. Khorlin, J. Larson, J.-M. Léon, Y. Li, A. Lloyd, V. Yushprakh, Megastore: providing scalable, highly available storage for interactive services, in: Proc. 5th Biennial Conf. Innovative Data Systems Research, 2011, pp. 223–234.
- [44] M. Balduzzi, J. Zaddach, D. Balzarotti, E. Kirda, S. Loureiro, A security analysis of Amazon’s elastic compute cloud service, in: Proc. 27th Annual ACM Symp. Applied Computing, 2012, pp. 1427–1434.
- [45] J. Baliga, R.W. Ayre, K. Hinton, R.S. Tucker, Green cloud computing: balancing energy in processing, storage, and transport, *Proceedings of the IEEE* 99 (1) (2011) 149–167.
- [46] A-L. Barabási, R. Albert, Emergence of scaling in random networks, *Science* 286 (5439) (1999) 509–512.
- [47] A-L. Barabási, R. Albert, H. Jeong, Scale-free theory of random networks; the topology of World Wide Web, *Physica A* 281 (2000) 69–77.
- [48] P. Barham, B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, R. Neugebauer, I. Pratt, A. Warfield, Xen and the art of virtualization, in: Proc. 19th ACM Symp. Operating Systems Principles, 2003, pp. 164–177.
- [49] B. Baron, M. Campista, P. Spathis, L.H. Costa, M. Dias de Amorim, O.C. Duarte, G. Pujolle, Y. Viniotis, Virtualizing vehicular node resources: feasibility study of virtual machine migration, *Vehicular Communications* 4 (2016) 39–46.
- [50] L.A. Barroso, J. Dean, U. Hözle, Web search for a planet: the Google cluster architecture, *IEEE MICRO* 23 (2) (2003) 22–28.
- [51] L.A. Barroso, U. Hözle, The case for energy-proportional computing, *IEEE Computer* 40 (12) (2007) 33–37.
- [52] L.A. Barosso, U. Hözle, P. Ranganathan, *The Datacenter as a Computer; an Introduction to the Design of Warehouse-Scale Machines*, third edition, Morgan and Claypool, 2013.
- [53] L. Bottou, Large-scale machine learning with stochastic gradient descent, in: Proc. Int. Conf. on Computational Statistics, Springer Verlag, Heidelberg, 2010, pp. 177–186.
- [54] A. Bavier, T. Voigt, M. Wawrzoniak, L. Peterson, P. Gunningberg, SILK: Scout paths in the Linux kernel, Technical Report 2002-009, Uppsala University, Department of Information Technology, February 2002.
- [55] G. Bell, Massively parallel computers: why not parallel computers for the masses?, in: Proc. 4th Symp. Frontiers of Massively Parallel Computing, IEEE, 1992, pp. 292–297.
- [56] R. Bell, M. Koren, C. Volinsky, The BellKor 2008 solution to the Netflix Prize, Technical report, AT&T Labs, 2008.
- [57] M. Ben-Yehuda, M.D. Day, Z. Dubitzky, M. Factor, N. Har’El, A. Gordon, A. Ligouri, O. Wasserman, B.-A. Yassour, The Turtles project: design and implementation of nested virtualization, in: Proc. 9th USENIX Conf. on OS Design and Implementation, 2010, pp. 423–436.
- [58] D. Bernstein, E. Ludvigson, K. Sankar, S. Diamond, M. Morrow, Blueprint for the Intercloud - protocols and formats for cloud computing interoperability, in: Proc. 4th Int. Conf. Internet and Web Applications and Services, ICIW ’09, 2009, pp. 328–336.
- [59] D. Bernstein, D. Vij, Intercloud security considerations, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 537–544.
- [60] D. Bernstein, D. Vij, S. Diamond, An Intercloud cloud computing economy - technology, governance, and market blueprints, in: Proc. SRII Global Conference, 2011, pp. 293–299.
- [61] D. Bertsekas, R. Gallagher, *Data Networks*, second edition, Prentice Hall, Upper Saddle River, NJ, 1992.
- [62] S. Bertram, M. Boniface, M. Surridge, N. Briscombe, M. Hall-May, On-demand dynamic security for risk-based secure collaboration in clouds, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 518–525.
- [63] S. Bhattacharjee, D.C. Marinescu, A cloud service for trust management in cognitive radio networks, *International Journal of Cloud Computing* 3 (14) (2014) 326–353.

- [64] M. Blackburn, A. Hawkins, Unused server survey results analysis, <https://silo.tips/download/unused-servers-survey-results-analysis>. (Accessed October 2021).
- [65] N.J. Boden, D. Cohen, R.E. Felderman, A.E. Kulawik, C.L. Seitz, J.N. Seizovic, W-K. Su, Myrinet – a gigabit-per-second local-area network, IEEE MICRO 5 (1) (1995) 29–36.
- [66] A. Boldyreva, N. Chenette, Y. Lee, A. O'Neill, Order-preserving symmetric encryption, in: Advances in Cryptology, Springer-Verlag, Heidelberg, 2009, pp. 224–241.
- [67] B. Bollobás, Random Graphs, Academic Press, London, 1985.
- [68] K. Boloor, R. Chirkova, Y. Viniotis, T. Salo, Dynamic request allocation and scheduling for context aware applications subject to a percentile response time SLA in a distributed cloud, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 464–472.
- [69] A. Boukerche, R.E. De Grande, Vehicular cloud computing: architectures, applications, and mobility, Computer Networks 135 (2018) 171–189.
- [70] I. Brandic, S. Dustdar, T. Ansett, D. Schumm, F. Leymann, R. Konrad, Compliant cloud computing (C3): architecture and language support for user-driven compliance management in clouds, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 244–251.
- [71] S. Brandt, S. Banachowski, C. Lin, T. Bisson, Dynamic integrated scheduling of hard real-time, soft real-time, and non-real-time processes, in: Proc. IEEE Real-Time Systems Symp., 2003, pp. 396–409.
- [72] S. Brin, L. Page, The anatomy of a large-scale hypertextual web search engine, Journal of Computer Networks and ISDN Systems 30 (1–7) (1998) 107–117.
- [73] N.F. Britton, Essential Mathematical Biology, Springer-Verlag, Heidelberg, 2004.
- [74] S. Brooks, C.E.A. Hoare, A.W. Roscoe, A theory of communicating sequential processes, Journal of the ACM 31 (3) (1984) 560–599.
- [75] C. Brooks, Enterprise NoSQL for Dummies, Wiley, New Jersey, NJ, 2014.
- [76] M. Buddhikot, K. Ryan, Spectrum management in coordinated dynamic spectrum access based cellular networks, in: Proc. IEEE Int. Symp. New Frontiers in Dynamic Spectrum Access Networks, 2005, pp. 299–307.
- [77] Y. Bu, B. Howe, M. Balazinska, M.D. Ernst, HaLoop: efficient iterative data processing on large clusters, Proceedings of the VLDB Endowment 3 (2010) 285–296.
- [78] C. Bunch, N. Chohan, C. Krintz, J. Chohan, J. Kupferman, P. Lakhina, Y. Li, Y. Nomura, An evaluation of distributed data stores using the AppScale cloud platform, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 305–312.
- [79] A.W. Burks, H.H. Goldstine, J. von Neumann, Preliminary discussion of the logical design of an electronic computer instrument, in: Report to the US Army Ordnance Department, 1946, also in: W. Asprey, A.W. Burks (Eds.), Papers of John von Neumann, MIT Press, Cambridge, MA, 1987, pp. 97–146.
- [80] B. Burns, B. Grant, D. Oppenheimer, E. Brewer, J. Wilkes, Borg, Omega, and Kubernetes; lessons learned from three container-management systems over a decade, ACM Queue 14 (1) (2016) 70–93.
- [81] M. Burrows, The Chubby lock service for loosely-coupled distributed systems, in: Proc. USENIX Symp. OS Design and Implementation, 2006, pp. 335–350.
- [82] B. Buth, J. Peleska, H. Shi, Combining methods for the livelock analysis of a fault-tolerant system, in: Proc. 7th Int. Conf. on Algebraic Methodology and Software Technology, 1998, pp. 124–139.
- [83] R. Buyya, R. Ranjan, R. Calheiros, Intercloud: utility-oriented federation of cloud computing environments for scaling of application services, in: Proc. 10th Int. Conf. Algorithms and Architectures for Parallel Processing, Part I, 2010, pp. 13–31.
- [84] C. Cacciari, F. D'Andria, M. Gonzalo, B. Hagemeier, D. Mallmann, J. Martrat, D.G. Perez, A. Rumpl, W. Ziegler, C. Zsigri, elasticLM: a novel approach for software licensing in distributed computing infrastructures, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 67–74.
- [85] D. Cash, S. Jarecki, C. Jutla, C.H. Krawczyk, M.C. Rosu, M. Steiner, Highly-scalable searchable symmetric encryption with support for Boolean queries, in: Advances in Cryptology, Springer-Verlag, Heidelberg, 2013, pp. 353–373.

- [86] A.J. Canty, A.C. Davison, D.V. Hinkley, V. Ventura, Bootstrap diagnostics and remedies, *Canadian Journal of Statistics* 34 (1) (2006) 5–27.
- [87] A.G. Carlyle, S.L. Harrell, P.M. Smith, Cost-effective HPC: the community or the cloud?, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 169–176.
- [88] E. Caron, F. Desprez, A. Muresan, Forecasting for grid and cloud computing on-demand resources based on pattern matching, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 456–463.
- [89] R. Cattell, Scalable SQL and NoSQL data stores, <http://cattell.net/databases/Datastores.pdf>, 2011. (Accessed August 2015).
- [90] C. Chambers, A. Raniwala, F. Perry, S. Adams, R.R. Henry, R. Bradshaw, N. Weizenbaum, FlumeJava: easy, efficient data-parallel pipelines, in: Proc. ACM SIGPLAN Conf. Programming Language Design and Implementation, 2010, pp. 363–375.
- [91] A. Chandra, P. Goyal, P. Shenoy, Quantifying the Benefits of Resource Multiplexing in on-Demand Data Centers, Computer Science Department Faculty Publication Series, vol. 20, 2003, http://scholarworks.umass.edu/cs_faculty_pubs/20. (Accessed January 2017).
- [92] T.D. Chandra, R. Griesemer, J. Redstone, Paxos made live: an engineering perspective, in: Proc. 26th ACM Symp. Principles of Distributed Computing, 2007, pp. 398–407.
- [93] K.M. Chandy, L. Lamport, Distributed snapshots: determining global states of distributed systems, *ACM Transactions on Computer Systems* 3 (1) (1985) 63–75.
- [94] F. Chang, J. Dean, S. Ghemawat, W.C. Hsieh, D.A. Wallach, M. Burrows, T. Chandra, A. Fikes, R.E. Gruber, Bigtable: a distributed storage system for structured data, in: Proc. USENIX Symp. on OS Design and Implementation, 2006, pp. 205–218.
- [95] V. Chang, G. Wills, D. De Roure, A review of cloud business models and sustainability, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 43–50.
- [96] F. Chang, J. Ren, R. Viswanathan, Optimal resource allocation in clouds, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 418–425.
- [97] L. Chang, Z. Wang, T. Ma, L. Jian, A. Golgshuv, L. Lonergan, J. Cohen, C. Welton, G. Sherry, M. Bhandarkar, HAWQ: a massively parallel processing SQL engine in Hadoop, in: Proc. ACM SIGMOD Int. Conf. Data Management, 2014, pp. 1223–1234.
- [98] K. Chard, S. Caton, O. Rana, K. Bubendorfer, Social cloud: cloud computing in social networks, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 99–106.
- [99] A. Chazalet, Service level checking in the cloud computing context, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 297–304.
- [100] P.M. Chen, B.D. Noble, When virtual is better than real, in: Proc. 8th Workshop Hot Topics in Operating Systems, 2001, pp. 133–141.
- [101] H. Chen, P. Liu, R. Chen, B. Zang, When OO meets system software: Rethinking the design of VMMs, Technical Report PPITR-2007-08003, Fudan University, Parallel Processing Institute, 2007, pp. 1–9.
- [102] T.M. Chen, S. Abu-Nimeh, Lessons from Stuxnet, *Computer* 44 (4) (2011) 91–93.
- [103] R. Chen, J.-M. Park, K. Bian, Robust distributed spectrum sensing in cognitive radio networks, in: Proc. IEEE Conf. Computer Communications, 2008, pp. 1876–1884.
- [104] Y. Chen, S. Alspaugh, R. Katz, Interactive analytical processing in big data systems: a cross-industry study of MapReduce workloads, *Proceedings of the VLDB Endowment* 5 (12) (2012) 1802–1813.
- [105] J. Cheney, L. Chiticariu, W.-C. Tan, Provenance in databases: why, how, and where, *Foundations and Trends in Databases* 1 (4) (2007) 379–474.
- [106] D. Chiu, G. Agarwal, Evaluating cache and storage options on the Amazon Web services cloud, in: Proc. IEEE/ACM Int. Symp. Cluster, Cloud and Grid Computing, 2011, pp. 362–371.
- [107] R. Chow, P. Golle, M. Jacobsson, E. Shi, J. Staddon, R. Masouka, J. Mollina, Controlling data on the cloud: outsourcing computations without outsourcing control, in: Proc. ACM Cloud Computing Security Workshop, 2009, pp. 85–90.

- [108] B. Clark, T. Deshane, E. Dow, S. Evabchik, M. Finlayson, J. Herne, J.N. Matthews, Xen and the art of repeated research, in: Proc. USENIX Annual Technical Conference, 2004, pp. 135–144.
- [109] R.A. Clarke, R.K. Knake, Cyber War: The Next Threat to National Security and What to do About It, Harper Collins, 2012.
- [110] C. Clos, A study of non-blocking switching networks, *The Bell System Technical Journal* 32 (2) (1953) 406–425.
- [111] E.F. Codd, A relational model of data for large shared data banks, *Communications of the ACM* 13 (6) (1970) 377–387.
- [112] E. Coffman, M.J. Elphick, S. Shoshani, System deadlocks, *Computing Surveys* 3 (2) (1971) 67–78.
- [113] L. Colitti, S.H. Gunderson, E. Kline, T. Refice, Evaluating IPv6 adoption in the Internet, in: Proc. Passive and Active Measurement Conf., in: Lecture Notes on Computer Science, vol. 6032, Springer-Verlag, Heidelberg, 2010, pp. 141–150.
- [114] P. Colp, M. Nanavati, J. Zhu, W. Aiello, G. Coker, T. Deegan, P. Loscocco, A. Warfield, Breaking up is hard to do: security and functionality in a commodity hypervisor, in: Proc. 23rd ACM Symp. Operating Systems Principles, 2011, pp. 189–202.
- [115] F.J. Corbatò, V.A. Vyssotsky, Introduction and overview of the MULTICS system, in: Proc. AFIPS, Fall Joint Computer Conf., 1965, pp. 185–196.
- [116] F.J. Corbatò, On building systems that will fail, Turing Award Lecture 1991, <http://larch-www.lcs.mit.edu:8001/~corbato/turing91/>, 2011.
- [117] J.C. Corbett, J. Dean, M. Epstein, A. Fikes, C. Frost, J.J. Furman, S. Ghemawat, A. Gubarev, C. Heiser, P. Hochschild, W. Hsieh, S. Kanthak, E. Kogan, H. Li, A. Lloyd, S. Melnik, D. Mwaura, D. Nagle, S. Quinlan, R. Rao, L. Rolig, Y. Saito, M. Szymaniak, C. Taylor, R. Wang, D. Woodford, Spanner: Google's globally-distributed database, in: Proc. USENIX Symp. on Operating Systems Design and Implementation, vol. 31(3), 2012, Paper #8.
- [118] P. Cramton, Y. Shoham, R. Steinberg (Eds.), Combinatorial Auctions, MIT Press, 2006.
- [119] F. Cristian, H. Aghili, R. Strong, D. Dolev, Atomic broadcast from simple message diffusion to Byzantine agreement, in: Proc. 15th Int. Symp. Fault Tolerant Computing, IEEE Press, 1985, pp. 200–206, also *Information and Computation* 118 (1) (1995) 158–179.
- [120] Cloud Security Alliance, Security guidance for critical areas of focus in cloud computing V2.1, <https://wikileaks.org/sony/docs/05/docs/Cloud/csaguide.pdf>, 2009.
- [121] Cloud Security Alliance, Top threats to cloud computing V1.0, <https://cloudsecurityalliance.org/toptreats/csathreats.v1.0.pdf>, 2010. (Accessed August 2015).
- [122] Cloud Security Alliance, Security guidance for critical areas of focus in cloud computing V3.0, <https://cloudsecurityalliance.org/guidance/csaguide.v3.0.pdf>, 2011. (Accessed August 2015).
- [123] Cloud Security Alliance, Security guidance for critical areas of focus in cloud computing v4.0, <https://cloudsecurityalliance.org/artifacts/security-guidance-v4/>. (Accessed April 2021).
- [124] E. Cuervo, A. Balasubramanian, D.-K. Cho, A. Wolman, S. Saroiu, R. Chandra, P. Bahl, MAUI: making smartphones last longer with code offload, in: Proc. 8th Int. Conf. on Mobile Systems, Applications, and Services, 2010, pp. 49–62.
- [125] L. Cuen, The debate about cryptocurrency and energy consumption, <https://techcrunch.com/2021/03/21/the-debate-about-cryptocurrency-and-energy-consumption/>, 2021. (Accessed April 2021).
- [126] C. Curino, D.E. Difallah, C. Douglas, S. Krishnan, R. Ramakrishnan, S. Rao, Reservation-based scheduling: if you're late don't blame us!, in: Proc. ACM Symp. Cloud Computing, 2014, pp. 1–14.
- [127] B. Das, Y.Z. Zhang, J. Kiszka, Nested virtualization; state of the art and future directions, in: KVM Forum, 2014, also <http://www.linux-kvm.org/images/3/33/02x03-NestedVirtualization.pdf>. (Accessed January 2017).
- [128] H.A. David, Order Statistics, Wiley, New York, NY, 1981.
- [129] J. Dean, S. Ghernawat, MapReduce: simplified data processing on large clusters, in: Proc. USENIX 6-th Symp. OS Design and Implementation, 2004, pp. 137–149.

- [130] J. Dean, L.A. Barroso, The tail at scale, *Communications of the ACM* 56 (2) (2013) 74–80.
- [131] J. Dean, The rise of cloud computing systems, in: Proc. ACM Symp. OS Principles, 2015, Article No. 12, also <http://sigops.org/sosp/sosp15/history/10-dean-slides.pdf>. (Accessed August 2016).
- [132] J. DeBrabant, A. Pavlo, S. Tu, M. Stonebraker, S. Zdonik, Anti-caching: a new approach to database management system architecture, *Proceedings of the VLDB Endowment* 6 (2013) 1942–1953.
- [133] G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Vosshall, W. Vogels, Dynamo: Amazon's highly available key-value store, in: Proc. 21st ACM/SIGOPS Symp. OS Principles, 2007, pp. 205–220.
- [134] C. Delimitrou, C. Kozyrakis, QoS Aware scheduling for heterogeneous data centers with Paragon, *ACM Transactions on Computer Systems* 31 (4) (2013) 12–24.
- [135] C. Delimitrou, C. Kozyrakis, The Netflix challenge: datacenter edition, *IEEE Computer Architecture Letters* 12 (1) (2013) 29–32.
- [136] C. Delimitrou, C. Kozyrakis, Quasar: resource-efficient and QoS-aware cluster management, in: Proc. ACM Int. Conf. on Architectural Support for Programming Languages and Operating Systems, 2014, pp. 127–144.
- [137] Y. Demchenko, C. de Laat, D.R. Lopes, Security services lifecycle management in on-demand infrastructure services provisioning, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 644–650.
- [138] A. Demers, S. Keshav, S. Shenker, Analysis and simulation of a fair queuing algorithm, in: Proc. ACM SIGCOMM'89 Symp. Communications Architectures and Protocols, 1989, pp. 1–12.
- [139] J.B. Dennis, General parallel computation can be performed with a cycle-free heap, in: Procs. 1998 Int. Conf. on Parallel Architectures and Compiling Techniques, IEEE, 1998, pp. 96–103.
- [140] J.B. Dennis, Fresh breeze: a multiprocessor chip architecture guided by modular programming principles, *ACM SIGARCH Computer Architecture News* 31 (1) (2003) 7–15.
- [141] J.B. Dennis, The fresh breeze model of thread execution, in: Proc. Workshop on Programming Models for Ubiquitous Parallelism, 2006, also <http://csg.csail.mit.edu/Users/dennis/pmup-final.pdf>. (Accessed January 2017).
- [142] D. Denisson, Continuous pipelines at Google, https://www.usenix.org/sites/default/files/continuouspipe_linesatgooglefinal.pdf, 2015. (Accessed May 2016).
- [143] M. Devarakonda, B. Kish, A. Mohindra, Recovery in the Calypso file system, *ACM Transactions on Computer Systems* 14 (3) (1996) 287–310.
- [144] D.J. DeWitt, R.V. Nehme, S. Shankar, J. Aguilar-Saborit, A. Avanes, M. Flasza, J. Gramling, Split query processing in polybase, in: Proc. ACM SIGMOD Int. Conf. on Management of Data, 2013, pp. 1255–1266.
- [145] M.D. Dikaiakos, D. Katsaros, P. Mehra, G. Pallis, A. Vakali, Cloud computing: distributed Internet computing for IT and scientific research, *IEEE Internet Computing* 13 (5) (2009) 10–13.
- [146] E.W. Dijkstra, Cooperating sequential processes, in: F. Genuys (Ed.), *Programming Languages*, Academic Press, New York, 1968, pp. 43–112, originally appeared in 1965, E.W. Dijkstra Archive: Cooperating sequential processes (EWD 123), <https://www.cs.utexas.edu/~EWD/transcriptions/EWD01xx/EWD123.html>.
- [147] E.W. Dijkstra, Self-stabilizing systems in spite of distributed control, *Communications of the ACM* 17 (11) (1974) 643–644.
- [148] H.T. Dinh, C. Lee, D. Niyato, P. Wang, A survey of mobile cloud computing: architecture, applications, and approaches, *Wireless Communications and Mobile Computing* 13 (2013) 1587–1611, <https://doi.org/10.1002/wcm.1203>.
- [149] DONA, Data Oriented Network Architecture, <https://www2.eecs.berkeley.edu/Research/Projects/Data/102146.html>. (Accessed December 2016).
- [150] P. Donnelly, P. Bui, D. Thain, Attaching cloud storage to a campus grid using Parrot, Chirp, and Hadoop, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 488–495.
- [151] T. Dörnemann, E. Juhnke, T. Noll, D. Seiler, B. Freileben, Data flow driven scheduling of BPEL workflows using cloud resources, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 196–203.

- [152] S. Drossopoulou, J. Noble, M.S. Miller, T. Murray, Reasoning about risk and trust in an open word, <http://static.googleusercontent.com/media/research.google.com/en/pubs/archive/44272.pdf>, 2015. (Accessed August 2016).
- [153] L. Ducas, D. Micciancio, Fhew: bootstrapping homomorphic encryption in less than a second, in: Advances in Cryptology, Springer-Verlag, Heidelberg, 2015, pp. 617–640.
- [154] K.J. Duda, R.R. Cheriton, Borrowed-Virtual-Time (BVT) scheduling: supporting latency-sensitive threads in a general-purpose scheduler, in: Proc. 17th ACM Symp. OS Principles, 1999, pp. 261–276.
- [155] N. Dukkipati, T. Refice, Y.-C. Cheng, J. Chu, T. Herbert, A. Agarwal, A. Jain, N. Sutin, An argument for increasing TCP's initial congestion window, ACM SIGCOMM Computer Communication Review 40 (3) (2010) 27–33.
- [156] X. Dutreild, N. Rivierre, A. Moreau, J. Malenfant, I. Truck, From data center resource allocation to control theory and back, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 410–417.
- [157] G. Dyson, Turing's Cathedral: The Origins of the Digital Universe, Pantheon, 2012.
- [158] D.L. Eager, E.D. Lazowska, J. Zahorjan, Adaptive load sharing in homogeneous distributed systems, IEEE Transactions on Software Engineering 12 (1986) 662–675.
- [159] D. Ebneter, S.G. Grivas, T.U. Kumar, H. Wache, Enterprise architecture frameworks for enabling cloud computing, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 542–543.
- [160] V.M. Eguiluz, K. Klemm, Epidemic threshold in structured scale-free networks, arXiv:cond-mat/0205439v1, 2002.
- [161] J. Ejarque, R. Sirvent, R.M. Badia, A multi-agent approach for semantic resource allocation, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 335–342.
- [162] M. Elhawary, Z.J. Haas, Energy-efficient protocol for cooperative networks, IEEE/ACM Transactions on Networking 19 (2) (2011) 561–574.
- [163] M. Eltoweissy, S. Olariu, M. Younis, Towards autonomous vehicular clouds, in: Proc. AdHocNets 2010, August 2010.
- [164] Enterprise Management Associates, How to make the most of cloud computing without sacrificing control, White Paper, prepared for IBM, September 2010, 18 pp., <http://www.siliconrepublic.com/reports/partner/26-ibm/report/311-how-to-make-the-most-of-clo/>. (Accessed August 2015).
- [165] P. Erdős, A. Rényi, On random graphs, Publicationes Mathematicae 6 (1959) 290–297.
- [166] R.M. Esteves, C. Rong, Social impact of privacy in cloud computing, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 593–596.
- [167] C. Evangelinos, C.N. Hill, Cloud computing for parallel scientific HPC applications: feasibility of running coupled atmosphere-ocean climate models on Amazon's EC2, in: Workshop on Cloud Computing and Its Applications, 2008.
- [168] S. Faber, S. Jarecki, H. Krawczyk, Q. Nguyen, M. Rosu, M. Steiner, Rich queries on encrypted data: beyond exact matches, in: Proc. 20th Euro. Symp. Research in Computer Security, in: Lecture Notes on Computer Science, vol. 9327, Springer-Verlag, Heidelberg, 2015, pp. 123–145.
- [169] A.D.H. Farwell, M.J. Sergot, M. Salle, C. Bartolini, D. Tresour, A. Christodoulou, Performance monitoring of service-level agreements for utility computing, in: Proc. IEEE. Int. Workshop Electronic Contracting, 2004, pp. 17–24.
- [170] M. Ferdman, B. Falsafi, A. Adileh, O. Kocberber, S. Volos, M. Alisafaee, D. Jevdjic, C. Kaynak, A.D. Popescu, A. Ailamaki, Clearing the clouds, in: Proc. 17th Int. Conf. on Architectural Support for Programming Languages and Operating Systems, ACM, New York, NY, 2012, pp. 37–48.
- [171] S. Ferretti, V. Ghini, F. Panzieri, M. Pellegrini, E. Turrini, QoS-aware clouds, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 321–328.
- [172] Federal Trade Comision, Privacy online: for information practice in the electronic marketplace, <https://www.ftc.gov/reports/privacy-online-fair-information-practices-electronic-marketplace-federal-trade-commission>, 2000. (Accessed August 2015).

- [173] A. Fikes, Storage architecture and challenges, <https://www.systutorials.com/3306/storage-architecture-and-challenges/>. (Accessed March 2017).
- [174] M. Fingerhuth, T. Babej, P. Wittek, Open source software in quantum computing, *PLoS ONE* 13 (12) (2018), also <https://doi.org/10.1371/journal.pone.0208561>, 2018. (Accessed July 2020).
- [175] M.J. Fischer, N.A. Lynch, M.S. Paterson, Impossibility of distributed consensus with one faulty process, *Journal of the ACM* 32 (2) (1985) 374–382.
- [176] A. Floratou, U.F. Minhas, F. Oğcan, SQLonHadoop: full circle back to shared-nothing database architectures, *Proceedings of the VLDB Endowment* 7 (12) (2014) 1295–1306.
- [177] R. Florin, P. Ghazizadeh, A. Ghazi Zadeh, S. El-Tawab, S. Olariu, Reasoning about job completion time in vehicular clouds, *IEEE Transactions on Intelligent Transportation Systems* 18 (7) (2017) 1762–1771.
- [178] R. Florin, S. Abolghasemi, A. Ghazi Zadeh, S. Olariu, Big data in the parking lot, in: *Big Data Management and Processing*, Taylor and Francis, 2017, pp. 425–449.
- [179] R. Florin, A. Ghazi Zadeh, P. Ghazizadeh, S. Olariu, Towards approximating job completion time in vehicular clouds, *IEEE Transactions on Intelligent Transportation Systems* 20 (7) (2019) 3168–3177.
- [180] R. Florin, A. Ghazi Zadeh, P. Ghazizadeh, S. Olariu, A tight estimate of job completion time in vehicular clouds, *IEEE Transactions on Cloud Computing* 8 (3) (2020) 721–734.
- [181] S. Floyd, Van Jacobson, Link-sharing and resource management models for packet networks, *IEEE/ACM Transactions on Networking* 3 (4) (1995) 365–386.
- [182] B. Ford, Icebergs in the clouds; the other risks of cloud computing, in: Proc. 4th USENIX Workshop Hot Topics in Cloud Computing, 2012, arXiv:1203.1979v2.
- [183] M. Franklin, A. Halevy, D. Maier, From databases to dataspaces: a new abstraction for information management, *SIGMOD Record* 34 (4) (2005) 27–33.
- [184] J. Franklin, K. Bowler, C. Brown, S. Edwards, N. McNab, M. Steele, Mobile device security - cloud and hybrid builds, *NIST Special Publication 1800-4b* (2015).
- [185] E. Gafni, D. Bertsekas, Dynamic control of session input rates in communication networks, *IEEE Transactions on Automatic Control* 29 (10) (1984) 1009–1016.
- [186] G. Ganesan, Y.G. Li, Cooperative spectrum sensing in cognitive radio networks, in: Proc. IEEE Symp. New Frontiers in Dynamic Spectrum Access Networks, 2005, pp. 137–143.
- [187] A.G. Ganek, T.A. Corbi, The dawning of the autonomic computing era, *IBM Systems Journal* 42 (1) (2003) 5–18.
- [188] M.R. Garey, D.S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, W.H. Freeman and Co., 1979.
- [189] S. Garfinkel, M. Rosenblum, When virtual is harder than real: security challenges in virtual machines based computing environments, in: Proc. 10th Conf. Hot Topics in Operating Systems, 2005, pp. 20–25.
- [190] A.F. Gates, O. Natkovich, S. Chopra, P. Kamath, S.M. Narayananamurthy, C. Olston, B. Reed, S. Srinivasan, U. Srivastava, Building a high-level dataflow system on top of MapReduce: the Pig experience, *Proceedings of the VLDB Endowment* 2 (2) (2009) 1414–1425.
- [191] M. Gell-Mann, Simplicity and complexity in the description of nature, *Engineering Sciences* 51 (3) (1988) 2–9, California Institute of Technology, <http://resolver.caltech.edu/CaltechES:51.3.Mann>.
- [192] C. Gentry, A fully homomorphic encryption scheme, Ph.D. Thesis, Stanford University, 2009.
- [193] C. Gentry, Fully homomorphic encryption using ideal lattices, in: Proc. 41st ACM Symp. Theory of Computing, 2009, pp. 169–178.
- [194] M. Gerla, Vehicular cloud computing, in: Proc. 1st Int. Workshop Vehicular Communication and Applications, IEEE, 2012, pp. 152–155, published in 11th Med. Ad Hoc Networking Workshop.
- [195] C. Gkantsidis, M. Mihail, A. Saberi, Random walks in peer-to-peer networks, *Performance Evaluation* 63 (3) (2006) 241–263.
- [196] P. Ghazizadeh, R. Florin, R.A. Ghazi Zadeh, S. Olariu, Reasoning about the mean time to failure in vehicular clouds, *IEEE Transactions on Intelligent Transportation Systems* 3 (17) (2016) 751–761.

- [197] S. Ghemawat, H. Gobioff, S.-T. Leung, The Google file system, in: Proc. 19th ACM Symp. OS Principles, SOSP 03, 2003, pp. 15–25.
- [198] S. Gilbert, N. Lynch, Perspectives on the CAP theorem, <https://groups.csail.mit.edu/tds/papers/Gilbert/Brewer2.pdf>. (Accessed April 2021).
- [199] D. Gmach, J. Rolia, L. Cerkasova, Satisfying service-level objectives in a self-managed resource pool, in: Proc. 3rd. Int. Conf. Self-Adaptive and Self-Organizing Systems, 2009, pp. 243–253.
- [200] R.P. Goldberg, Architectural principles for virtual computer systems, Thesis, Harvard University, 1973.
- [201] R.P. Goldberg, Survey of virtual machine research, IEEE Computer Magazine 7 (6) (1974) 34–45.
- [202] G. Gonnet, Expected length of the longest probe sequence in hash code searching, Journal of the ACM 28 (2) (1981) 289–304.
- [203] I. Goodfellow, Y. Bengio, A. Courville, Deep Learning, MIT Press, Boston, Mass, 2016.
- [204] Google, Google cloud platform, <https://cloud.google.com/>. (Accessed August 2016).
- [205] P. Goyal, X. Guo, H.M. Vin, A hierachial CPU scheduler for multimedia operating systems, in: Proc. USENIX Symp. OS Design and Implementation, 1996, pp. 107–121.
- [206] J. Gray, The transaction concept: virtues and limitations, in: Proc. 7th Int. Conf. Very Large Databases, 1981, pp. 144–154.
- [207] J. Gray, D. Patterson, A conversation with Jim Gray, ACM Queue 1 (4) (2003) 8–17.
- [208] G. Graefe, H. Volos, H. Kimura, H. Kuno, J. Tucek, In memory performance for Big Data, Proceedings of the VLDB Endowment 8 (1) (2014) 37–48.
- [209] T.G. Griffn, F.B. Shepherd, G. Wilfong, The stable paths problem and interdomain routing, IEEE/ACM Transactions on Networking 10 (2) (2002) 232–243.
- [210] M. Gittert, D.R. Cheriton, An architecture for content routing support in the Internet, in: Proc. USENIX Symp. on Internet Technologies and Systems, vol. 3, 2001, p. 4.
- [211] W. Gropp, E. Lusk, A. Skjellum, Using MPI, MIT Press, 1994.
- [212] N. Gruschka, M. Jensen, Attack surfaces: a taxonomy for attacks on cloud services, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 276–279.
- [213] T. Gunaratne, T.-L. Wu, J. Qiu, G. Fox, MapReduce in the clouds for science, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 565–572.
- [214] P.K. Gunda, L. Ravindranath, C.A. Thekkath, Y. Yu, L. Zhuang, Nectar: automatic management of data and computation in data centers, in: Proc. USENIX Symp. OS Design and Implementation, 2010, pp. 75–88.
- [215] I. Gupta, A.J. Ganesh, A.-M. Kermarrec, Efficient and adaptive epidemic-style protocols for reliable and scalable multicast, IEEE Transactions on Parallel and Distributed Systems 17 (7) (2006) 593–605.
- [216] V. Gupta, M. Harchol-Balter, Self-adaptive admission control policies for resource-sharing systems, in: Proc. 11th Int. Joint Conf. Measurement and Modeling Computer Systems, 2009, pp. 311–322.
- [217] A. Gupta, F. Yang, J. Govig, A. Kirsch, K. Chan, K. Lai, S. Wu, S.G. Dhoot, A.R. Kumar, A. Agiwal, S. Bhansali, M. Hong, J. Cameron, M. Siddiqi, D. Jones, J. Shute, A. Gubarev, S. Venkataraman, D. Agrawal, Mesa: geo-replicated, near real-time, scalable data warehousing, Proceedings of the VLDB Endowment 7 (12) (2014) 1259–1270.
- [218] A. Gupta, J. Shute, High-availability at massive scale: building Google’s data infrastructure for ads, in: Proc. 9th Workshop Business Intelligence for the Real Time Enterprise, Springer-Verlag, Heidelberg, 2015, pp. 81–89.
- [219] J.O. Gutierrez-Garcia, K.-M. Sim, Self-organizing agents for service composition in cloud computing, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 59–66.
- [220] P.J. Haas, Hoeffding inequalities for join-selectivity estimation and online aggregation, IBM Research Report RJ 10040 (90568), IBM Almaden Research, 1996.
- [221] T. Haig, M. Pристley, C. Rope, Los Alamos bets on ENIAC: nuclear Monte Carlo simulations, 1947–1948, IEEE Annals of the History of Computing 36 (3) (2014) 42–63.

- [222] F.A. Halderman, S.D. Schoen, N. Heninger, W. Clarkson, W. Paul, J.A. Calandino, A.J. Feldman, J. Appelbaum, E.W. Felten, Lest we remember: cold boot attacks on encryption keys, in: Proc. 17th Usenix, Security Symposium, 2008, pp. 45–60.
- [223] S. Halevi, V. Shoup, Algorithms in HElib, in: J.A. Garay, R. Gennaro (Eds.), Advances in Cryptology, in: Lecture Notes on Computer Science, vols. 8616 and 8617, Springer-Verlag, Heidelberg, 2014, pp. 554–571.
- [224] J.D. Halley, D.A. Winkler, Classification of emergence and its relation to self-organization, Complexity 13 (5) (2008) 10–15.
- [225] P.B. Hansen, The evolution of operating systems, in: P.B. Hansen (Ed.), Classic Operating Systems: from Batch Processing to Distributed Systems, Springer-Verlag, Heidelberg, 2001, pp. 1–36.
- [226] R.F. Hartl, S.P. Sethi, R.G. Vickson, Survey of the maximum principles for optimal control problems with state constraints, SIAM Review 37 (2) (1995) 181–218.
- [227] T. Härder, A. Reuter, Principles of transaction-oriented database recovery, ACM Computing Surveys 15 (4) (1983) 287–317.
- [228] S. Harizopoulos, D.J. Abadi, S. Madden, M. Stonebreaker, OLTP through the looking glass, and what we found there, in: Proc. ACM/SIGMOD Int. Conf. on Management of Data, 2008, pp. 981–992.
- [229] K. Hasebe, T. Niwa, A. Sugiki, K. Kato, Power-saving in large-scale storage systems with data migration, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 266–273.
- [230] R.L. Haskin, Tiger Shark - a scalable file system for multimedia, IBM Journal of Research and Development 42 (2) (1998) 185–197.
- [231] J.L. Hellerstein, Why feedback implementations fail: the importance of systematic testing, in: 5th Int. Workshop on Feedback Control Implementation and Design in Computing Systems and Networks, 2015.
- [232] J.L. Hennessy, D.A. Patterson, Computer Architecture - a Quantitative Approach, sixth edition, Morgan Kaufmann, Waltham, MA, 2017.
- [233] M. Herlihy, N. Shavit, The Art of Multiprocessor Programming, revised reprint, Morgan Kaufmann, Waltham, MA, 2012.
- [234] H. Herodotou, H. Lim, G. Luo, N. Borisov, L. Dong, F.B. Cetin, S. Babu, Starfish: a self-tuning system for big data analytics, in: Proc. 5th Conf. on Innovative Data Systems Research, 2011, pp. 261–272.
- [235] T. Hey, S. Tansley, K. Tolle, Jim Gray on eScience: a transformed scientific method, in: The Fourth Paradigm. Data-Intensive Scientific Discovery, Microsoft Research, 2009.
- [236] M. Hilbert, P. López, The world's technological capacity to store, communicate, and compute information, Science 332 (6025) (2011) 60–65.
- [237] M. Hill, S. Eggers, J.R. Larus, G. Taylor, G. Adams, B.K. Bose, G. Gibson, P. Hansen, J. Keller, S. Kong, C. Lee, D. Lee, J. Pendleton, S. Richie, D. Wood, B. Zorn, P. Hilfinger, D. Hodges, R. Katz, J.K. Osterhout, J. Petterson, Design decisions in SPUR, Computer 9 (11) (1986) 8–22.
- [238] M.D. Hill, M.R. Marty, Amdahl's law in the multicore era, IEEE Computer 41 (7) (2008) 33–38.
- [239] M. Hinckey, R. Sterritt, C. Rouff, J. Rash, W. Truszkowski, Swarm-based space exploration, ERCIM News 64 (2006).
- [240] B. Hindman, A. Konwinski, M. Zaharia, A. Ghodsi, A.D. Joseph, R. Katz, S. Shenker, I. Stoica, Mesos: a platform for fine-grained resource sharing in the data center, in: Proc. 8th USENIX Symp. Networked Systems Design and Implementation, 2011, pp. 295–308.
- [241] C.A.R. Hoare, Communicating sequential processes, Communications of the ACM 21 (8) (1978) 666–677.
- [242] U. Hölzle, Brawny cores still beat wimpy cores, most of the time, IEEE MICRO 30 (4) (2010) 3.
- [243] J. Hopfield, Neural networks and physical systems with emergent collective computational abilities, Proceedings of the National Academy of Sciences 79 (1982) 2554–2558.
- [244] C. Hopps, Analysis of an equal-cost multi-path algorithm, in: RFC 2992, Internet Engineering Task Force, 2000.
- [245] J.H. Howard, M.L. Kazer, S.G. Menees, D.A. Nichols, M. Satyanarayanan, R.N. Sidebotham, M.J. West, Scale and performance in a distributed file system, ACM Transactions on Computer Systems 6 (1) (1988) 51–81.

- [246] D.H. Hu, Y. Wang, C.-L. Wang, BetterLife 2.0: large-scale social intelligence reasoning on cloud, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 529–536.
- [247] C. Huang, R. Lu, H. Zhu, H. Hu, X. Lin, PTVC: achieving privacy-preserving trust-based verifiable vehicular cloud computing, in: Proc. IEEE Global Comm. Conf., GLOBECOM 2016, 2016, pp. 1–6.
- [248] G. Iachello, J. Hong, End-user privacy in human-computer interaction, Foundations and Trends in Human-Computer Interactions 1 (1) (2007) 1–137.
- [249] IBM Smart Business, Dispelling the vapor around the cloud computing. Drivers, barriers and considerations for public and private cloud adoption, White Paper, 2010, <ftp://software.ibm.com/common/ssi/ecm/en/ciw03062usen/CIW03062USEN.PDF>. (Accessed August 2015).
- [250] IBM Corporation, General parallel file systems (version 3, update 4). Documentation Updates, <http://publib.boulder.ibm.com/infocenter/clresctr/vxrx/topic/com.ibm.cluster.gpfs.doc>. (Accessed August 2015).
- [251] IBM Corporation, The evolution of storage systems, IBM Research, Almaden Research Center Publications, <http://www.almaden.ibm.com/storagesystems/pubs>. (Accessed August 2015).
- [252] IBM Corporation, Bringing big data to the enterprise, <https://www-01.ibm.com/software/in/data/bigdata/>, 2015. (Accessed May 2016).
- [253] Intel Corporation, Intel 64 and IA-32 architectures software developer manuals, <https://software.intel.com/en-us/node/699529#combined>, 2016. (Accessed January 2017).
- [254] M. Iorga, A. Karmel, Managing risk in a cloud ecosystem, IEEE Cloud Computing 2 (6) (2015) 51–57.
- [255] M. Isard, M. Budiu, Y. Yu, A. Birrell, D. Fetterly, Dryad: distributed data-parallel programs from sequential building blocks, in: Proc. 2nd ACM SIGOPS/EuroSys European Conf. Computer Systems, 2007, pp. 57–62.
- [256] K.R. Jackson, L. Ramakrishnan, K. Muriki, S. Canon, S. Cholia, J. Shalf, H. Wasserman, N.J. Wright, Performance analysis of high performance computing applications on the Amazon Web Services cloud, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 159–168.
- [257] D. Jaggar, A history of the ARM microprocessor, Google Talks, https://www.youtube.com/results?search_query=David+Jaggar, 2019.
- [258] M. Jelasity, A. Montresor, O. Babaoglu, Gossip-based aggregation in large dynamic networks, ACM Transactions on Computer Systems 23 (3) (2005) 219–252.
- [259] M. Jelasity, S. Voulgaris, R. Guerraoui, A.-M. Kermarrec, M. van Steen, Gossip-based peer sampling, ACM Transactions on Computer Systems 25 (3) (2007) 8.
- [260] M. Jensen, S. Schäge, J. Schwenk, Towards an anonymous access control and accountability scheme for cloud computing, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 540–541.
- [261] H. Jin, X.-H. Sun, Y. Chen, T. Ke, REMEM: REmote MEMory as checkpointing storage, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 319–326.
- [262] N.P. Jouppi, et al., In-datacenter performance analysis of a tensor processing unit, <https://arxiv.org/pdf/1704.04760.pdf>, 2017.
- [263] R. Kallman, H. Kimura, J. Atkins, A. Pavlo, A. Rasin, S. Zdonik, E. Jones, Y. Zhang, S. Madden, M. Stonebraker, J. Hugg, D.J. Abadi, H-Store: a high-performance, distributed main memory transaction processing system, Proceedings of the VLDB Endowment 1 (2) (2008) 1496–1499.
- [264] E. Kalyvianaki, T. Charalambous, S. Hand, Applying Kalman filters to dynamic resource provisioning of virtualized server applications, in: Proc. 3rd Int. Workshop Feedback Control Implementation & Design in Computing Systems & Networks, 2008, p. 6.
- [265] S. Kanev, J.P. Darago, K. Hazelwood, P. Ranganathan, T. Moseley, G-Y. Wei, D. Brooks, Profiling a warehouse-scale computer, in: Proc. 42nd Annual Int. Symp. Computer Architecture, ISCA, 2015, pp. 158–169.
- [266] R.M. Karp, M. Luby, F. Meyer auf der Heide, Efficient PRAM simulation on a distributed memory machine, in: Proc. 24th ACM Symp. on the Theory of Computing, 1992, pp. 318–326.
- [267] K. Kc, K. Anyanwu, Scheduling Hadoop jobs to meet deadlines, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 388–392.

- [268] J. Kephart, H. Chan, R. Das, D. Levine, G. Tesauro, F. Rawson, C. Lefurgy, Coordinating multiple autonomic managers to achieve specified power-performance tradeoffs, in: Proc. 4th Int. Conf. Autonomic Computing, 2007, pp. 100–109.
- [269] J. Kephart, The utility of utility: policies for autonomic computing, in: Proc. LCCC Workshop Control of Computing Systems, 2011.
- [270] W.O. Kermack, A.G. McKendrick, A contribution to the theory of epidemics, *Proceedings of the Royal Society of London. Series A* 115 (1927) 700–721.
- [271] W. Kim, Introduction to Object-Oriented Databases, MIT Press, 1990.
- [272] J. Kim, W.J. Dally, D. Abts, Flattened butterfly: a cost-efficient topology for high-radix networks, in: Proc. 34th Int. Symp. Computer Architecture, 2007, pp. 126–137.
- [273] J. Kim, W.J. Dally, D. Abts, Efficient topologies for large-scale cluster networks, in: Proc. 2010 Optical Fiber Communication Conf. and National Fiber Optic Engineers Conf., 2010, pp. 1–3.
- [274] S.T. King, P.M. Chen, Y.-M. Wang, C. Verbowski, H.J. Wang, J.R. Lorch, SubVirt: implementing malware with virtual machines, in: Proc. IEEE Symp. Security and Privacy, 2006, pp. 314–327.
- [275] A. Kleiner, A. Talwalkar, S. Agarwal, I. Stoica, M.I. Jordan, A general bootstrap performance diagnostic, in: Proc. 19th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining, ACM, New York, NY, 2013, pp. 419–427.
- [276] L. Kleinrock, Queuing Systems, vols. I and II, Wiley, New York, NY, 1965.
- [277] D.E. Knuth, The Art of Computer Programming I–III, 2nd edition, Addison–Wesley, Reading, Ma, 1973.
- [278] F. Koeppe, J. Schneider, Do you get what you pay for? Using proof-of-work functions to verify performance assertions in the cloud, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 687–692.
- [279] B. Koley, V. Vusirikala, C. Lam, V. Gill, 100Gb Ethernet and beyond for warehouse scale computing, in: Proc. 15th OptoElectronics and Communications Conf., 2010, pp. 106–107.
- [280] J.G. Koomey, S. Berard, M. Sanchez, H. Wong, Implications of historical trends in the energy efficiency of computing, *IEEE Annals of the History of Computing* 33 (3) (2011) 46–54.
- [281] H. Kopetz, Sparse time versus dense time in distributed real-time systems, in: Proc. 12 Int. Conf. on Distributed Computing Systems, IEEE Computer Society, 1992, pp. 460–467.
- [282] M. Kornacker, A. Behm, V. Bittorf, T. Bobrovitsky, A. Choi, J. Erickson, M. Grund, D. Hecht, M. Jacobs, I. Joshi, L. Kuff, D. Kuma, A. Leblang, N. Li, H. Robinson, D. Rorke, S. Rus, J. Russell, D. Tsirogiannis, S. Wanderman-Milne, M. Yoder, Impala: a modern, open-source SQL engine for Hadoop, in: Proc 7th Biennial Conf. on Innovative Data Systems Research, 2015 (online proceedings).
- [283] G. Koslovski, W.-L. Yeow, C. Westphal, T.T. Huu, J. Montagnat, P. Vicat-Blanc, Reliability support in virtual infrastructures, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 49–58.
- [284] P.R. Krugman, The Self-Organizing Economy, Blackwell Publishers, 1996.
- [285] J.F. Kurose, K.W. Ross, Computer Networking: A Top-down Approach, 6th edition, Addison–Wesley, Reading, Ma, 2013.
- [286] S. Kumar, L. Shi, N. Ahmed, S. Gil, D. Katabi, D. Rus, CarSpeak: a content-centric network for autonomous driving, *ACM SIGCOMM Computer Communication Review* 42 (4) (2012) 259–270.
- [287] D. Kusic, J.O. Kephart, N. Kandasamy, G. Jiang, Power and performance management of virtualized computing environments via lookahead control, in: Proc. 5th Int. Conf. Autonomic Computing, 2008, pp. 3–12.
- [288] A. Kivity, Y. Kamay, D. Laor, U. Lublin, A. Liguori, KVM: the Linux virtual machine monitor, in: Proc. Linux Symposium, Ottawa, 2007, pp. 225–230.
- [289] B.M. Leiner, V.G. Cerf, D.D. Clark, R.E. Khan, L. Kleinrock, D.C. Lynch, J. Postel, L.G. Roberts, S. Wolff, A brief history of the Internet, *ACM SIGCOMM Computer Communication Review* 1 (5) (2009) 22–39.
- [290] C.F. Lam, FTTH look ahead - technologies and architectures, in: Proc. 36th European Conf. Optical Communications, ECOC’10, 2010, pp. 1–18.
- [291] L. Lamport, P.M. Melliar-Smith, Synchronizing clocks in the presence of faults, *Journal of the ACM* 32 (1) (1985) 52–78.

610 Literature

- [292] L. Lamport, The part-time Parliament, *ACM Transactions on Computer Systems* 2 (1998) 133–169.
- [293] L. Lamport, Paxos made simple, *ACM SIGACT News* 32 (4) (2001) 51–58.
- [294] L. Lamport, Turing Lecture - the computer science of concurrency: the early years, *Communications of the ACM* 58 (6) (2015) 71–77.
- [295] B.W. Lampson, H.E. Sturmfis, Reflections on operating system design, *Communications of the ACM* 19 (5) (1976) 251–265.
- [296] P.A. Lascocco, S.D. Smalley, P.A. Muckelbauer, R.C. Taylor, S.J. Turner, J.F. Farrell, The inevitability of failure: the flawed assumption of security in modern computing environments, in: Proc. 21st Conf. on National Information System Security, 1998, pp. 303–314.
- [297] A.M. Law, W.D. Kelton, *Simulation Modeling and Analysis*, Mc Graw-Hill, New York, NY, 1982.
- [298] D. Lea, *Concurrent Programming in Java*, second edition, Addison-Wesley, Reading, Ma, 1999.
- [299] P.J. Leach, P. Levine, B. Douros, J. Hamilton, D. Nelson, B. Stumpf, The architecture of an integrated local area network, *IEEE Transactions on Selected Areas in Communication* 1 (5) (1983) 842–857.
- [300] E. Lee, E-K. Lee, M. Gerla, Vehicular cloud networking: architecture and design principles, *IEEE Communications Magazine* 52 (2) (2014) 142–155.
- [301] E. Le Sueur, G. Heiser, Dynamic voltage and frequency scaling: the laws of diminishing returns, in: Proc. Workshop Power Aware Computing and Systems, 2010, pp. 2–5.
- [302] H. Li, A. Ghodsi, M. Zaharia, S. Shenker, I. Stoica, Tachyon: reliable, memory speed storage for cluster computing frameworks, in: Proc. ACM Sym. Cloud Computing, 2014, pp. 1–15.
- [303] Z. Li, N.-H. Yu, Z. Hao, A novel parallel traffic control mechanism for cloud computing, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 376–382.
- [304] C. Li, A. Raghunathan, N.K. Jha, Secure virtual machine execution under an untrusted management OS, in: Proc. IEEE 3rd Int. Conf. Cloud Computing, 2010, pp. 172–179.
- [305] X. Li, J.B. Dennis, G.R. Gao, W. Lim, H. Wei, C. Yang, R. Pavel, FreshBreeze: a data flow approach for meeting DDDAS challenges, in: Proc. Int. Conf. on Computational Science, ICCS 2015, 2015, pp. 2573–2584.
- [306] H.C. Lim, S. Babu, J.S. Chase, S.S. Parekh, Automated control in cloud computing: challenges and opportunities, in: Proc. First Workshop Automated Control for Data Centers and Clouds, ACM Press, 2009, pp. 13–18.
- [307] N.M. Linke, D. Maslov, M. Roetteler, S. Debnath, C. Figgatt, K.A. Landsman, K. Wright, C. Monroe, Experimental comparison of two quantum computing architectures, *Proceedings of the National Academy of Sciences of the United States of America* 114 (2017) 3305–3310, also arXiv:1702.01852 [quant-ph].
- [308] X. Lin, Y. Lu, J. Deogun, S. Goddard, Real-time divisible load scheduling for cluster computing, in: Proc. 13th IEEE Real-Time and Embedded Technology and Applications Symp., 2007, pp. 303–314.
- [309] C. Lin, D.C. Marinescu, Stochastic high level Petri Nets and applications, *IEEE Transactions on Computers* C-37 (7) (1988) 815–825.
- [310] S. Liu, G. Quan, S. Ren, On-line scheduling of real-time services for cloud computing, in: Proc. 6th World Congress on Services, IEEE, 2010, pp. 459–464.
- [311] D. Lo, L. Cheng, R. Govindaraju, L-A. Barroso, C. Kozyrakis, Towards energy proportionality for large-scale latency-critical workloads, *Proceedings of the ACM SIGARCH Computer Architecture News* 42 (3) (2014) 301–312.
- [312] D. Lo, L. Cheng, R. Govindaraju, P. Ranganathan, Kozyrakis, Heracles: improving resource efficiency at scale, in: Proc. 42nd Annual Int. Symp. Computer Architecture, 2015, pp. 450–462.
- [313] G.K. Lockwood, Quick MPI cluster setup on Amazon EC2, <https://glenenklockwood.blogspot.com/2013/04/quick-mpi-cluster-setup-on-amazon-ec2.html>. (Accessed January 2017).
- [314] N. Loutas, V. Peristeras, T. Bouras, E. Kamateri, D. Zeginis, K. Tarabanis, Towards a reference architecture for semantically interoperable clouds, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 143–150.

- [315] C. Lu, J. Stankovic, G. Tao, S. Son, Feedback control real-time scheduling: framework, modeling and algorithms, *Journal of Real-Time Systems* 23 (1–2) (2002) 85–126.
- [316] W. Lu, J. Jackson, J. Ekanayake, R.S. Barga, N. Araujo, Performing large science experiments on Azure: pitfalls and solutions, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 209–217.
- [317] A. Luckow, S. Jha, Abstractions for loosely-coupled and ensemble-based simulations on Azure, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 550–556.
- [318] F. Lumineau, W. Wang, O. Schilke, Blockchain governance - a new way of organizing collaborations?, *Organization Science* 32 (2) (2020) 500–521.
- [319] J. Luna, N. Suri, M. Iorga, A. Karmel, Leveraging the potential of cloud security service level agreements through standards, *IEEE Cloud Computing* 2 (3) (2015) 32–40.
- [320] W.-Y. Ma, B. Shen, J. Brassil, Content service networks: the architecture and protocols, in: A. Bestavros, M. Rabinovich (Eds.), *Web Caching and Content Delivery*, Elsevier, 2001, pp. 83–101.
- [321] J. Madhavan, A. Halevy, S. Cohen, X. Dong, S.R. Jeffery, D. Ko, C. Yu, Structured data meets the web: a few observations, *IEEE Data Engineering Bulletin* 29 (3) (2006) 19–26.
- [322] D.J. Magenheimer, T.W. Christian, vBlades: optimized paravirtualization for the Itanium processor family, in: Proc. 3rd VM Research and Technology Workshop, San Jose, Ca, 2004, pp. 73–82.
- [323] S. Majumder, S. Rixner, Comparing Ethernet and Myrinet for MPI communication, in: Proc. 7th Workshop on Languages, Compilers, and Run-Time Support for Scalable Systems, 2004, pp. 1–7.
- [324] N. Malviya, A. Weisberg, S. Madden, M. Stonebraker, Rethinking main memory OLTP recovery, in: Proc. IEEE 30th Int. Conf. on Data Engineering, 2014, pp. 604–615.
- [325] D.C. Marinescu, *Internet-Based Workflow Management*, Wiley, New York, NY, 2002.
- [326] D.C. Marinescu, G.M. Marinescu, *Approaching Quantum Computing*, Prentice Hall, Upper Saddle River, NJ, 2004.
- [327] D.C. Marinescu, H.J. Siegel, J.P. Morrison, Options and commodity markets for computing resources, in: R. Buyya, K. Bubendorf (Eds.), *Market Oriented Grid and Utility Computing*, Wiley, New York, NY, ISBN 9780470287682, 2009, pp. 89–120.
- [328] D.C. Marinescu, C. Yu, G.M. Marinescu, Scale-free, self-organizing very large sensor networks, *Journal of Parallel and Distributed Computing* 50 (5) (2010) 612–622.
- [329] D.C. Marinescu, G.M. Marinescu, *Classical and Quantum Information*, Academic Press, Amsterdam, 2012.
- [330] D.C. Marinescu, A. Paya, J.P. Morrison, P. Healy, Distributed hierarchical control versus an economic model for cloud resource management, <http://arXiv.org/pdf/1503.01061.pdf>, 2015.
- [331] D.C. Marinescu, Cloud energy consumption, in: S. Muguresan, I. Bojanova (Eds.), *Encyclopedia of Cloud Computing*, Wiley, New York, NY, 2016, Chapter 25.
- [332] D.C. Marinescu, *Complex Systems and Clouds: A Self-Organization and Self-Management Perspective*, Morgan Kaufmann, Waltham, MA, 2016.
- [333] D.C. Marinescu, A. Paya, J.P. Morrison, A cloud reservation system for big data applications, *IEEE Transactions on Parallel and Distributed Computing* 28 (3) (2017) 606–618.
- [334] D.C. Marinescu, A. Paya, J.P. Morrison, S. Olariu, An approach for scaling cloud resource management, *Cluster Computing* 20 (1) (2017) 909–924, Springer Verlag, Heidelberg.
- [335] P. Marshall, K. Keahey, T. Freeman, Elastic site: using clouds to elastically extend site resources, in: Proc. IEEE Int. Symp. Cluster Computing and the Grid, 2010, pp. 43–52.
- [336] E. Masanet, A. Shehabi, N. Lei, S. Smith, J. Koomey, Recalibrating global data center energy-use estimates, *Science* 367 (6481) (2020) 984–986.
- [337] L. Mashayekhy, M.M. Nejad, D. Grosu, Cloud federations in the sky: formation game and mechanisms, *IEEE Transactions on Cloud Computing* 3 (1) (2015) 14–27.
- [338] F. Mattern, Virtual time and global states of distributed systems, in: Proc. Int. Workshop Parallel and Distributed Algorithms, Elsevier, New York, 1989, pp. 215–226.

- [339] J.M. May, Parallel I/O for High Performance Computing, Morgan Kaufmann, Waltham, MA, 2000.
- [340] M.W. Mayer, Architecting principles for system of systems, *Systems Engineering* 1 (4) (1998) 267–274.
- [341] S. McCartney, ENIAC; the Triumphs and Tragedies of the World's First Computer, Walker and Company Publishing House, New York, NY, 1999.
- [342] W. McCulloch, W. Pitts, A logical calculus of ideas immanent in nervous activity, *Bulletin of Mathematical Biophysics* 5 (4) (1943) 115–133.
- [343] P. McKenney, On the efficient implementation of fair queuing, *Internetworking: Research and Experience* 2 (1991) 113–131.
- [344] P. Mell, What is special about cloud security?, *IT Professional* 14 (4) (2012) 6–8, <http://doi.ieeecomputersociety.org/10.1109/MITP.2012.84>. (Accessed August 2015).
- [345] A. Menon, J.R. Santos, Y. Turner, G.J. Janakiraman, W. Zwaenepoel, Diagnosing performance overheads in Xen virtual machine environments, in: Proc. 1st ACM/USENIX Conf. Virtual Execution Environments, 2005.
- [346] A. Menon, A.L. Cox, W. Zwaenepoel, Optimizing network virtualization in Xen, in: Proc. USENIX Annual Technical Conf., 2006, pp. 15–28.
- [347] A.P. Miettinen, J.K. Miettinen, Energy efficiency of mobile clients in cloud computing, in: Proc. 2nd USENIX Conf. on Hot Topics in Cloud Computing, 2010, pp. 4–11.
- [348] R.A. Milner, *A Calculus of Communicating Systems*, Lecture Notes in Computer Science, vol. 92, Springer-Verlag, Heidelberg, 1980.
- [349] R.A. Milner, Lectures on a calculus for communicating systems, in: Proc. Seminar on Concurrency, in: Lecture Notes in Computer Science, vol. 197, Springer-Verlag, Heidelberg, 1984, pp. 197–220.
- [350] M. Miranda, When every atom counts, *IEEE Spectrum* 49 (7) (2012) 32–37.
- [351] J. Mitola, G.Q. Maguire, Cognitive radio: making software radios more personal, *IEEE Personal Communications* 6 (1999) 13–18.
- [352] J. Mitola, Cognitive radio: an integrated agent architecture for software defined radio, Ph.D. Thesis, KTH, Stockholm, 2000.
- [353] M. Mitzenmacher, The power of two choices in randomized load balancing, Ph.D. Dissertation, Computer Science Department, University of California at Berkeley, 1996.
- [354] M. Mitzenmacher, A.W. Richa, R. Sitaraman, The power of two random choices: a survey of techniques and results, in: S. Rajasekaran, P.M. Pardalos, J.H. Reif, J. Rolim (Eds.), *Handbook of Randomized Computing*, Kluwer Academic Publishers, 2001, pp. 255–312.
- [355] M. Mitzenmacher, The power of two choices in randomized load balancing, *IEEE Transactions on Parallel and Distributed Systems* 12 (10) (2001) 1094–1104.
- [356] T. Miyamoto, M. Hayashi, K. Nishimura, Sustainable network resource management system for virtual private clouds, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 512–520.
- [357] A. Mondal, S.K. Madria, M. Kitsuregawa, Abide: a bid-based economic incentive model for enticing non-cooperative peers in mobile P2P networks, in: Proc. 12th Int. Conf. Database Systems for Advanced Applications, 2007, pp. 703–714.
- [358] J.H. Morris, M. Satyanarayanan, M.H. Conner, M.H. Howard, D.S. Rosenthal, F.D. Smith, Andrew: a distributed personal computing environment, *Communications of the ACM* 29 (3) (1986) 184–201.
- [359] R.J.T. Morris, B.J. Truskowski, The evolution of storage systems, *IBM Systems Journal* 42 (2) (2003) 205–217.
- [360] J. Nagle, On packet switches with infinite storage, *IEEE Transactions on Communications* 35 (4) (1987) 435–438.
- [361] G. Neiger, A. Santoni, F. Leung, D. Rodgers, R. Uhlig, Intel virtualization technology: hardware support for efficient processor virtualization, *Intel Technology Journal* 10 (3) (2006) 167–177.
- [362] M. Nelson, Virtual memory for the Sprite operating system, Technical Report UCB/CSD 86/301, Computer Science Division (EECS), University of California, Berkeley, 1986.

- [363] M.N. Nelson, B.B. Welch, J.K. Osterhout, Caching in Sprite network file systems, ACM Transactions on Computer Systems 6 (1) (1988) 134–154.
- [364] A.J. Nicholson, S. Wolchok, B.D. Noble, Juggler: virtual networks for fun and profit, IEEE Transactions on Mobile Computing 9 (1) (2010) 31–43.
- [365] H. Nissenbaum, Can trust be secured online? A theoretical perspective, Etica e Politica I (2) (1999), 24 pp., Edizione Universita di Trieste, <http://hdl.handle.net/10077/5544>.
- [366] NIST–Reference Architecture Analysis Team, Cloud computing reference architecture - Strawman model V2, Document NIST CCRATWG 019, 2011, 28 pp., <http://www.ogf.org/pipermail/occi-wg/attachments/20110303/e63dee43/attachment-0001.pdf>. (Accessed February 2017).
- [367] Y. Nuevo, Cellular phones as embedded systems, in: Digest of Technical Papers, IEEE Solid-State Circuits Conference, 2004, pp. 32–37.
- [368] D. Nurmi, R. Wolski, C. Grzegorczyk, G. Obertelli, S. Soman, L. Youseff, D. Zagorodnov, The Eucalyptus open-source cloud-computing system, in: Proc. 9th IEEE/ACM Int Symp. Cluster Computing and the Grid, 2009, pp. 124–131.
- [369] S. Oikawa, R. Rajkumar, Portable RK: a portable resource kernel for guaranteed and enforced timing behavior, in: Proc. IEEE Real Time Technology and Applications Symp., June 1999, pp. 111–120.
- [370] T. Okuda, E. Kawai, S. Yamaguchi, A mechanism of flexible memory exchange in cloud computing environments, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 75–80.
- [371] S. Olariu, T. Hristov, G. Yan, The next paradigm shift: from vehicular networks to vehicular clouds, in: S. Basagni, et al. (Eds.), Mobile Ad Hoc Networking Cutting Edge Directions, second edition, Wiley and Sons, N.Y., 2013, pp. 645–700.
- [372] S. Olariu, R. Florin, Vehicular cloud research: what is missing?, in: Proc. 7th ACM Int. Symp. on Design and Analysis of Intelligent Vehicular Networks and Applications, DiVANET’2017, Nov. 2017.
- [373] S. Olariu, A survey of vehicular cloud research: trends, applications, and challenges, IEEE Transactions on Intelligent Transportation Systems 21 (6) (2020) 2648–2663.
- [374] D. Olmedilla, Security and privacy on the semantic web, in: M. Petkovic, W. Jonker (Eds.), Security, Privacy and Trust in Modern Data Management, Springer-Verlag, Heidelberg, 2006.
- [375] M. O'Neill, SaaS, PaaS, and IaaS: a security checklist for cloud models, <http://www.csoneonline.com/article/660065/saas-paas-and-iaas-a-security-checklist-for-cloud-models>. (Accessed August 2015).
- [376] OpenVZ, <http://wiki.openvz.org>. (Accessed August 2015).
- [377] A.M. Oprescu, T. Kielmann, Bag-of-tasks scheduling under budget constraints, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 351–359.
- [378] Oracle Corporation, Lustre file system, [http://en.wikipedia.org/wiki/Lustre_\(file_system\)](http://en.wikipedia.org/wiki/Lustre_(file_system)), 2010. (Accessed February 2017).
- [379] Oracle Corporation, Oracle NoSQL Database, <http://www.oracle.com/technetwork/database/nosqldb/learnmore/nosql-database-498041.pdf>, 2011. (Accessed February 2017).
- [380] OSA, SP-011: cloud computing pattern, <http://www.opensecurityarchitecture.org/cms/library/pattern-landscape/251-pattern-cloud-computing>. (Accessed February 2017).
- [381] D.L. Osisek, K.M. Jackson, P.H. Gum, ESA/390 interpretive-execution architecture, foundation for VM/ESA, IBM Systems Journal 30 (1) (1991) 34–51.
- [382] K. Ousterhout, P. Wendell, M. Zaharia, I. Stoica, Sparrow: distributed, low latency scheduling, in: Proc. 24th ACM Symp. on Operating Systems Principles, 2013, pp. 69–84.
- [383] N. Oza, K. Karppinen, R. Savola, User experience and security in the cloud - an empirical study in the Finnish Cloud Consortium, in: Proc. IEEE 2nd Int. Conf. Cloud Computing Technology and Science, 2010, pp. 621–628.
- [384] G. Pacifici, M. Spreitzer, A.N. Tantawi, A. Youssef, Performance management for cluster-based web services, IEEE Journal on Selected Areas in Communications 23 (12) (2005) 2333–2343.

- [385] P. Padala, X. Zhu, Z. Wang, S. Singhal, K.G. Shin, Performance evaluation of virtualization technologies for server consolidation, HP Technical Report HPL-2007-59, 2007, also <http://www.hpl.hp.com/techreports/2007/HPL-2007-59R1.pdf>. (Accessed August 2015).
- [386] P. Paillier, Public-key cryptosystems based on composite degree residuosity classes, in: Advances in Cryptology, Springer-Verlag, Heidelberg, 1999, pp. 223–238.
- [387] V. Pankratius, F. Lind, A. Coster, P. Erickson, J. Semeter, Space weather monitoring using multicore mobile devices, American Geographic Union (AGU) Fall Meeting Abstracts, <http://adsabs.harvard.edu/abs/2013AGUFMSA31B..06P>, also <https://mahali.mit.edu/sites/default/files/documents/Mahali-Overview.pdf>. (Accessed November 2016).
- [388] M. Paolino, A. Rigo, A. Spyridakis, J. Fanguède, P. Lalov, D. Raho, T-KVM: a trusted architecture for KVM ARM v7 and v8 virtual machines securing virtual machines by means of KVM, TrustZone, TEE, and SELinux, in: Proc. 6th Int. Conf. on Cloud Computing, GRIDs, and Virtualization, IARIA, Nice, France, 2015, pp. 39–45.
- [389] D.L. Parnas, On the criteria to be used in decomposing systems into modules, Communications of the ACM 15 (12) (1972) 1053–1058.
- [390] S. Parvin, S. Han, L. Gao, F. Hussain, E. Chang, Towards trust establishment for spectrum selection in cognitive radio networks, in: Proc. IEEE 24th Int. Conf. Advanced Information Networking and Applications, 2010, pp. 579–583.
- [391] A.M.K. Pathan, R. Buya, A taxonomy of content delivery networks, <http://cloudbus.org/reports/CDN-Taxonomy.pdf>, 2009. (Accessed August 2015).
- [392] B. Pawłowski, C. Juszezak, P. Staubach, C. Smith, D. Label, D. Hitz, NFS Version 3 design and implementation, in: Proc. Summer 1994 Usenix Technical Conference, 1994, pp. 137–151.
- [393] B. Pawłowski, S. Shepler, C. Beame, B. Callaghan, M. Eisler, D. Noveck, D. Robinson, R. Turlow, The NFS Version 4 protocol, in: Proc. 2nd Int. SANE Conf. System Administration and Network Engineering, 2000.
- [394] A. Paya, D.C. Marinescu, A cloud service for adaptive digital music streaming, in: Proc. 8th Int. Conf. Signal Image Technology and Internet Systems, 2012.
- [395] A. Paya, D.C. Marinescu, Energy-aware load balancing and application scaling for the cloud ecosystem, IEEE Transactions on Cloud Computing 5 (1) (2017) 15–27.
- [396] S. Pearson, A. Benameri, Privacy, security, and trust issues arising from cloud computing, in: Proc. Cloud Computing and Science, 2010, pp. 693–702.
- [397] M. Perrin, Distributed Systems: Concurrency and Consistency, Elsevier, Oxford, UK, 2017.
- [398] C.A. Petri, Kommunikation mit Automaten, Schriften des Rheinisch-Westfälisches Institutes für Instrumentelle Mathematik, vol. 2, 1962, Bonn.
- [399] C.A. Petri, Concurrency Theory, Lecture Notes in Computer Science, vol. 254, Springer-Verlag, Heidelberg, 1987, pp. 4–24.
- [400] A. Pnueli, The temporal logic of programs, in: Proc. 18th Annual IEEE Symp. Foundations of Computer Science, 1977, pp. 46–57.
- [401] R.A. Popa, C.M.S. Redfield, N. Zeldovich, H. Balakrishnan, Cryptdb: protecting confidentiality with encrypted query processing, in: Proc. 23 ACM Symp. on Operating Systems Principles, 2011, pp. 85–100.
- [402] G.J. Popek, R.P. Golberg, Formal requirements for virtualizable third generation architecture, Communications of the ACM 17 (7) (1974) 412–421.
- [403] D. Price, A. Tucker, Solaris Zones: operating systems support for consolidating commercial workloads, in: Proc. 18th Large Installation System Administration, USENIX, 2004, pp. 241–254.
- [404] M. Price, The paradox of security in virtual environments, Computer 41 (11) (2008) 22–28.
- [405] X. Pu, L. Liu, Y. Mei, S. Sivathanu, Y. Koh, C. Pu, Understanding performance interference of I/O workload in virtualized cloud environments, in: Proc. 3rd IEEE Int. Conf. Cloud Computing, 2010, pp. 51–58.
- [406] P. Radzikowski, SAN vs DAS: a cost analysis of storage in the enterprise (updated 2010), <http://capitalhead.com/articles/san-vs-das-a-cost-analysis-of-storage-in-the-enterprise.aspx>. (Accessed February 2017).

- [407] F.Y. Rashid, The dirty dozen: 12 cloud security threats, Infoworld, <https://www.computerworld.com/article/3043506/the-dirty-dozen-12-cloud-security-threats.html>, October, 2021.
- [408] V.J. Reddi, H. Yoon, A. Knies, Two billion devices and counting, IEEE MICRO 38 (38) (2018) 6–21.
- [409] T.K. Refaat, B. Kantarci, H.T. Mouftah, Virtual machine migration and management for vehicular clouds, Vehicular Communications 4 (2016) 47–56.
- [410] N. Regola, J.-C. Ducom, Recommendations for virtualization technologies in high performance computing, in: Proc. 2nd IEEE Int. Conf. Cloud Computing Technology and Science, 2010, pp. 409–416.
- [411] C. Reiss, J. Wilkes, J.L. Hellerstein, Google cluster-usage traces: format+schema, Technical report, Google, Mountain View, CA, 2011, <https://github.com/google/clusterdata>.
- [412] C. Reiss, A. Tumanov, G. Ganger, R. Katz, M. Kozuch, Heterogeneity and dynamicity of clouds at scale: Google trace analysis, in: Proc. ACM Symp. Cloud Computing, 2012, Article #7.
- [413] G. Ren, E. Tune, T. Moseley, Y. Shi, S. Rus, R. Hundt, Google-wide profiling: a continuous profiling infrastructure for data centers, IEEE MICRO 30 (4) (2010) 65–79, http://static.googleusercontent.com/external_content/untrusted_dlcp/research.google.com/en/us/pubs/archive/36575.pdf. (Accessed August 2016).
- [414] J. Reyes, D.C. Marinescu, E. Mucciolo, Simulation of quantum many-body systems on Amazon cloud, Computer Physics Communications 261 (April 2021) 107750.
- [415] M. Riandato, Jails Free BSD handbook, <http://www.freebsd.org/doc/handbook/jails.html>. (Accessed January 2017).
- [416] D.M. Ritchie, K. Thompson, The Unix time-sharing system, Communications of the ACM 17 (7) (1974) 365–375.
- [417] D.M. Ritchie, The evolution of the Unix time-sharing system, Bell Labs Technical Journal 63 (2.2) (1984) 1577–1593.
- [418] R. Rivest, A. Shamir, L. Adleman, A method for obtaining digital signatures and public-key cryptosystems, Communications of the ACM 120 (1978) 126.
- [419] A.A. Rocha, T. Salonidis, T. He, D. Towsley, sLRFU: a data streaming based least recently frequently used caching policy, <https://pdfs.semanticscholar.org/ea8d/df7b03109e18633c0b94a09a3dce03b18059.pdf>, 2015. (Accessed March 2017).
- [420] R. Rodrigues, P. Druschel, Peer-to-peer systems, Communications of the ACM 53 (10) (2010) 72–82.
- [421] J. Rolia, L. Cerkasova, M. Arlit, A. Andrzejak, A capacity management service for resource pools, in: Proc. 5th Int. Workshop on Software and Performance, ACM, 2005, pp. 229–237.
- [422] F. Rosenblatt, The perceptron: a probabilistic model for information storage and organization in the brain, Psychological Review 65 (6) (1958) 386–408.
- [423] M. Rosenblum, T. Garfinkel, Virtual machine monitors: current technology and future trends, Computer 38 (5) (2005) 39–47.
- [424] D.M. Rousseau, S.B. Sitkin, R.S. Burt, C. Camerer, Not so different after all: a cross-disciplinary view of trust, The Academy of Management Review 23 (3) (1998) 393–404.
- [425] D.E. Rumelhart, G.E. Hinton, R.J. Williams, Learning representations by back-propagating errors, Nature 323 (1986) 533–536.
- [426] K. Rzadka, et al., Autopilot: workload autoscaling at Google, in: Proc. EuroSys20, 2020, pp. 1–16.
- [427] P. Sadalage, M. Fowler, NoSQL Distilled; a Brief Guide to the Emerging World of Polyglot Persistence, Addison-Wesley, 2012.
- [428] H.F.-W. Sadrozinski, J. Wu, Applications of Field-Programmable Gate Arrays in Scientific Research, Taylor and Francis Inc., Bristol, PA, USA, 2010.
- [429] A. Salomaa, Public-Key Cryptography, Springer-Verlag, Heidelberg, 1990.
- [430] J.H. Saltzer, M.F. Kaashoek, Principles of Computer System Design, Morgan Kaufmann, Waltham, MA, 2009.
- [431] N. Samaan, A novel economic sharing model in a federation of selfish cloud providers, IEEE Transactions on Parallel and Distributed Systems 25 (1) (2014) 12–21.

- [432] R.R. Sambasivan, A.X. Zheng, M. De Rosa, E. Krevat, S. Whitman, M. Stroucken, W. Wang, L. Xy, G.R. Ganger, Diagnosing performance changes by comparing request flows, in: Proc. 8th USENIX Conf. Networked Systems Design and Implementation, 2011, 14 pp.
- [433] B. Sandberg, D. Goldberg, S. Kleiman, D. Walsh, B. Lyon, Design and implementation of the Sun network file system, in: Proc. Summer Usenix, Technical Conference, 1986, pp. 119–130.
- [434] T. Sandholm, K. Lai, Dynamic proportional share scheduling in Hadoop, in: Proc. 15th Workshop Job Scheduling Strategies for Parallel Processing, in: Lecture Notes in Computer Science, vol. 6253, Springer-Verlag, Heidelberg, 2010, pp. 110–131.
- [435] R. Sadhu, Good-enough security; toward a pragmatic business-driven discipline, *IEEE Internet Computing* 7 (1) (2003) 66–68.
- [436] M. Satyanarayanan, A survey of distributed file systems, CS Technical Report, CMU, <http://www.cs.cmu.edu/~satya/docdir/satya89survey.pdf>, 1989. (Accessed August 2015).
- [437] M. Satyanarayanan, P. Bahl, R. Caceres, N. Davies, The case for VM-based cloudlets in mobile computing, *IEEE Transactions on Pervasive Computing* 8 (4) (2009) 14–23.
- [438] M. Satyanarayanan, Mobile computing: the next decade, in: Proc. 1st ACM Workshop on Mobile Cloud Computing & Services: Social Networks and Beyond, 2010, pp. 1–5.
- [439] J. Savage, Models of Computation: Exploring the Power of Computing, Addison-Wesley, Reading, Ma, 1998.
- [440] F. Schmuck, R. Haskin, GFPS: a shared-disk file system for large computing clusters, in: Proc. Conf. File and Storage Technologies, USENIX, 2002, pp. 231–244.
- [441] P. Schuster, Nonlinear dynamics from Physics to Biology. Self-organization: an old paradigm revisited, *Complexity* 12 (4) (2007) 9–11.
- [442] M. Schwarzkopf, A. Konwinski, M. Abd-El-Malek, J. Wilkes, Omega: flexible, scalable schedulers for large compute clusters, in: Proc. 8th ACM European Conf. on Computer Systems, 2013, pp. 351–364.
- [443] S. Scott, D. Abts, J. Kim, W.J. Dally, The Blackwidow highradix Clos network, in: Proc. 33rd Annual Int. Symp. on Computer Architecture, IEEE, 2006, pp. 16–28.
- [444] D. Sehr, R. Muth, C. Biffle, V. Khimenko, E. Pasko, K. Schimpf, B. Yee, B. Chen, Adapting software fault isolation to contemporary CPU architectures, in: Proc. 19th USENIX Security Symposium, 2010, pp. 1–11.
- [445] P. Sempolinski, D. Thain, A comparison and critique of Eucalyptus, OpenNebula and Nimbus, in: Proc. 2nd IEEE Int. Conf. Cloud Computing Technology and Science, 2010, pp. 417–426.
- [446] J. Shneidman, C. Ng, D.C. Parkes, A. AuYoung, A.C. Snoeren, A. Vahdat, B. Chun, Why markets could (but don't currently) solve resource allocation problems in systems, in: Proc. 10th Worshop on Hot Topics in Operating Systems, 2005, pp. 7–14.
- [447] J. Shute, R. Vingralek, B. Samwel, B. Handy, C. Whipkey, E. Rollins, M. Oancea, K. Littlefield, D. Menestrina, S. Ellner, J. Cieslewicz, I. Rae, T. Stancescu, H. Apte, F1: a distributed SQL database that scales, *Proceedings of the VLDB Endowment* 6 (11) (2013) 1068–1079.
- [448] H.A. Simon, Administrative Behavior, Macmillan, New York, NY, 1947.
- [449] S. Sivathanu, L. Liu, M. Yiduo, X. Pu, Storage management in virtualized cloud environment, in: Proc. 3rd IEEE Int. Conf. Cloud Computing, 2010, pp. 204–211.
- [450] J.E. Smith, R. Nair, The architecture of virtual machines, *Computer* 38 (5) (2005) 32–38.
- [451] L. Snyder, Type architectures, shared memory, and the corollary of modest potential, *Annual Review of Computer Science* 1 (1986) 289–317.
- [452] B. Snyder, Server virtualization has stalled, despite the hype, <http://www.infoworld.com/article/2624771/server-virtualization-has-stalled-despite-the-hype.html>, 2010. (Accessed February 2017).
- [453] I. Sommerville, D. Cliff, R. Calinescu, J. Keen, T. Kelly, M. Kwiatowska, J. McDermid, R. Paige, Large-scale IT complex systems, *Communications of the ACM* 55 (7) (2012) 71–77.
- [454] M. Stecca, M. Maresca, An architecture for a mashup container in virtualized environments, in: Proc. 3rd IEEE nt. Conf. Cloud Computing, 2010, pp. 386–393.

- [455] R. Stoica, A. Ailamaki, Enabling efficient OS paging for main-memory OLTP databases, in: Proc. 9th Int. Workshop on Data Management on New Hardware, 2013, Article #6.
- [456] M. Stokely, J. Winget, E. Keyes, C. Grimes, B. Yolken, Using a market economy to provision compute resources across planet-wide clusters, in: Proc. IEEE Int. Symp. on Parallel and Distributed Processing, 2009, pp. 1–8.
- [457] I. Stoica, R. Morris, D. Karger, M.F. Kaashoek, H. Balakrishnan, Chord: a scalable peer-to-peer lookup service for Internet applications, in: Proc. ACM/SIGCOMM Conf. on Applications, Technologies, Architectures, and Protocols for Computer Communications, 2001, pp. 149–160.
- [458] M. Stonebreaker, The “NoSQL” has nothing to do with SQL, <http://cacm.acm.org/blogs/blog-cacm/50678-the-nosql-discussion-has-nothing-to-do-with-sql/fulltext>, 2009. (Accessed February 2017).
- [459] Stratokey, Cloud data protection, <https://www.stratokey.com/solutions/cloud-access-security-brokers?gclid=CM7bp5eOh80CFQ8kgQodYu4L5g>, 2015. (Accessed August 2016).
- [460] StreamingMedia, Only 18% using adaptive streaming, says Skyfire report, <http://www.streamingmedia.com/Articles/ReadArticle.aspx?ArticleID=79393>. (Accessed August 2015).
- [461] J. Stribling, J. Li, I.G. Councill, M.F. Kaashoek, R. Morris, Overcite: a distributed, cooperative citeseer, in: Proc. 3rd Symp. on Networked Systems Design and Implementation, 2006, pp. 69–79.
- [462] J. Sugerman, G. Venkitachalam, B. Lim, Virtualizing I/O devices on VMware Workstation’s hosted virtual machine monitor, in: Proc. USENIX Annual Technical Conference, 2001, pp. 1–14.
- [463] C. Sun, W. Zhang, K.B. Letaief, Cluster-based cooperative spectrum sensing for cognitive radio systems, in: Proc. IEEE Conf. on Communications, 2007, pp. 2511–2515.
- [464] C. Sun, W. Zhang, K.B. Letaief, Cooperative spectrum sensing for cognitive radios under BW constraints, in: Proc. IEEE Wireless Communications and Networking Conference, 2007, pp. 1–5.
- [465] Y. Sun, Z. Han, K.J. Ray Liu, Defense of trust management vulnerabilities in distributed networks, in: Special Issue, Security in Mobile Ad Hoc and Sensor Networks, IEEE Communications Magazine 46 (2) (2008) 112–119.
- [466] V. Sundaram, A. Chandra, P. Goyal, P. Shenoy, J. Sahni, H.M. Vin, Application performance in the QLinux multimedia operating system, in: Proc. 8th ACM Int. Conf. on Multimedia, 2000, pp. 127–136.
- [467] M. Steinder, I. Walley, D. Chess, Server virtualization in autonomic management of heterogeneous workloads, ACM SIGOPS Operating Systems Review 42 (1) (2008) 94–95.
- [468] T. Taleb, A. Ksentini, Follow me cloud: interworking distributed clouds & distributed mobile networks, IEEE Network 27 (5) (2013) 12–19.
- [469] D. Tancock, S. Pearson, A. Charlesworth, A privacy impact assessment tool for cloud computing, in: Proc. 2nd IEEE Int. Conf. Cloud Computing, 2010, pp. 667–674.
- [470] L. Tang, J. Dong, Y. Zhao, L.-J. Zhang, Enterprise cloud service architecture, in: Proc. 3rd IEEE Int. Conf. Cloud Computing, 2010, pp. 27–34.
- [471] J. Tang, Y. Cui, K. Ren, J. Liu, R. Buyya, Ensuring security and privacy preservation for cloud data services, ACM Computing Surveys 49 (1) (2016) 13.
- [472] J. Tate, F. Lucchese, R. Moore, Introduction to Storage Area Networks, IBM Redbooks, 2006, <http://www.redbooks.ibm.com/redbooks/pdfs/sg245470.pdf>. (Accessed August 2015).
- [473] D. Tennenhouse, Layered multiplexing considered harmful, in: H. Rudin, R.C. Williamson (Eds.), Protocols for High-Speed Networks, North Holland, 1989, pp. 143–148.
- [474] G. Tesauro, N.K. Jong, R. Das, M.N. Bennani, A hybrid reinforcement learning approach to autonomic resource allocation, in: Proc. IEEE Int Conf. on Autonomic Computing, 2006, pp. 65–73.
- [475] A. Thusoo, J.S. Sarma, N. Jain, Z. Shao, P. Chakka, S. Anthony, H. Liu, P. Wyckoff, R. Murthy, Hive a warehousing solution over a MapReduce framework, Proceedings of the VLDB Endowment 2 (2) (2009) 1626–1629.
- [476] M. Tirmazi, et al., Borg: the next generation, in: Proc. EuroSys, ACM, 2020, <https://doi.org/10.1145/3342195.3387517>, ISBN 978-1-4503-6882-7/20/04.

- [477] J. Timmermans, V. Ikonen, B.C. Stahl, E. Bozdag, The ethics of cloud computing. A conceptual review, in: Proc. 2nd IEEE Int. Conf. Cloud Computing Technology and Science, 2010, pp. 614–620.
- [478] Top 500 supercomputers, <http://top500.org/featured/top-systems/>. (Accessed January 2016).
- [479] TRIAD, Translating relaying Internet architecture integrating active directories, https://www.researchgate.net/publication/235184440_TRIAD_Translating_Relying_Internetwork_Architecture_Integrating_Active_Directories_Final_Report. (Accessed October 2021).
- [480] C. Tung, M. Steinder, M. Spreitzer, G. Pacifici, A scalable application placement controller for enterprise data centers, in: Proc. 16th Int. Conf. on WWW, 2007.
- [481] A.M. Turing, On computable numbers, with an application to the Entscheidungsproblem, Proceedings of the London Mathematical Society. Series 2 42 (1937) 230–265, and On computable numbers, with an application to the Entscheidungsproblem: a correction, Proceedings of the London Mathematical Society. Series 2 43 (1937) 544–546.
- [482] A.M. Turing, The chemical basis of morphogenesis, Philosophical Transactions of the Royal Society of London. Series B 237 (1952) 37–72.
- [483] R. Urgaonkar, S. Wang, T. He, M. Zafer, K. Chan, K.K. Leun, Dynamic service migration and workload scheduling in edge-clouds, Performance Evaluation 91 (C) (2015) 205–228.
- [484] USDOT, 2015–2019 strategic plan intelligent transportation systems (ITS), Joint Program Office (JPO) – FHWA-JPO-14-145, <http://www.its.dot.gov/strategicplan.pdf>, 2015.
- [485] L.G. Valiant, A bridging model for parallel computation, Communications of the ACM 33 (8) (1990) 103–112.
- [486] L.G. Valiant, A bridging model for multicore computing, in: Proc. 16th Annual European Symp. on Algorithms, in: Lecture Notes on Computer Science, vol. 5193, 2008, pp. 13–28.
- [487] J. Varia, Cloud architectures, <https://aws.amazon.com/articles/building-grepttheweb-in-the-cloud-part-1-cloud-architectures/>.
- [488] J. van Vliet, F. Paganelli, S. van Wel, D. Dowd, Elastic Beanstalk: Simple Cloud Scaling for Java Developers, O'Reilly Publishers, Sebastopol, California, 2011.
- [489] L.M. Vaquero, L. Rodero-Merino, R. Buyya, Dynamically scaling applications in the cloud, ACM SIGCOMM Computer Communication Review 41 (1) (2011) 45–52.
- [490] H.N. Van, F.D. Tran, J.-M. Menaud, Performance and power management for cloud infrastructures, in: Proc. IEEE Conf. Cloud Computing, 2010, pp. 329–336.
- [491] K. Varadhan, R. Govindan, D. Estrin, Persistent route oscillations in interdomain routing, Computer Networks 32 (1) (2000) 1–16.
- [492] S. Venkataraman, A. Panda, G. Ananthanarayanan, M.J. Franklin, I. Stoica, The power of choice in data-aware cluster scheduling, in: Proc. 11th USENIX Symp. Operating Systems Design and Implementation, 2014, pp. 301–314.
- [493] P. Veríssimo, L. Rodrigues, Aposteriori agreement for fault-tolerant clock synchronization on broadcast networks, in: Proc. 22nd Annual Int. Symp. on Fault-Tolerant Computing, IEEE Press, Los Alamitos, CA, 1992, pp. 527–536.
- [494] A. Verma, G. Dasgupta, T.K. Nayak, P. De, R. Kothari, Server workload analysis for power minimization using consolidation, in: Proc. USENIX Annual Technical Conference, 2009, p. 28.
- [495] A. Verma, L. Pedrosaz, M. Korupolu, D. Oppenheimer, E. Tune, J. Wilkes, Large-scale cluster management at Google with Borg, in: Proc. 10th European Conference on Computer Systems, 2015, Article No. 18.
- [496] VMware, VMware vSphere storage appliance, <https://www.vmware.com/files/pdf/techpaper/VM-vSphere-Storage-Appliance-Deep-Dive-WP.pdf>. (Accessed August 2015).
- [497] J. von Neumann, Probabilistic logic and synthesis of reliable organisms from unreliable components, in: C.E. Shannon, J. McCarthy (Eds.), Automata Studies, Princeton University Press, Princeton, NJ, 1956.
- [498] J. von Neumann, Fourth University of Illinois Lecture, in: A.W. Burks (Ed.), Theory of Self-Reproducing Automata, University of Illinois Press, Champaign, IL, 1966.

- [499] S.V. Vrbsky, M. Lei, K. Smith, J. Byrd, Data replication and power consumption in data grids, in: Proc. 2nd IEEE Int. Conf. Cloud Computing Technology and Science, 2010, pp. 288–295.
- [500] B. Walker, G. Popek, E. English, C. Kline, G. Thiel, The LOCUS distributed operating system, in: Proc. 9th ACM Symp. Operating Systems Principles, 1983, pp. 49–70.
- [501] M.E. Wall, A. Rechtsteiner, L.M. Rocha, Singular value decomposition and principal component analysis, in: D.P. Berrar, W. Dubitzky, M. Granzow (Eds.), *A Practical Approach to Microarray Data Analysis*, Kluwer, Norwell, MA, 2003.
- [502] K. Walsh, E.G. Sirer, Experience with an object reputation system for peer-to-peer file sharing, in: Proc. 3rd Symp. Networked Systems Design and Implementation, 2006, pp. 1–14.
- [503] L. Wang, L. Park, R. Pang, V.S. Pai, L. Peterson, Reliability and security in the CoDeeN content distribution network, in: Proc. USENIX Annual Technical Conference, 2004, 14 pp.
- [504] M. Wang, N. Kandasamy, A. Guez, M. Kam, Adaptive performance control of computing systems via distributed cooperative control: application to power management in computer clusters, in: Proc. 3rd IEEE Intl. Conf. on Autonomic Computing, 2006, pp. 165–174.
- [505] Y. Wang, I-R. Chen, D-C. Wang, A survey of mobile cloud computing applications: perspectives and challenges, *Wireless Personal Communications* 80 (4) (2015) 1607–1623.
- [506] S. Wang, R. Urgaonkar, M. Zafer, T. He, K. Chan, K.K. Leun, Dynamic service migration in edge-clouds, in: Proc. IFIP Networking Conf., 2015.
- [507] D.J. Watts, S.H. Strogatz, Collective-dynamics of small-world networks, *Nature* 393 (1998) 440–442.
- [508] J. Webster, Evaluating IBM's SVC and TPC for server virtualization, ftp://ftp.boulder.ibm.com/software/at/tivoli/analyst_paper_ibm_svc_tpc.pdf. (Accessed August 2016).
- [509] G. Weikum, G. Vossen, *Transactional Information Systems: Theory, Algorithms, and the Practice of Concurrency Control and Recovery*, Morgan Kaufmann, Waltham, MA, 2001.
- [510] J. Wilkes, Google cluster-usage traces v3, Technical report, Google, Mountain View, CA, 2019, <https://github.com/google/cluster-data>.
- [511] M. Whaiduzzaman, M. Sookhak, A. Gani, R. Buyya, A survey on vehicular cloud computing, *Journal of Network and Computer Applications* 40 (2014) 325–344.
- [512] S.E. Whang, H. Garcia-Molina, Managing information leakage, in: Proc. 5th Biennal Conf. on Innovative Data Systems Research, 2011, also <http://ilpubs.stanford.edu:8090/987/>. (Accessed August 2015).
- [513] A. Williams, *C++ Concurrency in Action: Practical Multithreading*, Manning Publications, Shelter Island, NY, 2012.
- [514] A. Whitaker, M. Shaw, S.D. Gribble, Denali; lightweight virtual machines for distributed and networked applications, Technical Report 02-0201, University of Washington, 2002.
- [515] V. Winkler, *Securing the Cloud: Cloud Computer Security Techniques and Tactics*, Elsevier Science and Technologies Books, 2011.
- [516] J.A. Winter, D.H. Albonesi, C.A. Shoemaker, Scalable thread scheduling and global power management for heterogeneous many-core architectures, in: Proc. 19th Int. Conf. Parallel Architectures and Compilation Techniques, 2010, pp. 29–40.
- [517] E.C. Withana, B. Plale, Usage patterns to provision for scientific experimentation in clouds, in: Proc. 2nd IEEE Int. Conf. Cloud Computing Technology and Science, 2010, pp. 226–233.
- [518] S.A. Wolf, A.Y. Chtchelkanova, D.M. Treger, Spintronics - a retrospective and perspective, *IBM Journal of Research and Development* 50 (1) (2006) 101–110.
- [519] Xen Wiki, https://wiki.xenproject.org/wiki/Main_Page, 2007.
- [520] Z. Xiao, D. Cao, A policy-based framework for automated SLA negotiation for Internet-based virtual computing environment, in: Proc. 16th IEEE Int. Conf. Parallel and Distributed Systems, 2010, pp. 694–699.
- [521] R.S. Xin, J. Rosen, M. Zaharia, M.J. Franklin, S. Shenker, I. Stoica, Shark: SQL and rich analytics at scale, in: ACM SIGMOD Int. Conf. on Management of Data, 2013, pp. 13–24.

- [522] G. Xylomenos, C.N. Ververidis, V.A. Siris, N. Fotiou, C. Tsilopoulos, X. Vasilakos, K.V. Katsaros, G.C. Polyzos, A survey of information-centric networking research, *IEEE Communications Surveys and Tutorials* 16 (2) (2014) 1024–1049.
- [523] G. Yan, D. Wen, S. Olariu, M.C. Weigle, Security challenges in vehicular cloud computing, *IEEE Transactions on Intelligent Transportation Systems* 4 (1) (2013) 6–16.
- [524] G. Yan, S. Olariu, J. Wang, S. Arif, Towards providing scalable and robust privacy in vehicular networks, *IEEE Transactions on Parallel and Distributed Systems* 25 (7) (2014) 1896–1906.
- [525] A.C. Yao, Some complexity questions related to distributed computing, in: Proc. 11th Symp. on the Theory of Computing, 1979, pp. 209–213.
- [526] A.C. Yao, How to generate and exchange secrets, in: Proc. 27th Annual Symp. on Foundations of Computer Science, 1986, pp. 162–167.
- [527] A. Yasin, A top-down method for performance analysis and counters architecture, in: Proc. IEEE Int. Symp. on Perf. Analysis Systems and Software, 2014, pp. 1–10.
- [528] A. Yasin, Y. Ben-Asher, A. Mendelson, Deep-dive analysis of the data analytics workload in CloudSuite, in: Proc. IEEE Int. Workshop/Symp. on Workload Characterization, 2014, pp. 1–10, Paper 67.
- [529] H. Yu, X. Bai, D.C. Marinescu, Workflow management and resource discovery for an intelligent grid, *Parallel Computing* 31 (7) (2005) 797–811.
- [530] Y. Yu, M. Isard, D. Fetterly, M. Budiu, U. Erlingsson, P. Kumar, G.J. Currey, DryadLINQ: a system for general-purpose distributed data-parallel computing using a high-level language, in: Proc. 8th USENIX Symp. Operating System Design and Implementation, 2009, pp. 1–14.
- [531] M. Zapf, A. Heinzl, Evaluation of process design patterns: an experimental study, in: W.M.P. van der Aalst, J. Desel, A. Oberweis (Eds.), Business Process Management, in: Lecture Notes on Computer Science, vol. 1806, Springer-Verlag, Heidelberg, 2000, pp. 83–98.
- [532] M. Zaharia, D. Borthakur, J.S. Sarma, K. Elmeleegy, S. Shenker, I. Stoica, Delay scheduling: a simple technique for achieving locality and fairness in cluster scheduling, in: Proc. 5th European Conf. on Computer Systems, 2010, pp. 265–278.
- [533] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M.J. Franklin, S. Shenker, I. Stoica, Resilient distributed datasets: a fault-tolerant abstraction for in-memory cluster computing, in: Proc. 9th USENIX Conf. on Networked Systems Design and Implementation, 2012, pp. 2–16.
- [534] M. Zaharia, T. Das, H. Li, S. Shenker, I. Stoica, Discretized streams: an efficient and fault-tolerant model for stream processing on large clusters, in: Proc. 4th USENIX Conference on Hot Topics in Cloud Computing, USENIX Association, 2012, pp. 10–16.
- [535] P. Zech, M. Felderer, R. Breu, Towards a model-based security testing approach in cloud computing environments, in: Proc. 6th IEEE Int. Conf. on Software Security and Reliability Companion, 2012, pp. 47–56.
- [536] Z.L. Zhang, D. Towsley, J. Kurose, Statistical analysis of the generalized processor sharing scheduling discipline, *IEEE Journal on Selected Areas in Communications* 13 (6) (1995) 1071–1080.
- [537] X. Zhang, J. Liu, B. Li, T.-S.P. Yum, CoolStreaming/DONet: a data-driven overlay network for peer-to-peer live media streaming, in: Proc. IEEE INFOCOM 2005, 2005, pp. 2102–2111.
- [538] C. Zhang, H.D. Sterck, CloudBATCH: A batch job queuing system on clouds with Hadoop and HBase, in: Proc. 2nd IEEE Int. Conf. Cloud Computing Technology and Science, 2010, pp. 368–375.
- [539] L. Zhang, A. Afanasyev, J. Burke, V. Jacobson, K.C. Claffy, P. Crowley, C. Papadopoulos, L. Wang, B. Zhang, Named data networks, *ACM SIGCOMM Computer Communication Review* 44 (3) (2014) 66–73.
- [540] T. Zhang, R.E. DeGrande, A. Boukerche, Urban traffic characterization for enabling vehicular clouds, in: Proc. IEEE Wireless Communications and Networking Conf., April 2016.
- [541] X. Zhao, K. Borders, A. Prakash, Virtual machine security systems, in: Advances in Computer Science and Engineering, 2011, pp. 339–365.

Glossary

A

Access Control List (ACL) list of pairs (subject, value) defining the set of access rights to an object; for example, read, write, and execute permissions for a file.

Advanced Configuration and Power Interface (ACPI) open standard for device configuration and power management by the operating system. It defines four Global “Gx” states and six Sleep “Sx” states. For example, “S3” is referred to as Standby, Sleep, or Suspend to RAM.

Advanced Microcontroller Bus Architecture (AMBA) open standard, on-chip interconnect specification for the connection and management of a large number of controllers and peripherals.

Amazon Machine Image (AMI) a unit of deployment, an environment including all information necessary to set up and boot an instance including: (1) a template for the root volume for the instance, e.g., an operating system, an application server, and applications; (2) launch permissions controlling all AWS accounts that can use the AMI to launch instances; and (3) a block device mapping specifying the volumes to be attached to the instance when launched.

Amdahl's law formula used to predict the theoretical maximum speedup for a program using multiple processors/cores. Informally, it states that the portion of the computation that cannot be parallelized determines the overall speedup.

Anti-entropy a process, often using Merkle trees, for comparing the data of all replicas and updating each replica to the newest version.

App Engine (AE) an ensemble of computer, storage, search, and networking services for building web and mobile applications and running them on Google servers.

Application Binary Interface (ABI) the projection of the computer system seen by a process or thread in execution. ABI allows the ensemble consisting of the application and the library modules to access the hardware. ABI does not include privileged system instructions; instead, it invokes system calls.

Application Program Interface (API) defines the set of instructions the hardware was designed to execute and gives the application access to the Instruction Set Architecture layer. It includes High Level Language (HLL) library calls that often invoke system calls. The API is the projection of the system from the perspective of the HLL program.

Application layer deployed software applications, targeted towards end-user software clients or other programs, and made available via the cloud.

Auction a sale where items are sold to the highest bidder.

Auditor party conducting independent assessments of cloud services, information system operations, the performance, and the security of the cloud implementation.

Audit systematic evaluation of a cloud system by measuring how well it conforms to a set of established criteria, e.g., security audit if the criteria is security, privacy-impact audit if the criteria is privacy assurance, and performance audit if the criteria is performance.

Authentication credential something that an entity is, has, or knows that allows that entity to prove its own identity to a system.

Auto Scaling AWS service providing automatic scaling of EC2 instances through grouping of instances, monitoring of the instances in a group, and defining *triggers*, pairs of CloudWatch alarms, and policies that allow the size of the group to be scaled up or down.

AWK utility utility for text processing based on a scripting language.

B

Bandwidth the number of operations per unit of time; for example, the bandwidth of a processor is expressed in Mips or Mflops, while the memory and I/O bandwidth are expressed in Mbps.

Basic Core Equivalent (BCE) quantity describing how resources of a multicore processor are allocated to the individual cores. For example, a symmetric core processor can be configured as sixteen 1-BCE cores, eight 2-BCE cores, four 4-BCE cores, two 6-BCE cores, or one 16-BCE cores. An asymmetric core processor may have ten 1-BCE cores and one 6-BCE core.

Basic input/output system (BIOS) system component invoked after a computer system is powered on to load the operating system and later to manage the data flow between the OS and devices, such as keyboard, mouse, disk, video adapter, and printer.

Bigquery fully managed enterprise data warehouse for large-scale data analytics on the Google cloud platform.

BigTable distributed storage system developed by Google to store massive amounts of data and to scale up to thousands of storage servers.

Bisection bandwidth the sum of the bandwidths of the minimal number of links that are cut when splitting the system into two parts.

Bit-level parallelism parallel computing based on increasing processor word size, thus lowering the number of instructions required to process larger size operands.

BitTorrent peer-to-peer communications protocol for file sharing.

Border Gateway Protocol (BGP) path-vector reachability protocol. It maintains a table of IP networks that designate network reachability among autonomous systems and makes the core routing decisions for the Internet based on path, network policies, and/or rule sets.

Borg management software for clusters consisting of tens of thousands of servers co-located and interconnected by a data center-scale network fabric.

Boundary value problem problem with conditions specified at the extremes of the independent variable(s).

Bounded input data defining property of batch processing. The computing engine has as input a dataset of known contents and size, as opposed to processing a continuous stream of incoming data.

Broker entity that manages the use, performance, and delivery of cloud services and negotiates relationships between cloud service providers and cloud users.

Buffer overflow anomaly in which a program, while writing data to a buffer, overruns the buffer's boundary and overwrites adjacent memory locations.

BusyBox software providing several stripped-down Unix tools in a single executable file and running in environments such as Linux, Android, FreeBSD, or Debian.

Bus.Device.Function (BDF) data used to describe PCI devices.

Byte-range tokens used to specify the range of read and write operations to data files.

Byzantine failure a fault presenting different symptoms to different observers. In a distributed system a Byzantine failure could be: an *omission failure*, e.g., a crash failure, failure to receive a request or to send a response; it could also be a *commission failure*, e.g., process a request incorrectly, corrupt the local state, and/or send an incorrect or inconsistent response to a request.

C

Callback executable code passed as an argument to other code; the callee is expected to execute the argument either immediately or at a later time for synchronous and, respectively, asynchronous callbacks.

Callstack data structure storing information about the active subprograms invoked during the execution of a program. Also called execution stack, program stack, control stack, or run-time stack.

Carrier a networking organization that provides connectivity and transports data between communicating entities. Also, a carrier signal is a transmitted electromagnetic pulse or wave at a steady base frequency on which information can be imposed by modulation.

Causal delivery extension of the First-In-First-Out (FIFO) delivery to the case in which a process receives messages from multiple sources.

Cell storage storage organization consisting of cells of the same size and objects fitting exactly in one cell.

Central Limit Theorem (CLT) statistical theory stating that the sum of a large number of independent random variables has a normal distribution.

Chaining in vector computers mechanisms allowing vector operations to start as soon as individual elements of vector source operands become available. Chaining operates on *convoy*, sets of vector instructions that can potentially be executed together.

Chameleon an NSF facility that is an OpenStack KVM experimental environment for large-scale cloud research.

Command Line Interface (CLI) provides the means for a user to interact with a program.

Client-server paradigm software organization enforcing modularity. It allows systems with different processor architecture, different operating systems, libraries, and other system software to cooperate.

Clock condition a strong clock condition in a distributed system requires an equivalence between the causal precedence and the ordering of the time stamps of messages.

Closed-box platforms systems with embedded cryptographic key that allow themselves to reveal their true identity to remote systems and authenticate the software running on them. Found on some cellular phones, game consoles, and ATMs.

Clos network multistage nonblocking network with an odd number of stages. In a Clos network, all packets overshoot their destinations and then hop back to it.

Cloud Bigtable high-performance NoSQL database service for large analytical and operational workloads on the Google cloud platform.

- Cloud Datastore** highly scalable NoSQL database for web and mobile applications on the Google cloud platform.
- CloudFormation** AWS service for creation of a stack describing the application infrastructure.
- Cloud Functions (CF)** a lightweight, event-based, asynchronous system to create single-purpose functions that respond to cloud events on the Google cloud platform.
- CloudLab** an NSF facility that serves as a testbed allowing researchers to experiment with cloud architectures and new applications.
- CloudWatch** AWS monitoring infrastructure used to collect and track metrics important for optimizing the performance of applications and for increasing the efficiency of resource utilization. Without installing any software, a user can monitor preselected metrics and then view graphs and statistics for these metrics.
- Coarse-grained parallelism** execution mode in which large blocks of code are executed before the concurrent threads/processes communicate with one another.
- Cognitive radio** wireless communication in which an intelligent transceiver detects which communication channels are not in use and uses them while avoiding channels in use.
- Cognitive radio trust** trust regarding the information received by a intelligent transceiver from other nodes.
- Combinatorial auction** auction in which participants can bid on combinations of items or packages.
- Community cloud** a cloud infrastructure shared by several organizations and supporting a specific community with shared concerns (e.g., mission, security requirements, policy, and compliance considerations).
- Communication channel** physical system allowing two entities to communicate with one another.
- Communication protocol** a communication discipline involving a finite set of messages exchanged among entities. A protocol implements error control, flow control, and congestion control mechanisms.
- Computation steering** interactively guiding a computational experiment towards a region of interest.
- Computer cloud** a collection of systems in a single administrative domain offering a set of computing and storage services; a form of utility computing.
- Computing grid** a distributed system consisting of a large number of loosely coupled, heterogeneous, and geographically dispersed systems in different administrative domains. The name is a metaphor for accessing computer power with similar ease as accessing electric power provided by the electric grid.
- Concurrency** activities are executed simultaneously.
- Concurrent write-sharing** multiple clients can modify the data in a file at the same time.
- Confidence interval** statistical measure offering a guarantee of the quality of a result. A procedure is said to generate confidence intervals with a specified coverage $\alpha \in [0, 1]$ if, on a proportion exactly α of the set of experiments, the procedure generates an interval that includes the answer. For example, a 95% confidence interval $[a, b]$ means that, in 95% of the experiments, the result will be in $[a, b]$.
- Conflict fraction** average number of conflicts per successful transactions in a transaction processing system.
- Congestion control** mechanism ensuring that the offered load of a network does not exceed the network capacity.
- Consistent hashing** hashing technique for reducing the number of keys to be remapped when a hash table is resized. On average, only K/n keys need to be remapped with K the number of keys and n the number of slots.
- Container Engine** cluster manager and orchestration system for Docker containers built on the Kubernetes system. It schedules and manages containers automatically according to user specifications on the Google cloud platform.
- Content** any type or volume of media, be it static or dynamic, monolithic or modular, live or stored, produced by aggregation, or mixed.
- Container** software system emulating a separate physical server; a container has its own files, users, process tree, IP address, shared memory, semaphores, and messages. Each container can have its own disk quotas.
- Control flow architecture** computer architecture when the program counter of a processor core determines the next instruction to be loaded in the instruction register and then executed.
- Control sensitive instructions** machine instructions changing either the memory allocation, or the execution to kernel mode.
- Cooperative spectrum sensing** mode of operation in which each node determines the occupancy of the spectrum based on its own measurements, combines it with information from its neighbors, and then shares its own spectrum occupancy assessment with its neighbors.
- Copy-on-write (COW)** mechanism used by virtual memory operating systems to minimize the overhead of copying the virtual memory of a process when a process creates a copy of itself.
- Cron** a job scheduler for Unix-like systems used to periodically schedule jobs; often used to automate system maintenance and administration.
- Cross-site scripting** the most popular form of attack against web sites; a browser permits the attacker to insert client-scripts into the web pages, and thus to bypass the access controls at the web site.

Compute Unified Device Architecture (CUDA) programming model invented by NVIDIA for using graphics processing units (GPUs) for general-purpose processing.

Cut subset of the local history of all processes of a process group. *The frontier of the cut* is an n -tuple consisting of the last event of every process included in the cut.

Cut-through (wormhole) network routing routing mechanism when a packet is forwarded to its next hop as soon as the header is received and decoded. The packet can experience blocking if the outgoing channel expected to carry it to the next node is in use; in this case, the packet has to wait until the channel becomes free.

D

Database as a Service (DBaaS) a cloud service where the database runs on the physical infrastructure of the cloud service provider.

Data Description Language (DDL) syntax similar to a computer programming language for defining data structures; it is widely used for database schemas.

Data Manipulation Language (DML) programming language used to retrieve, store, modify, delete, insert, and update data in database; SELECT, UPDATE, INSERT statements or query statements are examples of DML statements.

Dataflow architecture computer architecture in which operations are carried out at the time that their input becomes available.

Datagram basic transfer unit in a packet-switched network; it consists of a header containing control information necessary to transport its payload through the network.

Data hazards in pipelining potential danger situations in which the instructions in a pipeline are dependent upon one another.

Data-level parallelism an extreme form of coarse-grained parallelism, based on partitioning the data into chunks/blocks/segments and running concurrently either multiple programs or copies of the same program, each on a different data block.

Data portability the ability to transfer data from one system to another without being required to recreate or reenter data descriptions or to modify significantly the application being transported.

Data object a logical container of data that can be accessed over a network, e.g., a blob; it may be an archive, such as specified by the *tar* format.

Data-shipping allows fine-grained data sharing; an alternative to byte-range locking.

Deadlock synchronization anomaly occurring when concurrent processes or threads compete with one another for resources and reach a state from which none of them can proceed.

Denial-of-service attack (DoS attack) internet attack targeting a widely used network service that prevents legitimate access to the service. Forces the operating system of the targeted host(s) to fill the connection tables with illegitimate entries.

De-perimeterization process allowing systems to span the boundaries of multiple organizations and cross the security borders.

Direct Memory Access (DMA) hardware feature allowing I/O devices and other hardware subsystems direct access to the system memory without the CPU involvement. Also used for memory-to-memory copying and for offloading expensive memory operations, such as scatter-gather operations, from the CPU to the dedicated DMA engine. Intel includes I/O Acceleration Technology (I/OAT) on high-end servers.

Distributed system collection of computers interconnected via a network. Users perceive the system as a single, integrated computing facility.

Dynamic binary translation conversion of blocks of guest instructions from a portable code format to the instructions understood by a host system. Such blocks can be cached and reused to improve performance.

Dynamic instruction scheduling architectural feature of modern processors supporting out-of-order instruction execution. It can reduce the number of pipeline stalls but adds to circuit complexity.

Dynamic power range interval between the lower and the upper limit of the device power consumption. A large dynamic range means that the device is able to operate at a lower fraction of its peak power when its load is low.

Dynamic voltage scaling a power conservation technique; often used together with frequency scaling under the name *dynamic voltage and frequency scaling* (DVFS).

Dynamic voltage and frequency scaling (DVFS) power management technique to increase or decrease the operating voltage or the clock frequency of a processor to increase the instruction execution rate and, respectively, to reduce the amount of heat generated and to conserve power.

E

EC2 Placement Group a logical grouping of instances that allows the creation of a virtual cluster.

Elastic Beanstalk AWS service handling automatically the deployment, the capacity provisioning, the load balancing, the auto-scaling, and the application monitoring functions. It interacts with other AWS services, including EC2, S3, SNS, Elastic Load Balance, and AutoScaling.

Elastic Block Store (EBS) AWS service providing persistent block-level storage volumes for use with EC2 instances. EBS supports the creation of snapshots of the volumes attached to an instance and then uses them to restart an instance. The storage strategy provided by EBS is suitable for database applications, file systems, and applications using raw data devices.

Elastic Compute Cloud (EC2) AWS service for launching instances of an application under several operating systems, such as several Linux distributions, Windows, OpenSolaris, FreeBSD, and NetBSD.

Elastic IP address AWS feature enabling an EC2 user to mask the failure of an instance and remap a public IP address to any instance of the account, without the need to interact with the software support team.

Embarrassingly parallel application application when little or no effort is needed to extract parallelism and to run a number of concurrent threads with little communication among them.

Emergence generally understood as a property of a system that is not predictable from the properties of individual system components.

Energy proportional system the energy consumed by the system is proportional with its workload.

Enforced modularity software organization supported by the *client-server* paradigm when modules are forced to interact only by sending and receiving messages. The clients and the servers are independent modules and may fail separately. The servers are stateless; they do not have to maintain state information. Servers may fail and then come up without the clients being affected, or even noticing the failure.

Explicitly Parallel Instruction Computing (EPIC) processor architecture enabling the processor to execute multiple instructions in each clock cycle. EPIC implements a form of Very Long Instruction Word (VLIW) architecture.

Error bar a line segment through a point on a graph, parallel to one of the axes that represents the uncertainty or error of the corresponding coordinate of the point.

Event a change of state of a process or thread.

Event time the wall clock time when the event occurred.

Exception anomalous or exceptional conditions requiring special processing during the execution of a process. An exception breaks the normal flow of execution of a process/thread and executes a preregistered exception handler from a known memory location provided by the first-level interrupt handler (FLIV).

Exception behavior preservation condition required for dynamic instruction scheduling. Any change in instruction order must not change the order in which exceptions are raised.

F

Fabric controller a distributed Windows Azure application replicated across a group of machines that owns all resources in its environment and is aware of every application; it ensures scaling, load balancing, memory management, and reliability.

Facility layer heating, ventilation, air conditioning (HVAC), power, communications, and other aspects of the physical plant in a data center.

Failover-based software systems systems less affected by data center-level failures; such systems only run at one site, but checkpoints are created periodically and sent to backup data centers.

FedRAMP common security model enabling joint authorizations and continuous security monitoring services for government and commercial cloud computing systems intended for multiagency use. The use of this common security risk model provides a consistent baseline for cloud-based technologies and ensures that the benefits of cloud-based technologies are effectively integrated across a variety of cloud computing solutions. The risk model will enable the government to “approve once, and use often” by ensuring multiple agencies gain the benefit and insight of the FedRAMP’s authorization and access to service provider’s authorization packages.

Field-programmable gate array (FPGA) an integrated circuit designed to be configured, adapted, and programmed in the field to perform a well-defined function.

Fine-grained parallelism concurrency when only relatively small blocks of the code can be executed in parallel, without the need to communicate or synchronize with other threads or processes.

FISMA compliant environment environment that meets the requirements of the Federal Information Security Management Act of 2002. The law requires an inventory of information systems, the categorization of information and information systems according to risk level, security controls, a risk assessment, a system security plan, certification and accreditation of the system’s controls, and continuous monitoring.

First-In-First-Out delivery delivery rule requiring that messages are delivered in the same order they are sent.

First-level interrupt handler (FLIH) software component of the kernel of an operating system activated in case of an interrupt or exception. It saves the registers of current process in the PCB (Process Control Block), determines the source of interrupt, and initiates the service of the interrupt.

Flash crowds an event that disrupts the life of a very significant segment of the population, such as an earthquake in a very populated area, and dramatically increases the load of computing and communication service, e.g., an earthquake increases the phone and internet traffic.

Flynn's taxonomy classification of computer architectures proposed by Michael J. Flynn in 1966. Classifies the systems based on the number of control and data flows as: Single Instruction Single Data (SISD), Single Instruction Multiple Data (SIMD), or Multiple Instruction Multiple Data (MIMD).

Flow control mechanism used to control the traffic in a network. Feedback from the receiver forces the sender to transmit only the amount of data the receiver is able to buffer and then process.

Front-end system component of a server system tasked to dispatch the client requests to multiple *back-end* systems for processing.

Full virtualization type of virtualization in which each virtual machine runs on an exact copy of the actual hardware.

Future Internet a generic concept referring to all research and development activities involved in development of new architectures and protocols for the internet.

G

Geo replication operation in which a system runs at multiple sites concurrently.

Gather operation operation supported by vector processing units to deal with sparse vectors. It takes an index vector and fetches the vector elements at the addresses given by adding a base address to the offsets given by the index vector; as a result a dense vector is loaded into a vector register. In parallel computing it is an operation supported by MPI (Message Passing Interface) to take elements from many processes and gather them into one single process.

Global agreement on time a necessary condition to trigger actions that should occur concurrently.

Go or Golang open-source compiled, statically typed language like Algol and C; has garbage collection, limited structural typing, memory safety features, and CSP-style concurrent programming.

Guest operating system an operating system that runs under the control of a hypervisor, rather than directly on the hardware.

H

Hadoop extension of MapReduce with programming support for iterative applications and improved efficiency. Adds various caching mechanisms and makes the task scheduler loop-aware.

Hard deadline strict deadline with penalties, expressed precisely as milliseconds, or possibly seconds.

Hardware layer includes computers (CPU, memory), network (router, firewall, switch, network link, and interface) and storage components (hard disk), and other physical computing infrastructure elements.

Hash function function used to map data of arbitrary size to data of fixed size. For example, a function applied to the name of the file when the n low-order bits of the hash value give the block number of the directory where the file information can be found. *Extensible hashing* is used to add a new directory block.

Head-of-line blocking situation where a long-running task cannot be preempted and other tasks waiting for the same resource are blocked.

Hedged requests short-term tail-tolerant techniques; the client issues multiple replicas of the request to increase the chance of a prompt reply.

Hot standby a method to achieve redundancy. The primary and the secondary (backup) systems run simultaneously. The data is mirrored to the secondary system in real time so that both systems contain identical information.

Horizontal scaling application scaling by increasing the number of VMs as load increases and reducing this number when load decreases; most common form of cloud application scaling.

Hybrid cloud an infrastructure consisting of two or more clouds (private, community, or public) that remain unique entities but are bound together by standardized or proprietary technology that enables data and application portability.

HyperText Transfer Protocol (HTTP) application-level protocol built on top of the TCP transport protocol used by the web browser (the client) to communicate with the server.

HTTP-tunneling technique most often as a means of communication from network locations with restricted connectivity. Tunneling means the encapsulation of a network protocol. In this case, HTTP acts as a wrapper for the communication channel between the HTTP client and the HTTP server.

Hyper-convergence a software-centric architecture that tightly integrates compute, storage, networking, virtualization, and possibly other technologies into a commodity hardware box supported by a single vendor.

Hypervisor or virtual machine monitor (VMM) software that securely partitions the computer's resources of a physical processor into one or more virtual machines. Each virtual machine appears to be running on the bare hardware, giving the appearance of multiple instances of the same computer, but all are supported on a single physical system.

Hyper-threading term used to describe multiple execution threads possibly running concurrently but on a single-core processor.

|

Idempotent action action that repeated several times has the same effect as when the action is executed only once.

IEEE 754 Standard for Floating-Point Arithmetic defines arithmetic formats, interchange formats, rounding rules, operations, and exception handling for floating-point numbers.

Incommensurate scaling attribute of complex systems; when the size of the system, or when one of its important attributes such as speed increases, or when different system components are subject to different scaling rules.

InfiniBand switched fabric for supercomputer and data-center interconnects. The serial link can operate at several data rates: single (SDR), double (DDR), quad (QDR), fourteen (FDR), and enhanced (EDR). The highest speed supported is 300 Gbps.

Infrastructure as a Service (IaaS) cloud delivery model that supplies resources for processing, storage, and communication and allows the user to run arbitrary software, including operating systems and applications. The user does not manage or control the underlying cloud infrastructure, but has control over operating systems, storage, deployed applications, and possibly limited control of select networking components (e.g., host firewalls).

Initial value problem computational problem when all conditions are specified at the same value of the independent variable in the equation.

Input/Output Memory Management Unit (IOMMU) connects the main memory with a DMA-capable I/O bus; it maps device-visible virtual addresses to physical memory addresses and provides memory protection from misbehaving devices.

Instruction flow preservation preservation of the flow of data between the instructions producing results and the ones consuming these results.

Instruction-level parallelism simultaneous execution of independent instructions of an execution thread.

Instruction Set Architecture (ISA) interface between the computer software and hardware. It defines the valid instructions that a processor may execute. ISA allows the independent development of hardware and software.

Instruction pipelining technique implementing a form of parallelism called instruction-level parallelism within a single core or processor. A pipeline has multiple stages, and at any given time, several instructions are in different stages of processing. Each pipeline stage requires its own hardware.

Integrated Drive Electronics (IDE) interface for connecting disk drives; the drive controller is integrated into the drive, as opposed to a separate controller on, or connected to, the motherboard.

Intelligent Platform Management Interface (IPMI) standardized computer system interface developed by Intel and used by system administrators to manage a computer system and monitor its operation.

Interoperability capability to communicate, execute programs, or transfer data among various functional units under specified conditions.

Interrupt flag (IF) flag in the EFLAGS register used to control interrupt masking.

J

Jarvis short for *Just A Rather Very Intelligent Scheduler*; used to support Siri.

Java Database Connectivity (JDBC) API for Java defining how a client may access a database.

JobTracker and TaskTracker daemons handling processing of MapReduce jobs in Hadoop.

Journal storage storage for composite objects, such as records consisting of multiple fields.

Java Message Service (JMS) middleware of the Java Platform for sending messages between two or more clients.

L

Large-scale-dynamic-data data captured by sensing instruments and controlled in engineered, natural, and societal systems.

Last level cache (LLC) the cache called before accessing memory. Multicore processors have multiple level caches. Each core has its own L1 I-cache (instruction cache) and D-cache (data cache). Sometimes, two cores share the same unified (instruction+ data) L2 cache, and all cores share an L3 cache. In this case the highest shared LLC is L3.

Latch a counter that triggers an event when it reaches zero.

Late binding dynamical correlation of tasks with data, depending on the state of the cluster.

Latency the time elapsed from the instance an operation is initiated until the instance its effect is sensed. Latency is context dependent.

LRU (Least Recently Used), MRU (Most Recently Used), and LFU (Least Frequently Used) replacement policies used by memory hierarchies for caching and paging.

Livelock condition appearing when two or more processes/threads continually change their state in response to changes in the other processes and none of the processes can complete execution.

Logical clock abstraction necessary to ensure the clock condition in the absence of a global clock.

Loopback file system (LOFS) virtual file system providing an alternate path to an existing file system. When other file systems are mounted onto an LOFS file system, the original file system does not change.

M

MAC address unique identifier permanently assigned to a network interface by the manufacturer. MAC stands for Media Access Control.

Maintainability a measure of the ease of maintenance of a functional unit—synonymous with serviceability.

Malicious software (Malware) software designed to circumvent the authorization mechanisms and gain access to a computer system, gather private information, block access to a system, or disrupt the normal operation of a system; computer viruses, worms, spyware, and Trojan horses are examples of malware.

Man-in-the-middle attack attacker impersonates the agents at both ends of a communication channel making them believe that they communicate through a secure channel.

Mapping a computation assign suitable physical servers to the application.

Mashup application that uses and combines data, presentations, or functionality from two or more sources to create a service.

Megastore a scalable storage for online services.

Memcaching a general-purpose distributed memory system that caches objects in main memory.

Message-Digest Algorithm (MD5) cryptographic hash function used for checksums. MD5 produces a 128-bit hash value. SHA- i (Secure Hash Algorithm, $0 \leq i \leq 3$) is a family of cryptographic hash functions; SHA-1 is a 160-bit hash function resembling MD5.

Merkle tree hash tree in which leaves are hashes of the values of individual keys. Parent nodes higher in the tree are hashes of their respective children.

Message delivery rule an additional assumption about the channel-process interface. Establishes when a message received is actually delivered to the destination process.

Metering providing a measurement capability at some level of abstraction appropriate to the type of service.

Microkernel (μ -kernel) system software supporting only the basic functionality of an operating system kernel, including low-level address space management, thread management, and inter-process communication. Traditional operating system components, such as device drivers, protocol stacks, and file systems, are removed from the microkernel and run in user space.

Middleware software enabling computers of a distributed system to coordinate their activities and to share their resources.

Mode-sensitive instructions machine instructions whose behavior is different in the privileged mode.

Modularity basic concept in the design of man-made systems; a system is made out of components, or modules, with well-defined functions. A strong requirement for modularity is to define very clearly the interfaces between modules and to enable the modules to work together. Modularity can be *soft* or *enforced*.

Modularly divisible application application whose workload partitioning is decided a priori and cannot be changed.

Monitor process responsible for determining the state of a system.

Message Passing Interface (MPI) communication standard and communication library for a portable message-passing system.

Multi-homing a strategy to support high availability.

Multiple Instructions, Multiple Data architecture (MIMD) system with several processors/cores that function asynchronously and independently.

N

NAS Parallel Benchmarks benchmarks used to evaluate the performance of supercomputers. The original benchmark included five kernels: IS—Integer Sort; random memory access; EP—Embarrassingly Parallel; CG—Conjugate Gradient; MG—Multi-Grid on a sequence of meshes, long- and short-distance communication, memory intensive; FT—discrete 3D Fast Fourier Transform, and all-to-all communication.

Network bisection bandwidth network attribute, measures the communication bandwidth between the two partitions when a network is partitioned into two networks of the same size.

Network bisection width minimum number of links cut when dividing the network into two halves.

Network diameter average distance between all pairs of two nodes; if a network is fully connected, its diameter is equal to one.

Network Interface Controller (NIC) the hardware component connecting a computer to a Local Area Network (LAN); also known as a network interface card, network adapter, or LAN adapter.

Network layer layer of a communication network responsible for routing packets through a packet switched network from the source to the destination.

NMap a security tool running on most operating systems to map the network, that is, to discover hosts and services in the network. The systems include *Linux*, *Microsoft Windows*, *Solaris*, *HP-UX*, *SGI-IRIX*, and *BSD* variants, such as *Mac OS X*.

Nonce a random or pseudorandom number issued in an authentication protocol to ensure that old communications cannot be reused in replay attacks. Each time the authentication challenge response code is presented, the nonces are different, and so replay attacks are virtually impossible.

Nonprivileged instruction machine instruction executed in user mode.

O

Object Request Broker (ORB) the middleware that facilitates communication of networked applications.

Ontology branch of metaphysics dealing with the nature of being. Provides the means for knowledge representation within a domain; it consists of a set of domain concepts and the relationships among these concepts.

Open-box platforms traditional hardware designed for commodity operating systems that does not have the same facilities as the *closed-box platforms*.

Open Database Connectivity (ODBC) open standard application API for database access.

Overclocking technique-based on DVFS; increases the clock frequency of processor cores above the nominal rate when the workload increases.

Overlay network a virtual network superimposed over a physical network.

Over-provisioning investment in a larger infrastructure than the *typical* workload warrants.

Over-subscription the ratio of the worst-case achievable aggregate bandwidth among the servers to the total bisection bandwidth of an interconnect.

P

Packet-switched network network transporting data units called *packets* through a maze of *switches* where packets are queued and routed towards their destinations.

Pane a well-defined area within a window for the display of, or interaction with, a part of that window's application or output.

Paragon Intel family of supercomputers launched in 1992 based on the Touchstone Delta supercomputer installed at CalTech for the Concurrent Supercomputing Consortium.

Parallel slackness method of hiding communication latency by providing each processor with a large pool of ready-to-run threads, while other threads wait for either a message or for the completion of another operation.

Paravirtualization virtualization when each virtual machine runs on a slightly modified copy of the actual hardware; the reasons for paravirtualization: (i) some aspects of the hardware cannot be virtualized; (ii) to improve performance; (iii) to present a simpler interface.

Passphrase a sequence of words used to control access to a computer system; it is the analog of a password but provides added security.

Paxos protocols a family of protocols to reach consensus based on a finite-state machine approach.

Peer-to-Peer system (P2P) distributed computing system in which resources (storage, CPU cycles) are provided by participant systems.

Peripheral Component Interconnect (PCI) computer bus for attaching hardware devices to a computer. The PCI bus supports the functions found on a processor bus, but in a standardized format independent of any particular processor.

Perf profiler tool for Linux 2.6+ systems; it abstracts CPU hardware differences in Linux performance measurements.

Petri nets bipartite graphs used to model the dynamic behavior of systems.

Phase transition thermodynamics concept describing the transformation, often discontinuous, of a system from one phase/state to another, as a result of a change in the environment.

Phishing attacks aiming to gain information from a site database by masquerading as a trustworthy entity.

Physical data container storage device suitable for transferring data between cloud-subscribers and clouds.

Physical resource layer includes all physical resources used to provide cloud services.

Pinhole mapping between the pair (*external address*, *external port*) and the (*internal address*, *internal port*) tuple carried by the network address translation function of the router firewall.

Pipelining splitting of an instruction into a sequence of steps that can be executed concurrently by multiple circuitry on the chip.

Pipeline scheduling separates dependent instruction from the source instruction by the pipeline latency of the source instruction. Its effect is to reduce the number of stalls.

Pipeline stages execution units of a pipeline. A basic pipeline has five stages for instruction execution: IF = Instruction Fetch, ID = Instruction Decode, EX = Execute, MEM = Memory access, WB = Register write back.

Pipeline stall the delay in the execution of an instruction in an instruction pipeline to resolve a hazard. Such stalls could drastically affect the performance.

Platform architecture layer software layer consisting of compilers, libraries, utilities, and other software tools and development environments needed to implement applications.

Platform as a Service (PaaS) cloud delivery model supporting consumer-created or acquired applications created using programming languages and tools supported by the provider. The consumer does not manage or control the underlying cloud infrastructure, including network, servers, operating systems, or storage, but has control over the deployed applications.

Plesiochronous operation operation in which various parts of a system are almost, but not quite perfectly, synchronized, for example, when the core logic of a router operates at a frequency different from that of the I/O channels.

Pontryagin's principle method used in optimal control theory to find the best possible control that leads a dynamic system from one state to another, subject to a set of constraints.

Power consumption P of a CMOS-based circuit describes the power consumption function of the operating voltage frequency, $P = \alpha \cdot C_{eff} \cdot V^2 \cdot f$ with: α —the switching factor, C_{eff} —the effective capacitance, V —the operating voltage, and f —the operating frequency.

Privacy the assured, proper, and consistent collection, processing, communication, use, and disposition of personal information and personally identifiable information.

Private cloud infrastructure operated solely for the benefit of one organization; it may be managed by the organization or a third party and may exist on the premises or off the premises of the organization.

Privileged instructions machine instruction that can only be executed in kernel mode.

Process a program in execution.

Process group collection of cooperating processes.

Process/thread state ensemble of information needed to restart a process/thread after it was suspended.

Process or application virtual machine virtual machine running under the control of a normal OS and providing a platform-independent host for a single application, e.g., Java Virtual Machine (JVM).

Public-Key Infrastructure (PKI) model to create, distribute, revoke, use, and store digital certificates.

Pull paradigm distributed processing when resources are stored at the server site and the client pulls them from the server.

Q

Quick emulator (QEMU) a machine emulator; it runs unmodified OS images and emulates the guest architecture instructions on the host architecture it runs on.

R

Rapid provisioning automatically deploying cloud system based on the requested service, resources, and capabilities.

Recommender system system for predicting the preference of a user for an item by filtering information from multiple users regarding that item; used to recommend research articles, books, movies, music, news, and any imaginable item.

Red-black tree a self-balancing binary search tree where each node has a “color” bit (red or black) to ensure the tree remains approximately balanced during insertions and deletions.

Reference data infrequently used data such as archived copies of medical or financial records and customer account statements.

Reliability measure of the ability of a functional unit to perform a required function under given conditions for a given time interval.

Remote Procedure Call (RPC) procedure for inter-process communication. RPC allows a procedure on a system to invoke a procedure running in another address space, possibly on a remote system.

Resilience ability to reduce the magnitude and/or duration of the events disruptive to critical infrastructure.

Resilient Distributed Dataset (RDD) storage concept allowing a user to keep intermediate results and optimizes their placement in the memory of a large cluster; used for fault-tolerant, parallel data structures.

Resource abstraction and control layer software elements used to realize the infrastructure upon which a cloud service can be established, e.g., hypervisor, virtual machines, and virtual data storage.

Resource scale-out allocation of more servers to an application.

Resource scale up allocation of more resources to servers already allocated to an application.

Reservation station hardware used for dynamic instruction scheduling. A reservation station fetches and buffers an operand as soon as it becomes available. A pending instruction designates the reservation station it will send its output to.

Response time the time from the instance a request is sent until the response arrives.

Round-Trip Time (RTT) the time it takes a packet to cross the network from the sender to the receiver and back. Used to estimate the network load and detect network congestion.

Run total ordering of all the events in the global history of a distributed computation consistent with the local history of each participant process.

runC implementation of the Open Containers Runtime specification and the default executor bundled with Docker Engine.

S

Same Program Multiple Data (SPMD) parallel computing paradigm when multiple instances of one program run concurrently and each instance processes a distinct segments of the input data.

Scala general-purpose programming language with a strong static-type system and support for functional programming. Scala code is compiled as Java byte code and runs on JVM (Java Virtual Machine).

Scalability ability of a system to grow without affecting its global function(s).

Scatter operation vector processing operation, the inverse of a gather operation; it scatters the elements of a vector register to addresses given by the index vector and the base address. In distributed computing, MPI scatters data from one processor to a number of processors.

Searchable symmetric encryption (SSE) encryption method used when an encrypted databases is outsourced to a cloud or to a different organization. It supports conjunctive search and general Boolean queries on symmetrically encrypted data. SSE hides information about the database and the queries.

Secondary spectrum data falsification (SSDF) in software-defined radio, the occupancy report from a malicious node showing that channels used by the primary node are free.

Security accreditation the organization authorizes (i.e., accredits) the cloud system for processing before operations and updates the authorization when there is a significant change to the system.

Security assessment risk assessment of the management, operational, and technical controls of the cloud system.

Security certification certification for the accreditation of a cloud system.

Self-organization process in which some form of global order is the result of local interactions between parts of an initially disordered system. No single element acts as a coordinator, and the global patterns of behavior are distributed.

Semantic Web term coined by Tim Berners-Lee to describe “a web of data that can be processed directly and indirectly by machines.”

Sensitive instructions machine instructions behaving differently when executed in kernel and in user mode.

Sequential write-sharing condition when a file cannot be opened simultaneously for reading and writing by several clients.

Service aggregation operation when an aggregation brokerage service combines multiple services into one or more new services.

Service arbitrage/service aggregation grouping of cloud services. In service aggregation, the services being aggregated are not fixed. Arbitrage provides flexibility and opportunistic choices for the service aggregator, e.g., provides multiple e-mail services through one service provider, or provides a credit-scoring service that checks multiple scoring agencies and selects the best score.

Service deployment activities and organization needed to make a cloud service available.

Service intermediation operation when an intermediation broker provides a service that directly enhances a given service delivered to one or more service consumers.

Service interoperability the capability to communicate, execute programs, or transfer data among various cloud services under specified conditions.

Service layer defines the basic services provided by cloud providers.

Service Level Agreement (SLA) a negotiated contract between the customer and the service provider explaining expected quality of service and legal guarantees. An agreement usually covers: services to be delivered, performance, tracking and reporting, problem management, legal compliance and resolution of disputes, customer duties and responsibilities, security, handling of confidential information, and termination.

Shard a horizontal partitioning of a database; a row in a table structured data.

Shared channel architecture network organization when all physical devices share the same bandwidth; the higher the number of devices connected to the channel, the less bandwidth is available to each one of them.

Simple Object Access Protocol (SOAP) an application protocol developed in 1998 for web applications.

Singular Value Decomposition (SVD) given an $m \times n$ matrix $A = [a_{ij}]$, $1 \leq i \leq n$, $1 \leq j \leq m$ with entries either real or complex numbers, $a_{ij} \in \mathbb{R}$ or $a_{ij} \in \mathbb{C}$, there exists a factorization

$$A = U \Sigma V^* \quad (1)$$

with: U is an $m \times n$ unitary matrix, Σ is a diagonal $m \times n$ matrix with nonnegative real numbers on the diagonal, V is an $n \times n$ unitary matrix over the field, \mathbb{R} or \mathbb{C} , and V^* is the complex conjugate transpose of V .

- SLA management** the ensemble of activities related to SLAs, including SLA contract definition (basic schema with the quality of service parameters), SLA monitoring, and SLA enforcement.
- Service management** all service-related functions necessary for the management and operations of those services required by customers.
- Service orchestration** the arrangement, coordination, and management of cloud infrastructure to provide multiple cloud services to meet IT and business requirements.
- Service provider** entity responsible for making a service available to service consumers.
- Serverless computer service** AWS service when applications are triggered by conditions and/or events specified by the end user. Lambda is an example of such service.
- Shared cluster state** a resilient master copy of the state of all cluster resources.
- Sigmoid function $S(t)$** an “S-shaped” function defined as $S(t) = \frac{1}{1-e^{-t}}$. Its derivative can be expressed as a function of itself, $S'(t) = S(t)(1 - S(t))$.
- Simple DB** AWS nonrelational data store that allows developers to store and query data items via web services requests; it creates multiple geographically distributed copies of each data item and supports high-performance web applications.
- Simple Queue Service (SQS)** AWS service for hosted message queues. It allows multiple EC2 instances to coordinate their activities by sending and receiving SQS messages.
- Simple Storage System (S3)** AWS storage service for large objects. It supports a minimal set of functions: write, read, and delete. S3 allows an application to handle an unlimited number of objects, ranging in size from one byte to five terabytes.
- Single Instruction, Multiple Data architecture (SIMD)** computer architecture in which one instruction processes multiple data elements. Used in vector processing.
- Single Instruction, Single Data architecture (SISD)** computer architecture supporting the execution of a single thread or process at any given time. Individual cores of a modern multicore processor are SISD.
- Soft modularity** dividing a program into modules that call each other and communicate using shared memory or follow the procedure call convention. It hides the details of the implementation of a module. Once the interfaces of the modules are defined, the modules can be independently developed even in multiple programming languages, replaced, and tested.
- Software as a Service (SaaS)** cloud delivery model in which cloud applications are accessible from various client devices through a thin client interface, such as a web browser. The user does not manage or control the underlying cloud infrastructure.
- S/KEY** password system based on the Leslie Lamport scheme. The real password of the user is combined with a short set of characters and a counter that is decremented at each use to form a single-use password. Used by several operating systems, including Linux, OpenBSD, and NetBSD.
- Skype** communication system using a proprietary voice-over-IP protocol. The system was developed in 2003 and acquired by Microsoft in 2011. Nowadays, it is a hybrid P2P and client–server system. It allows close to 700 million registered users from many countries around the globe to communicate.
- Simultaneous multithreading (SMT)** architectural feature allowing instructions from more than one thread to be executed in any given pipeline stage at the same time.
- Simple Mail Transfer Protocol (SMTP)** application protocol defined in the early 1980s to support email services.
- Snapshot isolation** guarantee that all reads made in a transaction will see a consistent snapshot of the database.
- Soft deadlines** deadline that can be missed by fractions of the units. It is more of a guideline; no penalties are involved.
- Software development kit (SDK)** a set of software tools for the creation of applications in a specific software environment.
- Speed** term used informally to describe the maximum data transmission rate, or the capacity of a communication channel; this capacity is determined by the physical bandwidth of the channel, and this explains why the term channel “bandwidth” is also used to measure the channel capacity, or the maximum data rate.
- Speedup** measure of parallelization effectiveness.
- SQL injection** attack typically used against a web site; an SQL command entered in a web form causes the contents of a database used by the web site to be altered or to be dumped to the attacker site.
- ssh (Secure Shell)** network protocol that allows data to be exchanged using a secure channel between two networked devices; ssh uses public-key cryptography to authenticate the remote computer and allow the remote computer to authenticate the user. It also allows remote control of a device.
- Store-and-forward network** packet switched network where a router buffers a packet, verifies its checksum, and then forwards it to the next router along the path from its source to the destination.
- Streaming SIMD Extension (SSE)** SIMD instruction set extension to the x86 architecture introduced by Intel in 1999. Its latest expansion is SSE4. It supports floating-point operations and has a wider application than the MMX introduced in 1996.

Structural hazards in pipelining hazards occurring when a part of the processor hardware is needed by two or more instructions at the same time.

Structured overlay network network where each node has a unique key that determines its position in the structure. The keys are selected to guarantee a uniform distribution in a very large name space. Structured overlay networks use *key-based routing* (KBR); given a starting node v_0 and a key k , the function $KBR(v_0, k)$ returns the path in the graph from v_0 to the vertex with key k .

Structured Query Language (SQL) special-purpose language for managing structured data in a relational database system (RDBMS). SQL has three components: a data definition language, a data manipulation language, and a data control language.

Superscalar processor processor able to execute more than one instruction per clock cycle.

System history information about the past system evolution expressed as a sequence of events, each event corresponding to a change of the state of the system.

System portability the ability of a service to run on more than one type or size of cloud.

T

Task-level parallelism parallelism when the tasks of an application run concurrently on different processors. A job consists of multiple tasks scheduled either independently or co-scheduled when they need to communicate with one another.

TCP segmentation offload (TSO) procedure enabling a network adapter to compute the TCP checksum on transmit and receive; it saves the host CPU the overhead for computing the checksum; large packets have larger savings.

Tera Watt Hour (TWh) measure of energy consumption; one TWh is equal to 10^9 KWh.

Thread of execution the smallest unit of processing that can be scheduled by an operating system.

Thread-level parallelism term describing the data-parallel execution using a GPU. A thread is a subset of vector elements processed by one of the lanes of a multithreaded processor.

Thread block scheduler GPU control software assigning thread blocks to multithreaded SIMD processors.

Thread scheduler GPU control software running on each multithreaded SIMD processor to assign threads to the SIMD lanes.

Three-way handshake process to establish a TCP connection between the client and the server. The client provides an arbitrary initial sequence number in a special segment with the *SYN* control bit on; then, the server acknowledges the segment and adds its own arbitrarily chosen initial sequence number; finally, the client sends its own acknowledgment *ACK* and also the HTTP request, and the connection is established.

Threshold value of a parameter related to the system state that triggers a change in the system behavior.

Thrift framework for cross-language services.

Timestamps patterns used for event ordering using a global time-base constructed on local virtual clocks.

Top-Down methodology hierarchical organization of event-based metrics that identifies the dominant performance bottlenecks in an application.

TPC BenchmarkH (TPC-H) decision-support benchmark relevant for applications that examine large volumes of data and execute queries with high degree of complexity; it consists of a suite of business-oriented ad hoc queries and concurrent data modifications with broad industry-wide relevance.

TPC-DS de facto industry standard benchmark for assessing the performance of decision-support systems.

Trusted application application with special privileges for performing security-related functions.

Translation look-aside buffer (TLB) cache for dynamic address translation; it holds the physical address of recently used pages in virtual memory.

Transport layer network layer responsible for end-to-end communication, from the sending host to the destination host.

Trusted Computer Base (TCB) totality of protection mechanisms within a computer system, including hardware, firmware, and software, the combination of which is responsible for enforcing a security policy.

Turing complete computer model of computation equivalent to a universal Turing machine, except for memory limitations.

U

Ubuntu open-source operating system for personal computers. Ubuntu is an African humanist philosophy; “ubuntu” is a word in the Bantu language of South Africa meaning “humanity towards others.”

Unbounded input data concept related to data streaming; the computing engine processes a dynamic data set when one never knows if the set is complete because new records are continually added and old ones are retracted.

Usability extent to which a product can be used by specified users to achieve specified goals with effectiveness, efficiency, and satisfaction in a specified context of use.

Utility function relates the “benefits” of an activity or service with the “cost” to provide the service.

V

Vector computer computer operating on vector registers holding as many as 64 or 128 vector elements. Vector functional units carry out arithmetic and logic operations using data from vector registers as input and disperse the results back to memory.

Vector length register register of a SIMD processor for handling of vectors whose length is not a multiple of the length of the physical vector registers.

Vector mask register register of a SIMD processor used by conditional statements to disable/select vector elements.

Vertical scaling method to increase the resources of a cloud application. It keeps the number of VMs of an application constant, but increases the amount of resources allocated to each one of them.

Virtual Machine (VM) an isolated environment with access to a subset of the physical resources of a computer system. Each virtual machine appears to be running on the bare hardware, giving the appearance of multiple instances of the same computer, though all are supported by a single physical system.

Virtual Private Cloud (VPC) cloud organization providing a connection, via a Virtual Private Network, between an existing IT infrastructure of an organization and a set of isolated compute resources in the AWS cloud.

Virtual time warp abstraction allowing a thread to acquire an earlier effective virtual time, in other words, to borrow virtual time from its future CPU allocation.

Virtualization abstraction of hardware resources.

Virtualized infrastructure layer software elements, such as hypervisors, virtual machines, virtual data storage, and supporting middleware components, used to realize the infrastructure upon which a computing platform can be established. While virtual machine technology is commonly used at this layer, other means of providing the necessary software abstractions are not precluded.

W

WebSphere Extended Deployment (WXD) middleware supporting setting performance targets for individual web applications and for monitoring response time.

Where-provenance information describing the relationship between the source and the output locations of data in a database.

Why-provenance information describing the relationship between the source tuples and the output tuples in the result of a database query.

Wide Area Network (WAN) packet switched network connecting systems located throughout a very large area.

Witness of database record the subset of database records ensuring that the record is the output of a query.

Work-conserving scheduler scheduler with the goal keeping the resources busy if there is work to be done; a *non-work conserving scheduler* may leave resources idle while there is work to be done.

Work-conserving scheduling policy scheduling policy when the server cannot be idle while there is work to be done.

Workflow description of a complex activity involving an ensemble of multiple interdependent tasks.

Write-ahead database technique that writes updates to persistent storage only after the log records have been written.

X

x86-32, i386, x86 and IA-32 CISC-based instruction set architecture of Intel processors. Now supplanted by x86-64, which supports vastly larger physical and virtual address spaces. The x86-64 specification is distinct from Itanium, initially known as IA-64 architecture.

x86 architecture architecture of Intel processors supporting memory segmentation with a segment size of 64 K. The CR (code-segment register) points to the code segment. *MOV*, *POP*, and *PUSH* instructions serve to load and store segment registers, including CR.

Z

Zero-configuration networking (zeroconf) computer network based on the TCP/IP and characterized by automatic assignment of numeric network addresses for networked devices, automatic distribution and resolution of computer hostnames, and automatic location of network services.

ZooKeeper a distributed coordination service implementing a version of the Paxos consensus algorithm.

Index

A

Absolute deadline, 316
Acceptor, 398
Access Control List (ACL), 151, 227, 423
Access layer, 200
Access transparency, 73
Account, 261
Action, 106, 151
Active space, 421
Activities, 416
Actual virtual time, 311
Actuators, 334
Advanced Microcontroller Bus Architecture (AMBA), 65
Advanced Research Projects Agency (ARPA), 93
Advanced RISC Machine (ARM), 159
Advisory locks, 234
Aggregate layer, 200
Aggregation, 136
Aggregation function, 458
Algorithm
 allocation, 322
 edge, 126
 generic contraction, 556
 generic TLC, 554
 network resource management, 204
 Paxos, 236, 237
 rate monotonic, 298
Algorithmic communication, 354
Allocate thread, 385
Allocation algorithms, 322
Allocation manager, 230
Amazon, 15, 18, 23, 24
Amazon Machine Image (AMI), 19, 22, 137, 278, 430
Amazon Web Services (AWS), 1, 19
Amazon Web Services Licensing Agreement (AWSLA), 24
Amdahl's Law, 68
Amplitude amplifications, 511
Analytics, 298
Andrew File System (AFS), 216, 227
Ansible, 18
Anti-entropy, 250
Anticipated exceptions, 415
Apache
 capacity scheduler, 306
 Hadoop, 105, 320, 428, 431, 432
Aperiodic tasks, 316

Application

arbitrarily divisible, 329
coarse-grained data-parallel, 104
controllers, 335
data, 2, 6
data-parallel, 99
dynamic data-driven, 463
embarrassingly parallel, 68
layer, 16, 80
modularly divisible, 329
request-parallel, 99
trusted, 272
VM, 138
Application Binary Interface (ABI), 82
Application Program Interface (API), 82, 425
Application Specific Integrated Circuit (ASIC), 63
Application/logic tier, 4
Arbitrarily divisible application, 329
Arbitrarily divisible load sharing model, 412
Architecture, 42, 48
 ARM, 55, 57
 CISC, 43, 48
 computer, 42, 140
 taxonomy, 46
 control flow, 43
 data flow, 43, 92
 hourglass, 178
 Itanium, 142, 161, 173
 load-store, 43
 network, 154, 176
 ONF, 188
 scalable data center communication, 200
 SIMD, 58–60
 software, 4
 vector, 59
 von Neumann, 91
Arithmetic intensity (ArI), 60
ARM, 450
 architecture, 55, 57
 processors, 159, 160
ARM Trust Zone (ATZ), 160
Array, 98
 switch, 98
Artificial Intelligence (AI), 3, 11, 25, 501, 503
Assured rate (AR), 207
Asymmetric confusion, 364

- Asymmetric core processor, 70
 Atlimit, 207
 Atomic transaction, 111
 Atomicity, 394
 Atomicity log, 236
 Audit, 15
 Auditor, 14
 Auto Scaling, 22
 Automatic speech recognition (ASR), 507
 Autonomic performance managers, 336
 Autoscaling, 329
 Availability, 74, 247, 271, 524
 - cloud services, 9, 260
 Availability zone (AZ), 23
 Average memory access time (AMAT), 51
 AWS, 19, 24
 - networking, 23
 - security, 258, 289, 290
- B**
- Back-end bound, 103
 Back-propagation, 504
 Bad speculation, 103
 Balloon driver, 150
 Balloon process, 104
 Balls-and-bins model, 400
 Bandwidth, 42, 88, 204, 293
 - full bisection, 189
 - memory, 48, 60
 Bare metal hypervisor, 139
 Barrier synchronization, 352, 354, 355
 Base images, 129
 Basic Core Equivalent (BCE), 70
 Basic Paxos, 398
 Benchmarks, 47
 Best-effort delivery, 178
 Biased strategy, 53
 Bid, 205
 Bigquery, 26
 BigTable, 243, 255, 456
 Binary translation, 140
 Bipartite graph, 363
 Bisection width, 189
 Bit-by-bit Round-robin (BR), 204
 Bit-level parallelism, 45
 Black Widow topology, 193, 212
 Bloch sphere, 510
 Block, 61
 - cache, 50, 228
 - data, 50
 Blockchain, 89
 Bootstrap method, 492
 Bootstrap Performance Diagnostic (BPD), 493
 Bootstrap substitution principle, 492
 Border Gateway Protocol (BGP), 26
 Borg, 107
 Borglets, 108
 BorgMaster, 108–110
 Borrowed virtual time (BVT), 150, 311, 347
 Borrowing, 207
 Bound
 - back-end, 103
 - front-end, 103
 Branch condition speculation, 53
 Branch prediction, 53
 Branch target speculation, 53
 Breath first search (BFS), 449
 Bridge, 154
 Broker, 14
 Browny cores, 99
 Bucket, 21, 437, 460, 582
 Buffer space, 204
 Bulk synchronous parallel model, 360
 Bundles, 322
 Bundling, 278
 Business intelligence, 456
 Byte-range tokens, 230
 Byzantine attack, 530
- C**
- Cache, 49
 - banks, 51
 - block, 50, 228
 - capacity, 49, 50
 - data, 251
 - direct map, 50
 - file, 227, 228, 440
 - fully associative, 50
 - hits, 50, 51
 - implementation, 50
 - isolation, 141
 - L1D, 49
 - line, 50, 55
 - miss, 50, 136
 - miss rate, 51
 - n-way set associative, 50
 - on-chip L2, 49
 - row, 50
 - set-associative, 50
 Cache Allocation Technology (CAT), 118
 Cache Array Routing Protocol (CARP), 209
 Cache Monitoring Technology (CMT), 118
 Calculus-based trust, 267
 Callback, 106
 Callstack, 100
 CAP theorem, 88

- Carrier, 14
- Causal message delivery, 378
- Cause–effect relationships, 375
- CDN, 209
 - providers, 209
- Ceil rate (CR), 207
- Cell, 98, 108
 - storage, 220, 395, 396
- Central Limit Theorem (CLT), 495
- Central processing unit (CPU), 358
- Centralized control, 293
- Certificate Authority (CA), 32
- Chain, 404, 576
- Chaining, 59
- Chameleon, 38
- Checkpoint file, 110
- Checkpoint–restart procedures, 350
- Child images, 129
- Chroot, 97, 128, 164
- Chubby files, 236, 243
- Chubby locks, 236
- Chunk, 231
 - column, 439
 - handle, 231
 - servers, 232, 233
- Circular wait, 394
- CISC architecture, 43, 48
- Class of service (COS), 119
- Class-Based Queueing (CBQ), 206
- Classification engine, 115
- Classifier, 206
- Cleanup phase, 430
- Client
 - library, 461
 - thin, 11
- Client–server paradigm, 41, 75–77, 225, 413
- Clock
 - condition, 380
 - processor, 44
- Clos networks, 193
- Cloud
 - application development, 412
 - Bigtable, 26
 - CDN, 26
 - community, 4, 344
 - computer, 4
 - computing
 - delivery models, 14
 - infrastructure, 13, 63, 95, 246, 357, 412, 503
 - controller, 167, 335
 - controls matrix, 263
 - data encryption, 268
 - database, 222
 - DBaaS, 16
 - DNS, 26
 - elasticity, 22, 293, 325, 347, 411, 503
 - functionality, 502
 - Google, 25
 - hardware, 98
 - infrastructure, 455
 - hybrid, 4
 - IBM, 30
 - Identity and Access Management (IAM), 27
 - interconnection networks, 191, 503
 - interoperability, 13, 31, 32, 294, 344
 - KMS, 27
 - Machine Learning, 27, 506
 - Natural Language API, 27
 - oracle, 15
 - private, 4, 5, 167, 168
 - public, 4
 - scheduling subject, 294, 315
 - SDK, 27
 - security, 15, 30, 36, 258, 259, 263, 286, 289–291, 452
 - risks, 259
 - scanner, 27
 - service availability, 9, 260
 - source repositories, 27
 - speech API, 27
 - SQL, 26
 - storage, 26, 216, 255, 257, 289
 - diversity, 30
 - virtual, 316
 - vision API, 27
- Cloud Functions (CF), 26
- Cloud Load Balancing, 26
- Cloud Security Alliance (CSA), 261, 290
- Cloud Service Providers (CSP), 1, 32, 96, 135
- Cloud Virtual Network (CVN), 26
- CloudHSM, 290
- CloudLab, 38
- Cloudlets, 480
- CloudStore, 233, 431
- CloudTrail, 129, 289
- CloudWatch, 21, 22
- Cluster
 - controllers, 167
 - management, 107
 - software, 101
 - system, 110
 - scalable, 133
- Co-scheduling, 354
- Coarse control, 335
- Coarse-grained
 - data-parallel applications, 104
 - locks, 234

- parallelism, 354
- transformations, 121
- Cognitive Radio Networks (CRN), 529
- Cognito, 290
- Colored Petri Nets (CPN), 369
- Colossus File System (CFS), 461
- Column chunk, 439
- Column-family databases, 240
- Combinatorial auctions, 322
- Commission failure, 398
- Commit point, 395
- Committer, 459
- Common data bus (CDB), 52
- Common Object Request Broker Architecture (CORBA), 414
- Communicating sequential processes (CSP), 359, 408
- Communication, 353
 - algorithmic, 354
 - bandwidth, 118, 200, 362
 - channel, 369, 388
 - complexity, 354
 - events, 370, 371, 376, 377, 382
 - infrastructure, 175, 197, 200, 209
 - latency, 24, 87, 99, 353
 - location transparent, 191
 - protocols, 80, 84, 101, 117, 118, 349, 354, 374, 388, 525
 - quantum, 509
- Community cloud, 4
- Complex Instruction Set Computer (CISC), 45
- Compliance Level Agreements (CLA), 34
- Compliant Cloud Computing (C3), 34
- Composability bounds, 86
- Composite services, 450
- Composite task, 415
- Computational models, 358
- Compute Engine (CE), 26
- Computer
 - architecture, 42, 140
 - architecture taxonomy, 46
- Computing
 - grid, 85
 - network-centric, 10
 - utility, 1
 - volunteer, 449
- Concurrency, 349, 350, 353, 358, 361, 387, 388, 392
 - control, 228, 392, 460, 461
 - transparency, 73
- Concurrent events, 376
- Concurrent instructions, 46
- Confidentiality, 271
- Conflict fraction, 112
- Conflicting requirements, 502
- Confusion, 364
 - asymmetric, 364
 - symmetric, 364
- Congestion avoidance (CA), 184
- Congestion control, 180, 183, 185, 211, 344, 375
- Conjugate gradient (CG), 444
- Connection-oriented protocol, 183
- Connectionless protocols, 183
- Consensus service, 397
- Consistency
 - external time, 222
- Consistent cut, 382
- Consistent message delivery, 379
- Consumer service, 14
- Container, 129
- Container Engine (CntE), 26
- Container Registry, 26
- Containers, 95, 97, 128, 130, 434
 - Docker, 26, 128, 129
- Content, 10
 - network-centric, 10
- Content Delivery Network (CDN), 28, 207, 208
- Content service network (CSN), 211
- Content Store, 186
- Content-centric distribution, 211
- Content-Centric Networks (CCNs), 185, 210
- Context switch allowance, 311
- Control
 - coarse, 335
 - concurrency, 228, 392, 460, 461
 - congestion, 375
 - data, 30, 142
 - error, 374
 - fine, 335
 - flow, 43, 356, 374
 - granularity, 335
 - hazard, 46
 - hypervisor, 104
 - integral, 335
 - optimal, 330, 332
 - planes, 10, 161
 - system, 333, 334, 463
 - task, 331
 - theory, 296, 328–330, 333, 335, 346
 - third-party, 261
- Control Data Corporation (CDC), 72
- Control Unit (CU), 44
- Controller, 329, 333, 347, 458
 - application, 335
 - cloud, 167, 335
 - cluster, 167
 - fabric, 29
 - replication, 131
- Convolutional Neural Networks (CNN), 506
- Cooperative spectrum-sensing strategy, 530

- Coordinator, 421
 - Core
 - browny, 99
 - layer, 200
 - network, 176, 209
 - wimpy, 99
 - Correlating predictors, 54
 - Cost
 - memory, 48
 - network, 194
 - Costate equation, 331
 - CPU
 - optimization, 51
 - registers, 49
 - virtualization, 163
 - Critical section, 392, 396, 405, 407
 - Critical vulnerabilities, 54
 - Cron, 107
 - Crunch pipeline, 407
 - CSPs, 1, 15, 32, 96, 135, 335, 344, 345
 - Cumulative distribution function (CDF), 484
 - Current Privilege Level (CPL), 163
 - Current Program Status Register (CPSR), 56, 57
 - Customer, 181
 - Cut, 381
 - consistent, 382
 - frontier, 381
 - Cyberwarfare, 257
 - Cycle per Instruction (CPI), 44
- D**
- D-Wave, 518
 - Daemon, 105, 129
 - Data
 - application, 2, 6
 - block, 50
 - cache, 251
 - control, 30, 142
 - flow, 44, 180, 204, 233, 356, 357, 369, 409
 - architecture, 43, 92
 - generality, 75
 - hazards, 46
 - leakage, 261
 - link layer, 177
 - loss, 261
 - mashups, 450
 - memory flow, 52
 - networks, 210
 - sharing abstraction, 121
 - storage layer, 16, 498
 - streaming, 466, 469, 470
 - applications, 454, 466
 - pipeline, 483
 - services, 483
 - streams, 46
 - tier, 4
 - training, 504
 - warehouse, 456
 - Data Access Units (DAU), 464
 - Data center networks (DCN), 200
 - Data Description Language (DDL), 438
 - Data Manipulation Language (DML), 438
 - Data Mining (DM), 501
 - Data Oriented Network Architecture (DONA), 185
 - Data-aware scheduling, 303
 - Data-level parallelism, 67
 - Data-parallel applications, 99
 - Data-shipping, 230
 - Database, 438
 - cloud, 222
 - column-family, 240
 - document, 240
 - graph, 241
 - layer, 16
 - NoSQL, 222, 239–241, 243, 245, 247, 255, 269, 454
 - provenance, 216, 252
 - schema, 454
 - systems, 247, 269
 - management, 222
 - transactions, 255, 422
 - Database as a Service (DBaaS), 16
 - Datagram, 176
 - Datalog conjunctive query, 253
 - Dataspaces, 28, 499
 - De-perimeterization, 9
 - Deadline
 - absolute, 316
 - relative, 315
 - Deadlocks, 350, 367, 392
 - Decentralized consensus, 90
 - Decision models, 33
 - Decoherence, 513
 - Deep neural networks, 63
 - Deep Neural Networks (DNN), 506
 - Defense Information Systems Agency (DISA), 86
 - Deferred evaluation, 406
 - Degree, 189
 - Delay scheduling, 298, 300
 - Delay slots, 52
 - Delayer, 206
 - Delivery
 - best-effort, 178
 - consistent message, 379
 - First-In-First-Out (FIFO), 378
 - rule, 378
 - Denali hypervisor, 138

- Denial-of-service (DOS), 282, 575
- Deutsch–Jozsa problem, 515
- Device, 61
- DHCP server, 179
- Direct map cache, 50
- Directed acyclic graph (DAG), 125
- Disk locality, 251
- Disk maps, 230
- Disruptive computing paradigm, 2
- Disruptive technology, 39
- Distributed computing fallacies, 88
- Distributed Data Interface (DDI), 443
- Distributed sort, 426
- Distributed system, 73
- Distribution
 - content-centric, 211
 - heavy-tail, 484
 - network, 185
- DNS
 - cloud, 26
 - name, 19
- DNS-based routing, 210
- Docker containers, 26, 128, 129
- Dockerfile, 129
- Domain, 149
 - driver, 154
 - guest, 154
- Domain Specific Architectures (DSA), 62, 507
- DomU, 282
- Double data rate (DDR), 195
- DRAM, 49, 65, 216
 - bandwidth, 117, 118, 120
 - latency, 99
- Driver, 438
 - domain, 154
 - split, 151
- Dynamic
 - application scaling, 328
 - binary translation, 82
 - data-driven applications, 463
 - frequency scaling, 117
 - thresholds, 334
 - trust, 267
 - voltage scaling, 117
 - workflows, 421
- Dynamic Host Configuration Protocol (DHCP), 179, 212
- Dynamic Random Access Memory (DRAM), 216
- Dynamic voltage and frequency scaling (DVFS), 117
- DynamoDB, 22, 247
- E**
- Earliest Deadline First (EDF), 298
- Earliest Virtual Time (EVT), 311
- Ecological impact, 7
- ECS, 129
- Edge
 - algorithm, 126
 - network, 176
- Effective Machine Utilization (EMU), 133
- Effective virtual time, 311
- Elastic Beanstalk, 22
- Elastic Block Store (EBS), 21
- Elastic Compute Cloud (EC2), 19
- Elastic Container Service (ECS), 147
- Elastic Fabric Adapter (EFA), 61
- Elastic IP address, 19, 20
- Elastic Kubernetes Service (EKS), 147
- Elastic Load Balancer (ELB), 10, 22
- Elastic load balancing, 20
- Elastic MapReduce (EMR), 22
- ElastiCache, 22
- Elasticity, 325
- Embarrassingly parallel application, 68, 353
- Embarrassingly parallel problems, 67
- Emergence, 342
- Emulation, 136
 - processor, 137
 - system, 145
 - user-mode, 145
- Enactment engine, 416
- Enactment model, 416
- Energy efficiency, 325, 326
- Energy saving, 327
- Energy use, 7
- Energy-proportional systems, 326
- Enforced modularity, 75
- Enhanced data rated (EDR), 195
- Enterprise flash drive (EFD), 218
- Enterprise Mobile Management (EMM), 287
- Entity groups, 245, 246
- Equal Partitioning Rule (EPR), 319
- Error, 122
 - control, 374
 - sampling, 495
- Estimator, 206
- Eucalyptus, 167
- Event, 369, 376, 470
 - time, 471
- Eventual consistency, 249
- Exception bitmap, 144
- Exception Handler (EH), 388
- Excess vector, 323
- Execute instruction, 52
- Executer functions, 106
- Execution time
 - parallel, 68

- sequential, 68
- Executor, 106, 428
- Extended Nets, 367
- Extensibility, 414
- Extensible hashing techniques, 229
- External Interface, 438

- F**
- Fabric controller, 29
- Facts, 253
- Failure transparency, 73
- Fair Queuing (FQ), 204
- Fair share timeout, 302
- Fair-share scheduler, 127
- Fast Fourier Transform (FFT), 511
- Fast retransmit (FR), 184
- Fat-tree, 192, 193
- Fate sharing, 77
- FauxMaster, 110
- Feasibility, 109
- Feasible state, 332
- Federal Information Processing Standard (FIPS), 261
- Federal Information Security Management Act (FISMA), 261
- Fiber-to-the-home (FTTH), 182
- Fibre Channel (FC), 197
- Field Programmable Gate Arrays (FPGA), 67
- File
 - cache, 227, 228, 440
 - checkpoint, 110
 - Chubby, 236, 243
 - local, 225
 - lock, 236
 - remote, 225, 227
- Filter, 122
- Fine control, 335
- Fine-grained cluster resource sharing, 106
- Fine-grained locks, 234
- Fine-grained parallelism, 66, 354, 355
- Finite state machine (FSM), 53
- Firing sequence, 367
- First-In-First-Out (FIFO) delivery, 378
- Fixed parallel time per process, 69
- Fixed problem size, 69
- Flash crowd, 209
- Flash memory, 217
- Flattened butterfly network, 193
- Floating-point performance models, 60
- Flow
 - control, 183, 184, 189, 197, 296, 356, 374
 - data, 44, 180, 204, 233, 356, 357, 369, 409
 - memory, 52
 - register, 52
 - instruction, 52
- relations, 366
- Flow-control mechanism, 179
- Flowchart, 416
- FlumeJava, 406, 407, 428
 - library, 406
- Fork routing task, 416
- Forwarding Information Base, 187
- Forwarding plane, 178
- Front-end bound, 103
- Frontier of the cut, 381
- Full bisection bandwidth, 189
- Full virtualization, 135, 140
- Fully associative cache, 50
- Fully Homomorphic Encryption (FHE), 268
- Function
 - aggregation, 458
 - objective, 322
 - utility, 338
 - watch, 151

- G**
- Gather operation, 59
- General Parallel File System (GPFS), 229, 255
- Generalized Processor Sharing (GPS), 339
- Generic contraction algorithm, 556
- Generic TLC algorithm, 554
- Geo-replication, 457
- Geometric engine, 67
- GFS files, 231
- Global predicate evaluation problem (GPE), 374
- Global real-time clock, 379
- Gmail, 27
- Google
 - Base, 28
 - Calendar, 26, 27, 34, 37
 - Cloud, 25
 - infrastructure, 25, 98
 - Co-op, 28
 - Docs, 27
 - Drive, 26, 28, 39
 - Groups, 28
 - Maps, 28, 37
- Google App Engine, 26
- Google Container Engine (GKE), 129
- Google file system (GFS), 216, 256
- Gradient, 304
 - descent, 504
- Grand architectural complications, 48
- Graph
 - bipartite, 363
 - databases, 241
 - marked, 365, 368
- Graphics Processing Unit (GPU), 2, 60, 67

- GrepTheWeb application, 428, 451
 Ground truth, 493
 Guest
 - domain, 154
 - hypervisor, 142, 156, 158, 159
 - state, 142, 143
- H**
- Hadoop, 115, 431, 435
 - Apache, 105, 320, 428, 431, 432
 - master, 298
 - scheduler, 298, 302
 - system, 431
 - tasks, 299, 320
 - workloads, 115
- Hadoop Distributed File System (HDFS), 105
 Hadoop fair scheduler (HFS), 302
 Hard disk, 49
 Hard disk drive (HDD), 216, 217
 Hardware
 - cloud, 98
 - costs, 137, 171, 223
 - parallel, 68
 - parallelism, 42
 - support, 136, 141
 - virtual, 97, 137, 169
 - virtualization, 146, 171
- Hardware description language (HDL), 67
 Hashing
 - perfect random, 404
- Hazard
 - control, 46
 - data, 46
 - pipeline, 46
 - structural, 46
- HDFS, 105, 431
 Head node, 316
 Heavy-tail distributions, 484
 Heracles, 119
 Hidden states, 143
 Hiding the latency, 42, 43, 461
 Hierarchical coordination scheme, 421
 Hierarchical Token Buckets (HTB), 207
 Hierarchy, 79
 - of networks, 98
- High Performance Computing (HPC), 432
 History-based prediction, 53
 Hit time, 51
 Hive, 437
 Homomorphic encryption, 269
 Horizontal scaling, 242, 329
 Host, 176
 - hypervisor, 156
 - state, 142, 143
- Hosted, 139
 - independent hypervisors, 171
 - specialized hypervisors, 171
- Hosting
 - KVM, 145
 - Xen, 145
- Hostname, 19
 Hourglass architecture, 178
 HTTP redirection, 210
 HTTP server, 77
 HTTP-tunneling, 79
 Hub, 129
 Hybrid, 139
 - cloud, 4
- Hyper-threading, 47, 67
 Hypercalls, 150
 Hypervisor
 - bare metal, 139
 - Denali, 138
 - guest, 142, 156, 158, 159
 - host, 156
 - hosted independent, 171
 - hosted specialized, 171
 - native, 170
 - rogue, 169
 - Xen, 19, 138, 155, 164, 172, 281
- I**
- IBM clouds, 30
 Image, 129
 - base, 129
 - child, 129
 - official, 129
 - user, 129
- Impala, 438
 In-memory cluster computing, 120
 Incidence Matrix, 367
 Incommensurate scaling, 343
 Independence, 414
 Index vector, 59
 InfiniBand, 194, 195
 Information, 10
 - hiding, 74
- Information Management System (IMS), 254
 Information Processing Technology Office (IPTO), 93
 Infrastructure
 - cloud computing, 13, 63, 95, 246, 357, 412, 503
 - communication, 175, 197, 200, 209
 - Google Cloud, 25, 98
 - network, 82, 181
 - networking, 94, 191
 - virtualization, 136, 145

- Infrastructure as a Service (IaaS), 16
 - Infrastructure as Code (IaC), 18
 - Inhibitor arc, 367
 - Initial events, 359
 - Insecure APIs, 261
 - Inspector (IS), 289
 - Instance, 547
 - Instruction
 - concurrent, 46
 - flow, 52
 - issue, 52
 - privileged, 81, 162
 - privileged-sensitive, 163
 - sensitive, 81, 173
 - Instruction Access Units (IAU), 464
 - Instruction Level Parallelism (ILP), 101
 - Instruction Register (IR), 43, 163
 - Instruction Set Architecture (ISA), 43, 82
 - Instruction-level parallelism (ILP), 45
 - Instructions per Clock Cycle (IPC), 44
 - Integral control, 335
 - Integrated circuits (IC), 92
 - Intercloud, 32, 33, 345
 - exchange, 32
 - root, 32
 - nodes, 32
 - security, 291
 - Interconnection network, 47, 70, 85, 176, 189, 192, 194, 222, 404, 432, 464, 482
 - for computer clouds, 189
 - Interdependence, 267
 - Interface, 81
 - external, 438
 - network, 24, 178, 179, 189, 478, 480
 - Interface Definition Language (IDL), 414
 - Interface Message Processors (IMP), 93
 - Internal clock, 375
 - International Data Corporation (IDC), 7
 - Internet, 93, 175–178
 - transport protocols, 179
 - Internet Assigned Numbers Authority (IANA), 182
 - Internet Cache Protocol (ICP), 209
 - Internet exchange point (IXP), 181
 - Internet protocol (IP), 8
 - Internet Research Task Force (IRTF), 186
 - Internet Service Providers (ISP), 181
 - Interrupt virtualization, 143
 - Interval timer, 375
 - Intrusion detection, 268
 - systems, 575
 - Invariant behavior, 74
 - IP address, 178
 - elastic, 19, 20
 - private, 19
 - public, 19
 - IP anycasting, 210
 - IPsec, 183
 - Iso-latency policy, 116
 - Isolation
 - cache, 141
 - modules, 106
 - performance, 137
 - security, 37, 135, 137, 412
 - software fault, 450
 - Itanium
 - architecture, 142, 161, 173
 - paravirtualization, 136
 - processor, 162, 163
 - Iterator method, 121
 - IVA register, 163
- J**
- Java, 405
 - threads, 405
 - JAVA Message Service (JMS), 450
 - Java Virtual Machine (JVM), 273, 405
 - Job Manager (JM), 105
 - instance, 447
 - Job registry, 447
 - Job tracker, 298, 431, 432
 - Join, 121, 122
 - point, 465
 - routing task, 416
 - Journal storage, 220, 395
 - Just-in-time infrastructure, 411
- K**
- Kernel-based Virtual Machine (KVM), 145
 - Key Management Interoperability Protocol (KMIP), 33
 - Key Management Service (KMS), 268, 290
 - Key-based routing (KBR), 84
 - Key-value, 240, 241
 - Killer application, 37
 - Kimberlay, 481
 - Knowledge base (KB), 267
 - Kubernetes, 130–132
 - KVM hosting, 145
 - KVM virtualization, 148
- L**
- L1D cache, 49
 - Labels, 131
 - Lagrange multiplier method, 330
 - Lambda service, 23
 - Language INtegrated Query (LINQ), 105
 - Large Hadron Collider (LHC), 85

- Large-scale systems, 84
 - Last level cache (LLC), 118
 - Late binding, 303
 - Latency, 42, 293, 479, 484
 - communication, 24, 87, 99, 353
 - constraints, 474
 - DRAM, 99
 - hiding, 42, 43, 461
 - memory, 48, 67, 101, 103, 294, 443
 - Latency Critical (LC), 116
 - Latency-critical workloads, 116, 117
 - Layer, 81
 - access, 200
 - aggregate, 200
 - application, 16, 80
 - core, 200
 - data link, 177
 - data storage, 16, 498
 - database, 16
 - lineage, 127
 - network, 80, 176, 177
 - persistence, 127
 - processor abstraction, 163
 - transport, 80, 176
 - user interface, 16
 - Layering, 79, 176
 - Leader, 398, 423
 - participant, 460
 - Learner, 398
 - Learning rate, 114
 - Least Recently Used (LRU), 228
 - LHC Computing Grid (LCG), 85
 - Lineage, 121
 - layer, 127
 - Links, 122
 - Live media, 209
 - Live migration, 104, 146
 - Livelock condition, 394
 - Livelock freedom, 407
 - Liveness, 367, 416
 - Load balancing, 13
 - Load-store architecture, 43
 - Local
 - file, 225
 - memory, 251, 303, 462
 - relevance, 211
 - snapshot, 236
 - Locality, 48, 299, 455
 - disk, 251
 - rack, 299
 - servers, 299
 - spatial, 49, 356
 - temporal, 49, 356
 - Location transparency, 73, 227
 - Location transparent communication, 191
 - Lock, 392, 396
 - advisory, 234
 - Chubby, 236
 - coarse-grained, 234
 - file, 236
 - fine-grained, 234
 - granularity, 230
 - mandatory, 234
 - table, 460
 - tokens, 230
 - Locking, 234
 - Log data centers, 483
 - Logical clock (LC), 376
 - Logical organization, 223, 224
 - Logical Plan (LP), 498
 - Long Range Computer Study Group (LRCSG), 92
 - Long Short-Term Memory (LSTM), 506
 - Loopback file system (LOFS), 278
- M**
- Machine Intelligence (MI), 501
 - Machine Learning (ML), 11, 25, 63, 296, 501, 503, 506
 - Machine Perception (MP), 501
 - Machine-based automation, 90
 - Main memory, 49
 - Main Stack Pointer (MSP), 56
 - Malicious insiders, 261
 - Malicious software, 169
 - Malware, 169
 - Management
 - QoS-aware cluster, 113
 - resource, 127
 - shared state cluster, 111
 - Mandatory
 - locks, 234
 - security, 272
 - Map, 121, 122
 - Mapping a computation, 329
 - MapReduce, 105, 425, 426, 436, 449, 451
 - application, 113, 123, 133, 320, 433, 435, 452
 - scheduling, 320
 - framework, 433
 - Marked graph, 365, 368
 - Marking, 363
 - Markov decision processes, 488
 - Master, 232, 235, 236, 240, 425
 - Hadoop, 298
 - lease, 235
 - Master-slave replication, 240
 - Max-min fairness criterion, 297
 - Mechanism, 294

- flow-control, 179
 - parallel recovery, 467
 - Media Access Control (MAC), 80
 - Megastore, 245, 255, 438
 - Meltdown vulnerability, 54
 - Memcaching, 242
 - Memory
 - access, 54, 61, 138, 360, 361, 455
 - bandwidth, 48, 60
 - banks, 59
 - cost, 48
 - latency, 48, 67, 101, 103, 294, 443
 - local, 251, 303, 462
 - virtual, 51, 136, 155, 156, 171, 228, 250, 388, 409
 - virtualization, 104, 164, 172
 - Memory Management Unit (MMU), 57, 104, 161
 - Mesa, 106, 457, 458
 - Message, 356, 425
 - Message delivery
 - causal, 378
 - rules, 377
 - Message Passing Interface (MPI), 388
 - Metanode, 230
 - Metaphysical addressing, 164
 - Method
 - bootstrap, 492
 - iterator, 121
 - Lagrange multiplier, 330
 - Microservices, 18
 - Microsoft Azur Container Service, 130
 - Microsoft Windows Azure, 28
 - Middleware, 73, 85
 - Migration transparency, 74
 - MIMD, 47
 - Minimal witness basis, 253
 - Minimum charging units, 311
 - Minimum share timeout, 302
 - Minions, 130, 131
 - Miss penalty, 51
 - Mitigating cloud vulnerabilities, 287
 - Mobile Application Management (MAM), 287
 - Mobile Device Management (MDM), 287
 - Mode
 - overlay, 530
 - sensitive, 140
 - Model
 - arbitrarily divisible load sharing, 412
 - assumptions, 534
 - balls-and-bins, 400
 - bulk synchronous parallel, 360
 - computational, 358
 - decision, 33
 - enactment, 416
 - floating-point performance, 60
 - for multicore computing, 361
 - navigation, 254
 - NoSQL distribution, 240
 - roofline, 60
 - storage, 220
 - three-tier, 4
 - weak coordination, 421
 - Modularity, 74
 - Modularization, 176
 - Modularly divisible application, 329
 - Monitor, 334
 - Monitoring phase, 430
 - Multi-Layer Perceptrons (MLP), 506
 - Multicasting, 183
 - Multicloud, 39
 - Multicore processor speedup, 70
 - Multilevel nested virtualization, 156
 - Multiplexing, 136
 - Multistage interconnection networks, 193
 - Multitenancy, 259
 - Multithreaded SIMD processors, 61, 65
 - Mutual exclusion, 407
 - Myrinet, 194, 196
- N**
- N-way set associative cache, 50
 - Name node, 431, 432
 - Name Server (NS), 105
 - Named Data Networks (NDN), 186
 - National Science Foundation (NSF), 38
 - Native hypervisors, 170
 - Natural language processing (NLP), 507
 - Natural language understanding (NLU), 507
 - Navigational model, 254
 - Negotiation
 - object, 33
 - protocols, 33
 - Nested virtualization, 136, 142, 156, 158, 159, 172
 - multilevel, 156
 - single-level, 156
 - Network
 - adapter, 154
 - architecture, 154, 176
 - Clos, 193
 - congestion, 183, 204
 - content service, 211
 - content-centric, 185, 210
 - convolutional neural, 506
 - core, 176, 209
 - cost, 194
 - deep neural, 506
 - diameter, 189

- distribution, 185
- edge, 176
- fat-tree, 193
- flattened butterfly, 193
- infrastructure, 82, 181
- interconnection, 47, 70, 85, 176, 189
- Interconnection, 189
- interconnection, 189, 192, 194, 222, 404, 432, 464, 482
 - cloud, 191, 503
- interface, 24, 178, 179, 189, 478, 480
- layer, 80, 176, 177
- model, 238, 344
- multistage interconnection, 193
- overlay, 84, 188, 211
- packet-switched, 176
- resource management algorithms, 204
- software-defined, 187
- storage area, 197
- tensor, 553
- Tier 1, 181
- Tier 2, 181
- Tier 3, 181
- topology, 189, 231
- vehicular ad hoc, 211
- virtual, 84, 212, 480, 481
- virtualization, 154, 156
- Network Address Translation (NAT), 19, 182
- Network Element Control Protocol (NECP), 209
- Network File System (NFS), 76, 88, 216, 221, 223
- Network Interface Card (NIC), 153
- Network-centric computing, 10
- Network-centric content, 10
- Networking
 - AWS, 23
 - infrastructure, 94, 191
- Networks-on-a-chip (NoC), 65
- Neural networks (NN), 504
- Neutrality, 414
- Next Generation Enterprise Network (NGEN), 86
- Node, 189, 320
 - head, 316
 - name, 431, 432
 - worker, 316
- Non Secure World (NSW), 160
- Nonpreemption, 394
- Nonrepudiation, 32
- NoSQL databases, 222, 239–241, 243, 245, 247, 255, 269, 454
- NoSQL distribution models, 240
- O**
- Oak Ridge National Laboratory (ORNL), 72
- Object code, 82
- Object of negotiation, 33
- Object Request Broker (ORB), 414
- Object-relational impedance mismatch, 239
- Objective function, 322
- Official images, 129
- Omission failure, 398
- On-chip L1I, 49
- On-chip L2 cache, 49
- ONF architecture, 188
- Ontology, 32
- Opcode, 59
- Open Container Initiative (OCI), 130
- Open Network Foundation (ONF), 188
- Open Shortest Path First (OSPF), 187
- Open-Nebula, 168
- OpenStack, 168
- Operating system (OS), 57, 96, 135
 - security, 272
- Operation
 - gather, 59
 - scatter, 59
- Optimal control, 330, 332, 346, 490
- Optimal Partitioning Rule (OPR), 316
- Optimal resource management, 329
- Oracle, 515
- Order Preserving Encryption (OPE), 270
- Organization
 - kubernetes, 130
 - logical, 223, 224
 - physical, 223
 - scale-free, 344
- Over-provisioning, 325
- Overclocking, 118
- Overlay, 209
 - links, 84
 - networks, 84, 188, 211
- Overlimit, 207
- Overprovisioning, 6
- Oversubscription, 98, 191
- P**
- P2P replication, 240
- PaaS, 15, 16, 26, 34, 450
- Packet filtering, 575
- Packet-switched networks, 176
- Packets, 176, 425
- Pages, 51
- Paradigm
 - client-server, 41, 75–77, 225, 413
 - disruptive computing, 2
- Parallel
 - collection, 406
 - execution time, 68
 - file systems, 216, 222, 228

- hardware, 68
- recovery mechanism, 467
- slackness, 360
- Parallel File Systems (PFS), 221
- Parallel random-access machine (PRAM), 358
- Parallelism, 45
 - bit-level, 45
 - coarse-grained, 354
 - data-level, 67
 - fine-grained, 66, 354, 355
 - instruction-level, 45
 - task-level, 67
 - Thread-level, 67
- Paramount complication, 48
- Paravirtualization, 135, 140, 141, 148, 161, 164, 171, 274
 - Itanium, 136
 - strategy, 153
 - Xen, 149
- Partial Differential Equations (PDE), 66, 369
- Participant leader, 460
- Pass-through operators, 498
- Passphrase, 280
- Path, 224
 - trusted, 272
- Pattern, 151
 - workflow, 419
- Paxos, 397
 - algorithm, 236, 237
 - basic, 398
 - group, 460
- Peer-to-peer systems, 82
- Peering, 181
- Pending Interest Table (PIT), 187
- Perceptron, 504
 - rule, 505
- Perfect random hashing, 404
- Performance
 - metrics, 47, 113, 210, 338
 - transparency, 74
 - unpredictability, 137
 - VM, 153
 - WSC, 100
- Periodic task, 315
- Persistence layer, 127
- Persistent backlog, 207
- Persistent trust, 267
- Petri nets (PN), 363
- Phase
 - cleanup, 430
 - monitoring, 430
 - shutdown, 430
 - transition, 343
- Photon, 474
- Physical organization, 223
- Physical Plan (PP), 498
- Picasa, 27
- Pig, 436
- Pinhole, 575
- Pipeline
 - hazards, 46
 - scheduling, 46
 - stall, 46
- Pipelining, 45
- Plane
 - forwarding, 178
 - routing, 179
- Planning, 416
- Platform as a Service (PaaS), 15
- Pod, 201
- Point-of-presence (POP), 181
- Policy, 268, 294
- Port, 179, 180
- Precision, 375
- Precursors, 223
- Predecessor task, 415
- Prediction
 - branch, 53
 - history-based, 53
- Presentation tier, 4
- Preserve exception behavior, 48
- Preserve instruction flow, 48
- Preserve security, 48
- Pricing, 322
- Primitive processes, 359
- Primitive task, 415
- Priority bands, 109
- Priority inheritance, 127
- Priority inversion, 394
- Priority-based scheduling, 127
- Privacy
 - impact, 291
 - risks, 266, 278, 280
- Privacy Impact Assessment (PIA), 264, 266
- Private
 - cloud, 4, 5, 167, 168
 - IP address, 19
- Privilege leaking, 162
- Privilege rings, 162
- Privileged instructions, 81, 162
- Privileged-sensitive instructions, 163
- Probability density function (PDF), 484
- Problem
 - Deutsch–Jozsa, 515
 - embarrassingly parallel, 67
- Process, 82, 359, 369
 - balloon, 104

- description, 416
 Markov decision, 488
 primitive, 359
 VM, 138
- Process Stack Pointer (PSP), 56
- Processor
 abstraction layer, 163
 asymmetric core, 70
 bandwidth, 43, 307, 356
 clock, 44
 emulation, 137
 Itanium, 162, 163
 management, 81
 multithreaded SIMD, 61, 65
 resources, 117
 sharing, 96
 superscalar, 45, 47
 symmetric core, 70
 utilization, 96, 132, 336, 353, 361
 virtualization, 80, 96, 135, 137, 140, 172, 388, 409
- Producer-consumer synchronization, 407
- Program counter (PC), 43
- Proportional thresholding, 335
- Proposal, 398, 425
- Proposer, 398
- Protocol, 176, 209, 371
 connection-oriented, 183
 connectionless, 183
 negotiation, 33
- Provider
 CDN, 209
 service, 14
- Public cloud, 4
- Public IP address, 19
- Q**
- QEMU, 144
- QoS, 113
 crosstalk, 137
- QoS-aware cluster management, 113
- Quad Data Rate (QDR), 195, 445
- Quantum communication, 509
- Quantum Computing (QC), 501, 509
 Playground, 516
- Quantum Information Processing (QIP), 509
- Quantum Inspire, 518
- Quantum parallelism, 510
- Quasar, 114
- Qubit, 509
 control, 511
- Quorum, 398
- Quota system, 109
- R**
- Rack, 98
 locality, 299
- Random-access machine (RAM), 358
- Rank, 121, 122
- Rate Monotonic Algorithms (RMA), 298
- Reachability, 367
- Read after Write (RAW), 46
- Real-time bandwidth, 118
- Recurrent Neural Networks (RNN), 506
- Recursive construction, 75
- Red Hat Virtualization (RHV), 170
- Regions, 23
- Register
 CPU, 49
 data flow, 52
 stack engine, 163
 vector length, 59
 vector mask, 59
- Register File (RF), 44
- Registration Authority (RA), 32
- Regularization factor, 114
- Relational algebra, 238
- Relational Database Service (RDS), 10
- Relational trust, 267
- Relative deadline, 315
- Remote File System (RFS), 223, 255
- Remote Procedure Call (RPC), 73, 76, 225
- Replanning, 415
- Replay log, 236
- Replication
 controllers, 131
 master-slave, 240
 P2P, 240
 transparency, 73
- Representational State Transfer (REST), 414
- Reputation, 268
- Request for quotation (RFQ), 93
- Request-parallel applications, 99
- Reservation stations (RS), 52
- Resilient Distributed Dataset (RDD), 121, 466
- Resource
 allocation, 113
 assignment, 113
 isolation, 116
 management, 127, 328
 multiplexing, 2
 processor, 117
 scale out, 115
 scale up, 115
 virtualization, 136
- Resource Description Framework (RDF), 85
- Response time, 77, 78, 165

- Ring, 151
 - aliasing, 143
 - compression, 143, 162
 - deprivileging, 142
 - privilege, 162
 - Risk, 267
 - privacy, 266, 278, 280
 - profile, 261
 - security, 259, 271, 279–281, 287, 291
 - Rogue hypervisor, 169
 - Roofline model, 60
 - Root, 97, 576
 - Rootkit, 169
 - Routing
 - DNS-based, 210
 - plane, 179
 - task, 415
 - Row
 - cache, 50
 - groups, 439
 - Rule, 253, 576
 - Run, 381
- S**
- Safety, 416
 - Same Program Multiple Data (SPMD), 67, 425
 - Sample estimate result, 495
 - Sampling, 121
 - Sampling error, 495
 - distribution, 495
 - Scalability, 86, 239, 240
 - Scalable data center communication architectures, 200
 - Scale-free organization, 344
 - Scaled speedup, 69
 - Scaling
 - horizontal, 242, 329
 - incommensurate, 343
 - transparency, 74
 - vertical, 328
 - Scan consolidation, 498
 - Scatter operation, 59
 - Scheduler
 - Apache capacity, 306
 - fair-share, 127
 - Hadoop, 298, 302
 - Spark, 121, 122, 306
 - thread, 63, 311
 - workload, 111
 - Scheduler virtual time (SVT), 311
 - Scheduling, 316
 - data-aware, 303
 - delay, 298
 - Schema, 454
 - Schema-on-write, 454
 - Scoring, 109
 - Search time, 217
 - Searchable symmetric encryption (SSE), 270
 - Secondary spectrum-data falsification (SSDF), 530
 - Secure arguments, 75
 - Secure Sockets Layer (SSL), 183, 262
 - Secure Virtual Machine (SVM), 142
 - Security, 258
 - assessment, 264
 - audit, 15
 - AWS, 258, 289, 290
 - cloud, 15, 30, 36, 258, 259, 263, 286, 289–291, 452
 - Intercloud, 291
 - isolation, 37, 135, 137, 412
 - mandatory, 272
 - of database services, 270
 - operating system, 272
 - risks, 259, 271, 278–281, 287, 291
 - virtual machine, 273
 - virtualization, 257, 275
 - VM, 257, 273
 - Security as a service (SecaaS), 264
 - Security Extensions (SE), 161
 - Seek time, 217
 - Segments, 179
 - Selector, 206
 - Self-organization, 342, 344
 - Semantic integration, 28
 - Semantic Web, 85
 - Sensitive instructions, 81, 173
 - Sensitive machine instructions, 140
 - Sensors, 329, 333, 334
 - Sequencer, 234
 - Sequential execution time, 68
 - Serialization, 392
 - Server
 - chunk, 232, 233
 - HTTP, 77
 - locality, 299
 - stateless, 413
 - thrift, 438
 - utilization, 127, 210, 227, 259, 304, 325
 - virtual, 18, 19, 488
 - WSC, 99
 - Service, 131
 - composite, 450
 - consensus, 397
 - consumer, 14
 - hijacking, 261
 - Lambda, 23
 - provider, 14
 - Service Level Agreement (SLA), 6, 13, 33, 36, 260

- Service Level Objectives (SLO), 116
- Servless computer service, 23
- Sessions, 470
- Set-associative cache, 50
- Sharding, 240
- Shared cluster state, 111
- Shared memory system (SMS), 464
- Shared security responsibility (SSR), 17, 289
- Shared state cluster management, 111
- Shared technology, 261
- Shared write lock, 230
- Shutdown phase, 430
- Sigmoid neuron (SN), 505
- SIMD, 47
 - architectures, 58–60
- Simple DB, 21
- Simple Object Access Protocol (SOAP), 414
- Simple Queue Service (SQS), 21
- Simple Storage System (S3), 20
- Simple Workflow Service (SWS), 22
- Simultaneous multithreading (SMT), 103
- Single data rate (SDR), 195
- Single-level nested virtualization, 156
- Singular Value Decomposition (SVD), 113
- SISD processors, 46
- Slack, 27, 119, 329
- Slave threads, 465
- Sliding, 470
- Slots, 106, 320
 - sticky, 299
- Socket, 88, 179
- Soft modularity, 75
- Soft-state, 241
- Software, 87
 - architecture, 4
 - malicious, 169
 - virtualization, 136, 170
 - vulnerability, 278
- Software as a Service (SaaS), 15
- Software fault isolation (SFI), 450
- Software-defined Network (SDN), 187
- Software-Defined Wide Area Network (SD-WAN), 188
- Solid-state disks, 217
- Space
 - active, 421
 - tuple, 421
- Space–time diagram, 377, 381
- Spanner, 459, 460
- Spark, 121, 122
 - dependencies, 122
 - driver programs, 122
 - scheduler, 121, 122, 306
 - streaming, 466, 467, 470, 498
- Spatial locality, 49, 356
- Spectre, 54
- Speculation
 - branch condition, 53
 - branch target, 53
- Speedup, 68, 353
- Spintronics, 216
- Split, 427
 - drivers, 151
- Sprite File System (SFS), 216, 227, 255
- SRAM, 49
- Stack pointer, 387
- Start-time fair queuing (SFQ), 307, 347
- State, 369
 - feasible, 332
- State Machine, 368
- State of a case, 416
- Stateless address auto-configuration (SLAAC), 183
- Stateless servers, 413
- Static content, 209
- Sticky slots, 299
- Stochastic Fairness Queuing (SFQ), 206
- Stochastic Gradient Descent (SGD), 114
- Stochastic Petri Nets (SPN), 369
- Storage
 - cell, 220, 395, 396
 - cloud, 26, 216, 255, 257, 289
 - journal, 220, 395
- Storage area network (SAN), 197, 221, 255
- Storage manager, 396
- Storage models, 220
- Straggler, 303
- Strategy
 - biased, 53
 - cooperative spectrum-sensing, 530
- Streaming SIMD extensions (SSE), 59
- Stride, 59
- Strong clock condition, 380
- Structural hazards, 46
- Structured Query Language (SQL), 238
- Stubs, 414
- Successor task, 415
- Supernodes, 84
- Superscalar processor, 45, 47
- Supersteps, 360
- Symmetric confusion, 364
- Symmetric core processor, 70
- Synchronization, 350
 - barrier, 352, 354, 355
 - producer–consumer, 407
- System
 - distributed, 73
 - emulation, 145

- energy-proportional, 326
 - intrusion detection, 575
 - large-scale, 84
 - management
 - cluster, 110
 - resource, 75
 - model, 316
 - peer-to-peer, 82
 - quota, 109
 - VM, 138
- T**
- Table schema, 458
 - Tablet, 459
 - Tachyon, 124, 127
 - Tachyon Master (TM), 125
 - Tag, 50
 - Target qubit, 511
 - Task, 360, 415, 416
 - aperiodic, 316
 - composite, 415
 - periodic, 315
 - predecessor, 415
 - primitive, 415
 - routing, 415
 - fork, 416
 - join, 416
 - successor, 415
 - tracker, 298, 431, 432
 - Task-level parallelism, 67
 - Task-priority register (TPR), 143
 - TCP congestion control, 183
 - Temporal locality, 49, 356
 - Tensor, 62
 - contraction, 552, 553
 - network, 553
 - Tensor Processing Units (TPU), 2, 506, 507
 - Terraform, 18
 - Tez, 434
 - Thin clients, 11, 41
 - Third Party Auditors (TPA), 271
 - Third-party control, 261
 - Thread, 67
 - allocate, 385
 - blocks scheduler, 61
 - groups, 354, 369
 - pool, 405
 - scheduler, 63, 311
 - block, 63
 - slave, 465
 - Thread-level parallelism, 67
 - Three-tier model, 4
 - Three-way handshake, 77
 - Thresholds, 335
 - Thrift Server, 438
 - Thunk, 145
 - Tier
 - application/logic, 4
 - data, 4
 - presentation, 4
 - Tier 1 network, 181
 - Tier 2 network, 181
 - Tier 3 network, 181
 - Time
 - intervals, 375
 - response, 77, 78, 165
 - search, 217
 - seek, 217
 - skew, 470
 - stamp, 375, 379, 468
 - Time-slice, 298
 - Timers, 468
 - Timing, 204
 - Token manager, 230
 - Tokens, 366
 - byte-range, 230
 - lock, 230
 - Total surplus, 322
 - TPUs, 2, 42, 63, 506, 507
 - Tracker
 - job, 298, 431, 432
 - task, 298, 431, 432
 - Traditional threats, 260
 - Training data, 504
 - Transaction manager, 460
 - Transit, 181
 - centers, 23
 - Transition, 363
 - phase, 343
 - Transparency
 - access, 73
 - concurrency, 73
 - failure, 73
 - location, 73, 227
 - migration, 74
 - performance, 74
 - replication, 73
 - scaling, 74
 - Transport layer, 80, 176
 - Transport Layer Security (TLS), 183, 282, 287
 - Tread scheduler, 61
 - Triggers, 22, 471
 - Trust, 267, 530
 - calculus-based, 267
 - dynamic, 267
 - indexes, 32

- persistent, 267
- relational, 267
- Trust Zone Secure World (TZSW), 161
- Trusted applications, 272
- Trusted Computer Base (TCB), 38, 281
- Trusted Execution Environment (TEE), 160
- Trusted paths, 272
- Tunneling, 200
- Tuple calculus, 255
- Tuple space, 421
- Two-level resource allocation strategy, 333

- U**
- Undelete files, 278
- Underlimit, 207
- Union, 121
- Universal gates, 505
- Universal Serial Bus (USB), 153
- UNIX File System (UFS), 223, 236
- User
 - challenges, 34
 - images, 129
- User-mode emulation, 145
- Utility computing, 1
- Utility function, 338
- Utility-based, 296

- V**
- Validation Authority (VA), 32
- Values, 457
- Variability, 502
- Vast amounts, 476
- Vector
 - architectures, 59
 - clocks, 249
 - length registers, 59
 - mask registers, 59
- Vehicular
 - clouds, 518, 519
 - networks, 519
- Vehicular ad hoc network (VANET), 211
- Vehicular Virtual Machine Migration (VVMM), 522
- Vendor lock-in, 9, 30
- Vertical scaling, 242, 328
- Very Long Instruction Word (VLIW), 161
- Virtual
 - address, 51
 - cloud, 316
 - hardware, 97, 137, 169
 - memory, 51, 136, 155, 156, 171, 228, 250, 388, 409
 - network, 84, 212, 480, 481
 - server, 19
 - servers, 18, 488
 - time, 311
- Virtual Block Device (VBD), 151
- Virtual circuits (VC), 199
- Virtual Machine Control Structure (VMCS), 142
- Virtual machine (VM), 18, 20, 95, 135, 138, 164, 213
- Virtual Memory Manager (VMM), 388
- Virtual Network Computing (VNC), 481
- Virtual Network Interfaces (VNI), 151
- Virtual Private Cloud (VPC), 22
- Virtual Private Network (VPN), 22, 26
- Virtual-Machine Based Rootkit (VMBR), 169, 170
- Virtualization
 - by aggregation, 95
 - by multiplexing, 95
 - CPU, 163
 - full, 135, 140
 - hardware, 146, 171
 - support, 141
 - interrupt, 143
 - KVM, 148
 - memory, 104, 164, 172
 - nested, 136, 142, 156, 158, 159, 172
 - multilevel, 156
 - single-level, 156
 - network, 154, 156
 - processor, 80, 96, 135, 137, 140, 172, 388, 409
 - resource, 136
 - security, 257, 275
 - software, 136, 170
- VM
 - application, 138
 - import, 20, 547
 - performance, 153
 - process, 138
 - security, 257, 273
 - system, 138
- Volunteer computing, 449
- von Neumann architecture, 91
- Vulnerability
 - critical, 54
 - software, 278

- W**
- Wait time, 302
- Warehouse Scale Computers (WSC), 95, 96, 98
- Watch function, 151
- Wave-particle duality principle, 508
- Weak coordination models, 421
- Web Application Firewall (WAF), 290
- Web browser, 77
- Web Cache Coordination Protocol (WCCP), 209
- Web page, 77

- Web Services Description Language (WSDL), 414
Wide Area Networks (WAN), 73
Wimpy cores, 99
Windowing, 470
Windows Azure, 28
Wireless local area networks (WLAN), 478
Witness tree, 404
Worker node, 316
Worker role, 29
Workflow, 415
 dynamic, 421
 pattern, 419
Working set, 49, 120, 455
Workload
 derivative, 316
 latency-critical, 117
 scheduler, 111
Write after Read (WAR), 46
Write after Write (WAW), 46
Write-ahead, 229
Write-write conflict, 240
WSC
 performance, 100
 servers, 99
- X**
- XDP, 77
Xen, 148
 hosting, 145
 hypervisor, 19, 138, 155, 164, 172, 281
 paravirtualization, 149
XenStore, 151, 282, 284, 285, 291
Xoar, 283
- Y**
- Yarn, 433
- Z**
- Znode stores, 425
Zonemaster, 459