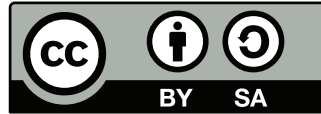


# Operating Systems and Middleware: Supporting Controlled Interaction

Max Hailperin  
Gustavus Adolphus College

Revised Edition 1.1  
July 27, 2011

Copyright © 2011 by Max Hailperin.



This work is licensed under the Creative Commons Attribution-ShareAlike 3.0 Unported License. To view a copy of this license, visit

*<http://creativecommons.org/licenses/by-sa/3.0/>*

or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.

To my family



# Contents

<b>Preface</b>	<b>xi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Chapter Overview . . . . .	1
1.2 What Is an Operating System? . . . . .	2
1.3 What is Middleware? . . . . .	6
1.4 Objectives for the Book . . . . .	8
1.5 Multiple Computations on One Computer . . . . .	9
1.6 Controlling the Interactions Between Computations . . . . .	11
1.7 Supporting Interaction Across Time . . . . .	13
1.8 Supporting Interaction Across Space . . . . .	15
1.9 Security . . . . .	17
<b>2 Threads</b>	<b>21</b>
2.1 Introduction . . . . .	21
2.2 Example of Multithreaded Programs . . . . .	23
2.3 Reasons for Using Concurrent Threads . . . . .	27
2.4 Switching Between Threads . . . . .	30
2.5 Preemptive Multitasking . . . . .	37
2.6 Security and Threads . . . . .	38
<b>3 Scheduling</b>	<b>45</b>
3.1 Introduction . . . . .	45
3.2 Thread States . . . . .	46
3.3 Scheduling Goals . . . . .	49
3.3.1 Throughput . . . . .	51
3.3.2 Response Time . . . . .	54
3.3.3 Urgency, Importance, and Resource Allocation . . . . .	55
3.4 Fixed-Priority Scheduling . . . . .	61

3.5	Dynamic-Priority Scheduling . . . . .	65
3.5.1	Earliest Deadline First Scheduling . . . . .	65
3.5.2	Decay Usage Scheduling . . . . .	66
3.6	Proportional-Share Scheduling . . . . .	71
3.7	Security and Scheduling . . . . .	79
<b>4</b>	<b>Synchronization and Deadlocks</b>	<b>93</b>
4.1	Introduction . . . . .	93
4.2	Races and the Need for Mutual Exclusion . . . . .	95
4.3	Mutexes and Monitors . . . . .	98
4.3.1	The Mutex Application Programming Interface . . . . .	99
4.3.2	Monitors: A More Structured Interface to Mutexes . . . . .	103
4.3.3	Underlying Mechanisms for Mutexes . . . . .	106
4.4	Other Synchronization Patterns . . . . .	110
4.4.1	Bounded Buffers . . . . .	113
4.4.2	Readers/Writers Locks . . . . .	115
4.4.3	Barriers . . . . .	116
4.5	Condition Variables . . . . .	117
4.6	Semaphores . . . . .	123
4.7	Deadlock . . . . .	124
4.7.1	The Deadlock Problem . . . . .	126
4.7.2	Deadlock Prevention Through Resource Ordering . . . . .	128
4.7.3	Ex Post Facto Deadlock Detection . . . . .	129
4.7.4	Immediate Deadlock Detection . . . . .	132
4.8	The Interaction of Synchronization with Scheduling . . . . .	134
4.8.1	Priority Inversion . . . . .	135
4.8.2	The Convoy Phenomenon . . . . .	137
4.9	Nonblocking Synchronization . . . . .	141
4.10	Security and Synchronization . . . . .	145
<b>5</b>	<b>Atomic Transactions</b>	<b>159</b>
5.1	Introduction . . . . .	159
5.2	Example Applications of Transactions . . . . .	162
5.2.1	Database Systems . . . . .	163
5.2.2	Message-Queuing Systems . . . . .	167
5.2.3	Journalled File Systems . . . . .	172
5.3	Mechanisms to Ensure Atomicity . . . . .	174
5.3.1	Serializability: Two-Phase Locking . . . . .	174
5.3.2	Failure Atomicity: Undo Logging . . . . .	183
5.4	Transaction Durability: Write-Ahead Logging . . . . .	186

5.5	Additional Transaction Mechanisms . . . . .	190
5.5.1	Increased Transaction Concurrency: Reduced Isolation	191
5.5.2	Coordinated Transaction Participants: Two-Phase Com- mit . . . . .	193
5.6	Security and Transactions . . . . .	196
<b>6</b>	<b>Virtual Memory</b>	<b>207</b>
6.1	Introduction . . . . .	207
6.2	Uses for Virtual Memory . . . . .	212
6.2.1	Private Storage . . . . .	212
6.2.2	Controlled Sharing . . . . .	213
6.2.3	Flexible Memory Allocation . . . . .	216
6.2.4	Sparse Address Spaces . . . . .	219
6.2.5	Persistence . . . . .	219
6.2.6	Demand-Driven Program Loading . . . . .	220
6.2.7	Efficient Zero Filling . . . . .	221
6.2.8	Substituting Disk Storage for RAM . . . . .	222
6.3	Mechanisms for Virtual Memory . . . . .	223
6.3.1	Software/Hardware Interface . . . . .	225
6.3.2	Linear Page Tables . . . . .	229
6.3.3	Multilevel Page Tables . . . . .	234
6.3.4	Hashed Page Tables . . . . .	239
6.3.5	Segmentation . . . . .	242
6.4	Policies for Virtual Memory . . . . .	247
6.4.1	Fetch Policy . . . . .	248
6.4.2	Placement Policy . . . . .	250
6.4.3	Replacement Policy . . . . .	252
6.5	Security and Virtual Memory . . . . .	259
<b>7</b>	<b>Processes and Protection</b>	<b>269</b>
7.1	Introduction . . . . .	269
7.2	POSIX Process Management API . . . . .	271
7.3	Protecting Memory . . . . .	281
7.3.1	The Foundation of Protection: Two Processor Modes	282
7.3.2	The Mainstream: Multiple Address Space Systems . .	285
7.3.3	An Alternative: Single Address Space Systems . . . .	287
7.4	Representing Access Rights . . . . .	289
7.4.1	Fundamentals of Access Rights . . . . .	289
7.4.2	Capabilities . . . . .	295
7.4.3	Access Control Lists and Credentials . . . . .	299

7.5	Alternative Granularities of Protection . . . . .	307
7.5.1	Protection Within a Process . . . . .	308
7.5.2	Protection of Entire Simulated Machines . . . . .	309
7.6	Security and Protection . . . . .	313
<b>8</b>	<b>Files and Other Persistent Storage</b>	<b>329</b>
8.1	Introduction . . . . .	329
8.2	Disk Storage Technology . . . . .	332
8.3	POSIX File API . . . . .	336
8.3.1	File Descriptors . . . . .	336
8.3.2	Mapping Files Into Virtual Memory . . . . .	341
8.3.3	Reading and Writing Files at Specified Positions . . . . .	344
8.3.4	Sequential Reading and Writing . . . . .	344
8.4	Disk Space Allocation . . . . .	346
8.4.1	Fragmentation . . . . .	347
8.4.2	Locality . . . . .	350
8.4.3	Allocation Policies and Mechanisms . . . . .	352
8.5	Metadata . . . . .	354
8.5.1	Data Location Metadata . . . . .	355
8.5.2	Access Control Metadata . . . . .	364
8.5.3	Other Metadata . . . . .	367
8.6	Directories and Indexing . . . . .	367
8.6.1	File Directories Versus Database Indexes . . . . .	367
8.6.2	Using Indexes to Locate Files . . . . .	369
8.6.3	File Linking . . . . .	370
8.6.4	Directory and Index Data Structures . . . . .	374
8.7	Metadata Integrity . . . . .	375
8.8	Polymorphism in File System Implementations . . . . .	379
8.9	Security and Persistent Storage . . . . .	380
<b>9</b>	<b>Networking</b>	<b>391</b>
9.1	Introduction . . . . .	391
9.1.1	Networks and Internets . . . . .	392
9.1.2	Protocol Layers . . . . .	394
9.1.3	The End-to-End Principle . . . . .	397
9.1.4	The Networking Roles of Operating Systems, Middle- ware, and Application Software . . . . .	398
9.2	The Application Layer . . . . .	399
9.2.1	The Web as a Typical Example . . . . .	399



9.2.2	The Domain Name System: Application Layer as Infrastructure . . . . .	402
9.2.3	Distributed File Systems: An Application Viewed Through Operating Systems . . . . .	405
9.3	The Transport Layer . . . . .	407
9.3.1	Socket APIs . . . . .	408
9.3.2	TCP, the Dominant Transport Protocol . . . . .	414
9.3.3	Evolution Within and Beyond TCP . . . . .	417
9.4	The Network Layer . . . . .	418
9.4.1	IP, Versions 4 and 6 . . . . .	418
9.4.2	Routing and Label Switching . . . . .	421
9.4.3	Network Address Translation: An End to End-to-End? . . . . .	422
9.5	The Link and Physical Layers . . . . .	425
9.6	Network Security . . . . .	427
9.6.1	Security and the Protocol Layers . . . . .	428
9.6.2	Firewalls and Intrusion Detection Systems . . . . .	430
9.6.3	Cryptography . . . . .	431
<b>10</b>	<b>Messaging, RPC, and Web Services</b>	<b>443</b>
10.1	Introduction . . . . .	443
10.2	Messaging Systems . . . . .	444
10.3	Remote Procedure Call . . . . .	447
10.3.1	Principles of Operation for RPC . . . . .	448
10.3.2	An Example Using Java RMI . . . . .	451
10.4	Web Services . . . . .	455
10.5	Security and Communication Middleware . . . . .	463
<b>11</b>	<b>Security</b>	<b>473</b>
11.1	Introduction . . . . .	473
11.2	Security Objectives and Principles . . . . .	474
11.3	User Authentication . . . . .	480
11.3.1	Password Capture Using Spoofing and Phishing . . . . .	481
11.3.2	Checking Passwords Without Storing Them . . . . .	483
11.3.3	Passwords for Multiple, Independent Systems . . . . .	483
11.3.4	Two-Factor Authentication . . . . .	483
11.4	Access and Information-Flow Controls . . . . .	486
11.5	Viruses and Worms . . . . .	491
11.6	Security Assurance . . . . .	495
11.7	Security Monitoring . . . . .	497
11.8	Key Security Best Practices . . . . .	500

<b>A Stacks</b>	<b>511</b>
A.1 Stack-Allocated Storage: The Concept . . . . .	512
A.2 Representing a Stack in Memory . . . . .	513
A.3 Using a Stack for Procedure Activations . . . . .	514
<b>Bibliography</b>	<b>517</b>

# Preface

Suppose you sit down at your computer to check your email. One of the messages includes an attached document, which you are to edit. You click the attachment, and it opens up in another window. After you start editing the document, you realize you need to leave for a trip. You save the document in its partially edited state and shut down the computer to save energy while you are gone. Upon returning, you boot the computer back up, open the document, and continue editing.

This scenario illustrates that computations interact. In fact, it demonstrates at least three kinds of interactions between computations. In each case, one computation provides data to another. First, your email program retrieves new mail from the server, using the Internet to bridge space. Second, your email program provides the attachment to the word processor, using the operating system's services to couple the two application programs. Third, the invocation of the word processor that is running before your trip provides the partially edited document to the invocation running after your return, using disk storage to bridge time.

In this book, you will learn about all three kinds of interaction. In all three cases, interesting software techniques are needed in order to bring the computations into contact, yet keep them sufficiently at arms length that they don't compromise each other's reliability. The exciting challenge, then, is supporting controlled interaction. This includes support for computations that share a single computer and interact with one another, as your email and word processing programs do. It also includes support for data storage and network communication. This book describes how all these kinds of support are provided both by operating systems and by additional software layered on top of operating systems, which is known as middleware.

## Audience

If you are an upper-level computer science student who wants to understand how contemporary operating systems and middleware products work and why they work that way, this book is for you. In this book, you will find many forms of balance. The high-level application programmer's view, focused on the services that system software provides, is balanced with a lower-level perspective, focused on the mechanisms used to provide those services. Timeless concepts are balanced with concrete examples of how those concepts are embodied in a range of currently popular systems. Programming is balanced with other intellectual activities, such as the scientific measurement of system performance and the strategic consideration of system security in its human and business context. Even the programming languages used for examples are balanced, with some examples in Java and others in C or C++. (Only limited portions of these languages are used, however, so that the examples can serve as learning opportunities, not stumbling blocks.)

## Systems Used as Examples

Most of the examples throughout the book are drawn from the two dominant families of operating systems: Microsoft Windows and the UNIX family, including especially Linux and Mac OS X. Using this range of systems promotes the students' flexibility. It also allows a more comprehensive array of concepts to be concretely illustrated, as the systems embody fundamentally different approaches to some problems, such as the scheduling of processors' time and the tracking of files' disk space.

Most of the examples are drawn from the stable core portions of the operating systems and, as such, are equally applicable to a range of specific versions. Whenever Microsoft Windows is mentioned without further specification, the material should apply to Windows NT, Windows 2000, Windows XP, Windows Server 2003, Windows Vista, Windows 2008, and Windows 7. All Linux examples are from version 2.6, though much of the material applies to other versions as well. Wherever actual Linux source code is shown (or whenever fine details matter for other reasons), the specific subversion of 2.6 is mentioned in the end-of-chapter notes. Most of the Mac OS X examples originated with version 10.4, also known as Tiger, but should be applicable to other versions.

Where the book discusses the protection of each process's memory, one

additional operating system is brought into the mix of examples, in order to illustrate a more comprehensive range of alternative designs. The IBM iSeries, formerly known as the AS/400, embodies an interesting approach to protection that might see wider application within current students' lifetimes. Rather than giving each process its own address space (as Linux, Windows, and Mac OS X do), the iSeries allows all processes to share a single address space and to hold varying access permissions to individual objects within that space.

Several middleware systems are used for examples as well. The Oracle database system is used to illustrate deadlock detection and recovery as well as the use of atomic transactions. Messaging systems appear both as another application of atomic transactions and as an important form of communication middleware, supporting distributed applications. The specific messaging examples are drawn from the IBM WebSphere MQ system (formerly MQSeries) and the Java Message Service (JMS) interface, which is part of Java 2 Enterprise Edition (J2EE). The other communication middleware examples are Java RMI (Remote Method Invocation) and web services. Web services are explained in platform-neutral terms using the SOAP and WSDL standards, as well as through a J2EE interface, JAX-RPC (Java API for XML-Based RPC).

## Organization of the Text

Chapter 1 provides an overview of the text as a whole, explaining what an operating system is, what middleware is, and what sorts of support these systems provide for controlled interaction.

The next nine chapters work through the varieties of controlled interaction that are exemplified by the scenario at the beginning of the preface: interaction between concurrent computations on the same system (as between your email program and your word processor), interaction across time (as between your word processor before your trip and your word processor after your trip), and interaction across space (as between your email program and your service provider's email server).

The first of these three topics is controlled interaction between computations operating at one time on a particular computer. Before such interaction can make sense, you need to understand how it is that a single computer can be running more than one program, such as an email program in one window and a word processing program in another. Therefore, Chapter 2 explains the fundamental mechanism for dividing a computer's attention

between concurrent computations, known as threads. Chapter 3 continues with the related topic of scheduling. That is, if the computer is dividing its time between computations, it needs to decide which ones to work on at any moment.

With concurrent computations explained, Chapter 4 introduces controlled interactions between them by explaining synchronization, which is control over the threads' relative timing. For example, this chapter explains how, when your email program sends a document to your word processor, the word processor can be constrained to read the document only after the email program writes it. One particularly important form of synchronization, atomic transactions, is the topic of Chapter 5. Atomic transactions are groups of operations that take place as an indivisible unit; they are most commonly supported by middleware, though they are also playing an increasing role in operating systems.

Other than synchronization, the main way that operating systems control the interaction between computations is by controlling their access to memory. Chapter 6 explains how this is achieved using the technique known as virtual memory. That chapter also explains the many other objectives this same technique can serve. Virtual memory serves as the foundation for Chapter 7's topic, which is processes. A process is the fundamental unit of computation for protected access, just as a thread is the fundamental unit of computation for concurrency. A process is a group of threads that share a protection environment; in particular, they share the same access to virtual memory.

The next three chapters move outside the limitations of a single computer operating in a single session. First, consider the document stored before a trip and available again after it. Chapter 8 explains persistent storage mechanisms, focusing particularly on the file storage that operating systems provide. Second, consider the interaction between your email program and your service provider's email server. Chapter 9 provides an overview of networking, including the services that operating systems make available to programs such as the email client and server. Chapter 10 extends this discussion into the more sophisticated forms of support provided by communication middleware, such as messaging systems, RMI, and web services.

Finally, Chapter 11 focuses on security. Because security is a pervasive issue, the preceding ten chapters all provide some information on it as well. Specifically, the final section of each chapter points out ways in which security relates to that chapter's particular topic. However, even with that coverage distributed throughout the book, a chapter specifically on security

is needed, primarily to elevate it out of technical particulars and talk about general principles and the human and organizational context surrounding the computer technology.

The best way to use these chapters is in consecutive order. However, Chapter 5 can be omitted with only minor harm to Chapters 8 and 10, and Chapter 9 can be omitted if students are already sufficiently familiar with networking.

## Relationship to Computer Science Curriculum 2008

Operating systems are traditionally the subject of a course required for all computer science majors. In recent years, however, there has been increasing interest in the idea that upper-level courses should be centered less around particular artifacts, such as operating systems, and more around cross-cutting concepts. In particular, the *Computing Curricula 2001* (CC2001) and its interim revision, *Computer Science Curriculum 2008* (CS2008), provide encouragement for this approach, at least as one option. Most colleges and universities still retain a relatively traditional operating systems course, however. Therefore, this book steers a middle course, moving in the direction of the cross-cutting concerns while retaining enough familiarity to be broadly adoptable.

The following table indicates the placement within this text of knowledge units from CS2008's computer science body of knowledge. Those knowledge units designated as core units within CS2008 are listed in italics. The book covers all core operating systems (OS) units, as well as one elective OS unit. The overall amount of coverage for each unit is always at least that recommended by CS2008, though sometimes the specific subtopics don't quite correspond exactly. Outside the OS area, this book's most substantial coverage is of Net-Centric Computing (NC); another major topic, transaction processing, comes from Information Management (IM). In each row, the listed chapters contain the bulk of the knowledge unit's coverage, though

some topics may be elsewhere.

<b>Knowledge unit (italic indicates core units in CS2008)</b>	<b>Chapter(s)</b>
<i>OS/OverviewOfOperatingSystems</i>	1
<i>OS/OperatingSystemPrinciples</i>	1, 7
<i>OS/Concurrency</i>	2, 4
<i>OS/SchedulingAndDispatch</i>	3
<i>OS/MemoryManagement</i>	6
<i>OS/SecurityAndProtection</i>	7, 11
OS/FileSystems	8
NC/Introduction	9
NC/NetworkCommunication (partial coverage)	9
NC/NetworkSecurity (partial coverage)	9
NC/WebOrganization (partial coverage)	9
NC/NetworkedApplications (partial coverage)	10
IM/TransactionProcessing	5

## Your Feedback is Welcome

Comments, suggestions, and bug reports are welcome; please send email to [max@gustavus.edu](mailto:max@gustavus.edu). Bug reports in particular can earn you a bounty of \$2.56 apiece as a token of gratitude. (The great computer scientist Donald Knuth started this tradition. Given how close to bug-free his publications have become, it seems to work.) For purposes of this reward, the definition of a bug is simple: if as a result of your email the author chooses to make a change, then you have pointed out a bug. The change need not be the one you suggested, and the bug need not be technical in nature. Unclear writing qualifies, for example.

## Features of the Text

Each chapter concludes with five standard elements. The last numbered section within the chapter is always devoted to security matters related to the chapter's topic. Next comes three different lists of opportunities for active participation by the student: exercises, programming projects, and exploration projects. Finally, the chapter ends with historical and bibliographic notes.

The distinction between exercises, programming projects, and exploration projects needs explanation. An exercise can be completed with no



outside resources beyond paper and pencil: you need just this textbook and your mind. That does not mean all the exercises are cut and dried, however. Some may call upon you to think creatively; for these, no one answer is correct. Programming projects require a nontrivial amount of programming; that is, they require more than making a small, easily identified change in an existing program. However, a programming project may involve other activities beyond programming. Several of them involve scientific measurement of performance effects, for example; these exploratory aspects may even dominate over the programming aspects. An exploration project, on the other hand, can be an experiment that can be performed with no real programming; at most you might change a designated line within an existing program. The category of exploration projects does not just include experimental work, however. It also includes projects that require you to do research on the Internet or using other library resources.

## Supplemental Resources

The author of this text is making supplemental resources available on his own web site. Additionally, the publisher of the earlier first edition commissioned additional resources from independent supplement authors, which may still be available through the publisher's web site and would largely still apply to this revised edition. The author's web site, <http://gustavus.edu/+max/os-book/>, contains at least the following materials:

- Full text of this revised edition
- Source code in Java, C, or C++ for all programs that are shown in the text
- Artwork files for all figures in the text
- An errata list that will be updated on an ongoing basis

## About the Revised Edition

Course Technology published the first edition of this book in January of 2006 and in October of 2010 assigned the copyright back to the author, giving him the opportunity to make it freely available. This revised edition closely follows the first edition; rather than being a thorough update, it is aimed at three narrow goals:

- All errata reported in the first edition are corrected.
- A variety of other minor improvements appear throughout, such as clarified explanations and additional exercises, projects, and end-of-chapter notes.
- Two focused areas received more substantial updates:
  - The explanation of Linux’s scheduler was completely replaced to correspond to the newer “Completely Fair Scheduler” (CFS), including its group scheduling feature.
  - A new section, 4.9, was added on nonblocking synchronization.

In focusing on these limited goals, a key objective was to maintain as much compatibility with the first edition as possible. Although page numbering changed, most other numbers stayed the same. All new exercises and projects were added to the end of the corresponding lists for that reason. The only newly added section, 4.9, is near the end of its chapter; thus, the only changed section number is that the old Section 4.9 (“Security and Synchronization”) became 4.10. Only in Chapter 4 did any figure numbers change.

It is my hope that others will join me in making further updates and improvements to the text. I am releasing it under a Creative Commons license that allows not just free copying, but also the freedom to make modifications, so long as the modified version is released under the same terms. In order to such modifications practical, I’m not just releasing the book in PDF form, but also as a collection of LaTeX source files that can be edited and then run through the `pdflatex` program (along with `bibtex` and `makeindex`). The source file collection also includes PDF files of all artwork figures; Course Technology has released the rights to the artwork they contracted to have redrawn.

If you produce a modified version of this text, the Creative Commons license allows you considerable flexibility in how you make your modified version available. I would urge you to send it back to me ([max@gustavus.edu](mailto:max@gustavus.edu)) so that I can add your version to the main web site—we will all benefit from having a central repository of progress. Separate materials to supplement the text would also be welcome. One category that occurs to me is animations or screencasts; the static figures in the text are rather limited. Another worthwhile project would be to transform the text into a more contribution-friendly form, such as a wiki.

## Acknowledgments

This book was made possible by financial and logistical support from my employer, Gustavus Adolphus College, and moral support from my family. I would like to acknowledge the contributions of the publishing team, especially developmental editor Jill Batistick and Product Manager Alyssa Pratt. I am also grateful to my students for doing their own fair share of teaching. I particularly appreciate the often extensive comments I received from the following individuals, each of whom reviewed one or more chapters: Dan Cosley, University of Minnesota, Twin Cities; Allen Downey, Franklin W. Olin College of Engineering; Michael Goldweber, Xavier University; Ramesh Karne, Towson University; G. Manimaran, Iowa State University; Alexander Manov, Illinois Institute of Technology; Peter Reiher, University of California, Los Angeles; Rich Salz, DataPower Technology; Dave Schulz, Wisconsin Lutheran College; Sanjeev Setia, George Mason University; and Jon Weissman, University of Minnesota, Twin Cities. Although I did not adopt all their suggestions, I did not ignore any of them, and I appreciate them all.

In preparing the revised edition, I took advantage of suggestions from many readers. I would like to thank all of them, even those I've managed to lose track of, to whom I also apologize. Those I can thank by name are Joel Adams, Michael Brackney, Jack Briner, Justin Delegard, Ben Follis, MinChan Kim, Finn Kuusisto, Matt Lindner, Milo Martin, Gabe Schmidt, Fritz Sieker, and Alex Wauck.

