

Operating Systems and Middleware

Supporting Controlled Interaction

Max Hailperin

Free re-release of first edition

Operating Systems and Middleware: Supporting Controlled Interaction
by Max Hailperin

The commercially published version of this work (ISBN 0-534-42369-8) was Copyright © 2007 by Thomson Course Technology, a division of Thomson Learning, Inc., pursuant to an assignment of rights from the author.

This free re-release is Copyright © 2005-2010 by Max Hailperin, pursuant to an assignment of the rights back to him by Course Technology, a division of Cengage Learning, Inc., successor-in-interest to the publisher. Rights to illustrations rendered by the publisher were also assigned by Course Technology to Max Hailperin and those illustrations are included in the license he grants for this free re-release.

This work is licensed under the Creative Commons Attribution-ShareAlike 3.0 United States License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-sa/3.0/us/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.

The free re-release was prepared from final page proofs and should be completely identical to the commercially published version. In particular, all the errata listed on the web site still apply. (The author intends to release subsequent versions that incorporate the corrections as well as updates and improvements. Subsequent versions may also be in a more easily modifiable form to encourage participation by other contributors. Please email suggestions to max@gustavus.edu.)

Credits from the commercially published version:

Senior Product Manager: Alyssa Pratt
Managing Editor: Mary Franz
Development Editor: Jill Batistick
Senior Marketing Manager: Karen Seitz
Associate Product Manager: Jennifer Smith
Editorial Assistant: Allison Murphy
Senior Manufacturing Coordinator: Justin Palmeiro
Cover Designer: Deborah VanRooyen
Compositor: Interactive Composition Corporation

To
My Family

BRIEF
Contents

PREFACE	xv
1. Introduction	1
2. Threads	19
3. Scheduling	39
4. Synchronization and Deadlocks	72
5. Atomic Transactions	123
6. Virtual Memory	165
7. Processes and Protection	220
8. Files and Other Persistent Storage	269
9. Networking	324
10. Messaging, RPC, and Web Services	369
11. Security	395
APPENDIX A. Stacks	428

TABLE OF
Contents

PREFACE	xv
1. Introduction	1
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	7
1.5 Multiple Computations on One Computer	9
1.6 Controlling Interactions Between Computations	10
1.7 Supporting Interaction Across Time	12
1.8 Supporting Interaction Across Space	14
1.9 Security	15
Exercises	16
Programming Project	17
Exploration Projects	17
Notes	18
2. Threads	19
2.1 Introduction	19
2.2 Example of Multi-Threaded Programs	22
2.3 Reasons for Using Concurrent Threads	24
2.4 Switching Between Threads	27
2.5 Preemptive Multitasking	33
2.6 Security and Threads	34
Exercises	35
Programming Projects	36
Exploration Projects	37
Notes	38

3. Scheduling	39
3.1 Introduction	39
3.2 Thread States	40
3.3 Scheduling Goals	44
3.3.1 Throughput	44
3.3.2 Response Time	47
3.3.3 Urgency, Importance, and Resource Allocation	48
3.4 Fixed-Priority Scheduling	52
3.5 Dynamic-Priority Scheduling	56
3.5.1 Earliest Deadline First Scheduling	56
3.5.2 Decay Usage Scheduling	57
3.6 Proportional-Share Scheduling	62
3.7 Security and Scheduling	65
Exercises	67
Programming Projects	69
Exploration Projects	70
Notes	70
4. Synchronization and Deadlocks	72
4.1 Introduction	72
4.2 Races and the Need for Mutual Exclusion	74
4.3 Mutexes and Monitors	76
4.3.1 The Mutex Application Programming Interface	77
4.3.2 Monitors: A More Structured Interface to Mutexes	81
4.3.3 Underlying Mechanisms for Mutexes	82
4.4 Other Synchronization Patterns	89
4.4.1 Bounded Buffers	89
4.4.2 Readers/Writers Locks	91
4.4.3 Barriers	92
4.5 Condition Variables	94
4.6 Semaphores	98
4.7 Deadlock	101
4.7.1 The Deadlock Problem	101
4.7.2 Deadlock Prevention Through Resource Ordering	104
4.7.3 Ex Post Facto Deadlock Detection	105
4.7.4 Immediate Deadlock Detection	106
4.8 The Interaction of Synchronization with Scheduling	109
4.8.1 Priority Inversion	109
4.8.2 The Convoy Phenomenon	111

4.9 Security and Synchronization	114
Exercises	115
Programming Projects	117
Exploration Projects	119
Notes	121
5. Atomic Transactions	123
5.1 Introduction	123
5.2 Example Applications of Transactions	126
5.2.1 Database Systems	127
5.2.2 Message-Queuing Systems	130
5.2.3 Journalled File Systems	135
5.3 Mechanisms to Ensure Atomicity	137
5.3.1 Serializability: Two-Phase Locking	137
5.3.2 Failure Atomicity: Undo Logging	144
5.4 Transaction Durability: Write-Ahead Logging	147
5.5 Additional Transaction Mechanisms	151
5.5.1 Increased Transaction Concurrency: Reduced Isolation	151
5.5.2 Coordinated Transaction Participants: Two-Phase Commit	154
5.6 Security and Transactions	156
Exercises	159
Programming Project	161
Exploration Projects	162
Notes	163
6. Virtual Memory	165
6.1 Introduction	165
6.2 Uses for Virtual Memory	170
6.2.1 Private Storage	170
6.2.2 Controlled Sharing	171
6.2.3 Flexible Memory Allocation	174
6.2.4 Sparse Address Spaces	176
6.2.5 Persistence	177
6.2.6 Demand-Driven Program Loading	178
6.2.7 Efficient Zero Filling	178
6.2.8 Substituting Disk Storage for RAM	179
6.3 Mechanisms for Virtual Memory	180
6.3.1 Software/Hardware Interface	182
6.3.2 Linear Page Tables	185

x ► Table of Contents

6.3.3	Multilevel Page Tables	190
6.3.4	Hashed Page Tables	194
6.3.5	Segmentation	197
6.4	Policies for Virtual Memory	201
6.4.1	Fetch Policy	202
6.4.2	Placement Policy	204
6.4.3	Replacement Policy	205
6.5	Security and Virtual Memory	212
	Exercises	213
	Programming Projects	215
	Exploration Projects	215
	Notes	218
7.	Processes and Protection	220
7.1	Introduction	220
7.2	POSIX Process Management API	222
7.3	Protecting Memory	230
7.3.1	The Foundation of Protection: Two Processor Modes	231
7.3.2	The Mainstream: Multiple Address Space Systems	234
7.3.3	An Alternative: Single Address Space Systems	235
7.4	Representing Access Rights	237
7.4.1	Fundamentals of Access Rights	237
7.4.2	Capabilities	241
7.4.3	Access Control Lists and Credentials	245
7.5	Alternative Granularities of Protection	252
7.5.1	Protection Within a Process	252
7.5.2	Protection of Entire Simulated Machines	253
7.6	Security and Protection	257
	Exercises	262
	Programming Projects	264
	Exploration Projects	265
	Notes	266
8.	Files and Other Persistent Storage	269
8.1	Introduction	269
8.2	Disk Storage Technology	272
8.3	POSIX File API	275
8.3.1	File Descriptors	275
8.3.2	Mapping Files into Virtual Memory	279

8.3.3 Reading and Writing Files at Specified Positions	281
8.3.4 Sequential Reading and Writing	282
8.4 Disk Space Allocation	283
8.4.1 Fragmentation	284
8.4.2 Locality	287
8.4.3 Allocation Policies and Mechanisms	288
8.5 Metadata	290
8.5.1 Data Location Metadata	291
8.5.2 Access Control Metadata	300
8.5.3 Other Metadata	301
8.6 Directories and Indexing	303
8.6.1 File Directories Versus Database Indexes	303
8.6.2 Using Indexes to Locate Files	304
8.6.3 File Linking	305
8.6.4 Directory and Index Data Structures	308
8.7 Metadata Integrity	310
8.8 Polymorphism in File System Implementations	313
8.9 Security and Persistent Storage	315
Exercises	317
Programming Projects	319
Exploration Projects	319
Notes	322
9. Networking	324
9.1 Introduction	324
9.1.1 Networks and Internets	325
9.1.2 Protocol Layers	327
9.1.3 The End-to-End Principle	330
9.1.4 The Networking Roles of Operating Systems, Middleware, and Application Software	331
9.2 The Application Layer	332
9.2.1 The Web as a Typical Example	332
9.2.2 The Domain Name System: Application Layer as Infrastructure	334
9.2.3 Distributed File Systems: An Application Viewed Through Operating Systems	337
9.3 The Transport Layer	339
9.3.1 Socket APIs	340
9.3.2 TCP, The Dominant Transport Protocol	344
9.3.3 Evolution Within and Beyond TCP	347

xii ► Table of Contents

9.4	The Network Layer	349
9.4.1	IP, Versions 4 and 6	349
9.4.2	Routing and Label Switching	351
9.4.3	Network Address Translation: An End to End-to-End?	352
9.5	The Link and Physical Layers	355
9.6	Network Security	356
9.6.1	Security and the Protocol Layers	357
9.6.2	Firewalls and Intrusion Detection Systems	359
9.6.3	Cryptography	361
	Exercises	365
	Programming Projects	366
	Exploration Projects	367
	Notes	368
10.	Messaging, RPC, and Web Services	369
10.1	Introduction	369
10.2	Messaging Systems	370
10.3	Remote Procedure Call	373
10.3.1	Principles of Operation for RPC	374
10.3.2	An Example Using Java RMI	377
10.4	Web Services	382
10.5	Security and Communication Middleware	387
	Exercises	390
	Programming Projects	392
	Exploration Projects	394
	Notes	394
11.	Security	395
11.1	Introduction	395
11.2	Security Objectives and Principles	396
11.3	User Authentication	401
11.3.1	Password Capture Using Spoofing and Phishing	402
11.3.2	Checking Passwords Without Storing Them	403
11.3.3	Passwords for Multiple, Independent Systems	404
11.3.4	Two-Factor Authentication	405
11.4	Access and Information-Flow Controls	406
11.5	Viruses and Worms	411
11.6	Security Assurance	414
11.7	Security Monitoring	417

11.8 Key Security Best Practices	419
Exercises	421
Programming Projects	424
Exploration Projects	425
Notes	426
APPENDIX A. Stacks	428
A.1 Stack-Allocated Storage: The Concept	429
A.2 Representing a Stack in Memory	430
A.3 Using a Stack for Procedure Activations	431
BIBLIOGRAPHY	434
INDEX	444

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 applies to Windows NT, Windows 2000, Windows XP, Windows Server 2003, and (so far as can be determined from pre-release information) Windows Vista. 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. All Mac OS X examples are from version 10.4, also known as Tiger. However, other than the description of the Spotlight feature for indexed file search, all the material is applicable to earlier 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 one 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 Computing Curricula 2001

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 recently adopted *Computing Curricula 2001* (CC2001) provides 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 CC2001's computer science body of knowledge. Those knowledge units designated as core units within CC2001 are listed in italics. The book covers all core operating systems (OS) units, as well as two elective OS units. The overall amount of coverage for each unit is always at least that recommended by CC2001, 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.

Knowledge unit (<i>italic indicates core units in CC2001</i>)	Chapter(s)
<i>OS1 Overview of operating systems</i>	1
<i>OS2 Operating system principles</i>	1, 7
<i>OS3 Concurrency</i>	2, 4
<i>OS4 Scheduling and dispatch</i>	3
<i>OS5 Memory management</i>	6
OS7 Security and protection	7, 11
OS8 File systems	8
<i>NC1 Introduction to net-centric computing</i>	9
<i>NC2 Communication and networking (partial coverage)</i>	9
<i>NC3 Network security</i>	9
<i>NC4 The web as an example of ... (partial coverage)</i>	9
NC5 Building web applications (partial coverage)	10
IM7 Transaction processing	5

Your Feedback is Welcome

Comments, suggestions, and bug reports are welcome; please send email to 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 has commissioned additional resources from independent supplement authors and is making them available through the Thomson Course Technology web site.

Author's Supplements

The author's web site, <http://www.gustavus.edu/+max/os-book/>, will contain at least the following materials:

- 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

Publisher's Supplements

The publisher's web site, *www.course.com*, will contain the same Java, C, and C++ program files that are available on the author's site and printed in the text. The publisher will provide other supplements as well; the author of each independently created supplement will be listed in the preface of the Instructor's Manual. The following descriptions were provided by the publisher:

Electronic Instructor's Manual The Instructor's Manual that accompanies this textbook includes additional instructional material to assist in class preparation, including Sample Syllabi, Chapter Outlines, Technical Notes, Lecture Notes, Quick Quizzes, Teaching Tips, Discussion Topics, and Key Terms.

ExamView® This objective-based test generator lets the instructor create paper, LAN, or Web-based tests from testbanks designed specifically for this Thomson Course Technology text. Instructors can use the QuickTest Wizard to create tests in fewer than five minutes by taking advantage of Thomson Course Technology's question banks, or they can create customized exams.

PowerPoint Presentations Microsoft PowerPoint slides are included for each chapter. Instructors might use the slides in a variety of ways, including as teaching aids during classroom presentations or as printed handouts for classroom distribution. Instructors can modify the slides provided or include slides of their own for additional topics introduced to the class.

Solutions Solutions to Exercises and Projects are provided on the Teaching Tools CD-ROM and may also be found on the Thomson Course Technology Web site at *www.course.com*. The solutions are password protected.

Figure Files Electronic figure files for all art in the text are available on the Teaching Tools CD-ROM.

Distance Learning Thomson Course Technology is proud to present online test banks in WebCT and Blackboard to provide the most complete and dynamic learning experience possible. For more information on how to access the online test bank, contact your local Thomson Course Technology sales representative.

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

xxii ► Preface

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.