

ann mciver mehoer · ida m. flynn

understanding

# operating systems

sixth  
edition

This is an electronic version of the print textbook. Due to electronic rights restrictions, some third party content may be suppressed. Editorial review has deemed that any suppressed content does not materially affect the overall learning experience. The publisher reserves the right to remove content from this title at any time if subsequent rights restrictions require it. For valuable information on pricing, previous editions, changes to current editions, and alternate formats, please visit [www.cengage.com/highered](http://www.cengage.com/highered) to search by ISBN#, author, title, or keyword for materials in your areas of interest.

# Understanding Operating Systems

*Sixth Edition*

*Ann McIver McHoes*

*Ida M. Flynn*



COURSE TECHNOLOGY  
CENGAGE Learning™

---

Australia • Canada • Mexico • Singapore • Spain • United Kingdom • United States

**Understanding Operating Systems,  
Sixth Edition****Ann McIver McHoes and Ida M. Flynn**

Executive Editor: Marie Lee

Acquisitions Editor: Amy Jollymore

Senior Product Manager: Alyssa Pratt

Editorial Assistant: Zina Kresin

Content Project Manager: Jennifer Feltri

Art Director: Faith Brosnan

Print Buyer: Julio Esperas

Cover Designer: Night &amp; Day Design

Cover Photos: iStockphoto

Proofreader: Suzanne Huizenga

Indexer: Ann McIver McHoes

Compositor: Integra

© 2011 Course Technology, Cengage Learning

ALL RIGHTS RESERVED. No part of this work covered by the copyright herein may be reproduced, transmitted, stored, or used in any form or by any means graphic, electronic, or mechanical, including but not limited to photocopying, recording, scanning, digitizing, taping, Web distribution, information networks, or information storage and retrieval systems, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without the prior written permission of the publisher.

For product information and technology assistance, contact us at  
**Cengage Learning Customer & Sales Support, 1-800-354-9706**

For permission to use material from this text or product,  
submit all requests online at [www.cengage.com/permissions](http://www.cengage.com/permissions)  
Further permissions questions can be e-mailed to  
[permissionrequest@cengage.com](mailto:permissionrequest@cengage.com)

Library of Congress Control Number: 2010920344

ISBN-13: 978-1-4390-7920-1

ISBN-10: 1-4390-7920-x

**Course Technology**20 Channel Center Street  
Boston, MA 02210  
USA

Some of the product names and company names used in this book have been used for identification purposes only and may be trademarks or registered trademarks of their respective manufacturers and sellers.

Any fictional data related to persons, or companies or URLs used throughout this book is intended for instructional purposes only. At the time this book was printed, any such data was fictional and not belonging to any real persons or companies.

Course Technology, a part of Cengage Learning, reserves the right to revise this publication and make changes from time to time in its content without notice.

Cengage Learning is a leading provider of customized learning solutions with office locations around the globe, including Singapore, the United Kingdom, Australia, Mexico, Brazil and Japan. Locate your local office at:  
[www.cengage.com/global](http://www.cengage.com/global)

Cengage Learning products are represented in Canada by  
Nelson Education, Ltd.

To learn more about Course Technology, visit  
[www.cengage.com/coursestechnology](http://www.cengage.com/coursestechnology)

Purchase any of our products at your local college store or at our preferred online store [www.CengageBrain.com](http://www.CengageBrain.com)

*Dedicated to an award-winning teacher and a wonderful  
friend, Ida Moretti Flynn; her love for teaching lives on.*

AMM

# Contents

<b>Part One</b>	<b>Operating Systems Concepts</b>	<b>1</b>
<b>Chapter 1</b>	<b>Introducing Operating Systems</b>	<b>3</b>
	Introduction	4
	What Is an Operating System?	4
	Operating System Software	4
	Main Memory Management	6
	Processor Management	6
	Device Management	7
	File Management	7
	Network Management	7
	User Interface	7
	Cooperation Issues	8
	A Brief History of Machine Hardware	9
	Types of Operating Systems	12
	Brief History of Operating System Development	14
	1940s	14
	1950s	16
	1960s	18
	1970s	19
	1980s	20
	1990s	21
	2000s	22
	Threads	24
	Object-Oriented Design	25
	Conclusion	26
	Key Terms	27
	Interesting Searches	29
	Exercises	29
<b>Chapter 2</b>	<b>Memory Management: Early Systems</b>	<b>31</b>
	Single-User Contiguous Scheme	32
	Fixed Partitions	34

	Dynamic Partitions	36
	Best-Fit Versus First-Fit Allocation	38
	Deallocation	44
	Case 1: Joining Two Free Blocks	45
	Case 2: Joining Three Free Blocks	46
	Case 3: Deallocating an Isolated Block	47
	Relocatable Dynamic Partitions	48
	Conclusion	54
	Key Terms	54
	Interesting Searches	56
	Exercises	56
<b>Chapter 3</b>	<b>Memory Management: Virtual Memory</b>	<b>63</b>
	Paged Memory Allocation	64
	Demand Paging	71
	Page Replacement Policies and Concepts	76
	First-In First-Out	77
	Least Recently Used	79
	The Mechanics of Paging	82
	The Working Set	84
	Segmented Memory Allocation	86
	Segmented/Demand Paged Memory Allocation	89
	Virtual Memory	92
	Cache Memory	94
	Conclusion	98
	Key Terms	100
	Interesting Searches	102
	Exercises	102
<b>Chapter 4</b>	<b>Processor Management</b>	<b>107</b>
	Overview	108
	About Multi-Core Technologies	110
	Job Scheduling Versus Process Scheduling	110
	Process Scheduler	111
	Job and Process Status	113
	Process Control Blocks	114
	PCBs and Queueing	115
	Process Scheduling Policies	116
	Process Scheduling Algorithms	118

First-Come, First-Served	118
Shortest Job Next	120
Priority Scheduling	121
Shortest Remaining Time	122
Round Robin	124
Multiple-Level Queues	127
Case 1: No Movement Between Queues	128
Case 2: Movement Between Queues	128
Case 3: Variable Time Quantum Per Queue	128
Case 4: Aging	129
A Word About Interrupts	129
Conclusion	130
Key Terms	131
Interesting Searches	134
Exercises	134
<b>Chapter 5</b>	<b>Process Management</b>
<b>Chapter 5</b>	<b>139</b>
Deadlock	141
Seven Cases of Deadlock	142
Case 1: Deadlocks on File Requests	142
Case 2: Deadlocks in Databases	143
Case 3: Deadlocks in Dedicated Device Allocation	145
Case 4: Deadlocks in Multiple Device Allocation	145
Case 5: Deadlocks in Spooling	146
Case 6: Deadlocks in a Network	147
Case 7: Deadlocks in Disk Sharing	148
Conditions for Deadlock	149
Modeling Deadlocks	150
Strategies for Handling Deadlocks	153
Starvation	161
Conclusion	163
Key Terms	164
Interesting Searches	165
Exercises	165
<b>Chapter 6</b>	<b>Concurrent Processes</b>
<b>Chapter 6</b>	<b>171</b>
What Is Parallel Processing?	172
Evolution of Multiprocessors	174
Introduction to Multi-Core Processors	174



Typical Multiprocessing Configurations	175
Master/Slave Configuration	175
Loosely Coupled Configuration	176
Symmetric Configuration	177
Process Synchronization Software	178
Test-and-Set	179
WAIT and SIGNAL	180
Semaphores	180
Process Cooperation	183
Producers and Consumers	183
Readers and Writers	185
Concurrent Programming	187
Applications of Concurrent Programming	187
Threads and Concurrent Programming	190
Thread States	191
Thread Control Block	193
Concurrent Programming Languages	193
Java	194
Conclusion	196
Key Terms	197
Interesting Searches	198
Exercises	198
<b>Chapter 7    Device Management</b>	<b>203</b>
Types of Devices	204
Sequential Access Storage Media	205
Direct Access Storage Devices	208
Fixed-Head Magnetic Disk Storage	208
Movable-Head Magnetic Disk Storage	209
Optical Disc Storage	211
CD and DVD Technology	213
Blu-ray Disc Technology	215
Flash Memory Storage	215
Magnetic Disk Drive Access Times	216
Fixed-Head Drives	216
Movable-Head Devices	218
Components of the I/O Subsystem	219
Communication Among Devices	222
Management of I/O Requests	225
Device Handler Seek Strategies	226
Search Strategies: Rotational Ordering	230

RAID	232
Level Zero	234
Level One	234
Level Two	236
Level Three	236
Level Four	236
Level Five	237
Level Six	238
Nested RAID Levels	238
Conclusion	239
Key Terms	240
Interesting Searches	243
Exercises	243
<b>Chapter 8</b>	<b>249</b>
File Management	249
The File Manager	250
Responsibilities of the File Manager	250
Definitions	251
Interacting with the File Manager	252
Typical Volume Configuration	253
Introducing Subdirectories	255
File-Naming Conventions	256
File Organization	258
Record Format	259
Physical File Organization	259
Physical Storage Allocation	263
Contiguous Storage	263
Noncontiguous Storage	264
Indexed Storage	265
Access Methods	267
Sequential Access	268
Direct Access	268
Levels in a File Management System	269
Access Control Verification Module	272
Access Control Matrix	273
Access Control Lists	274
Capability Lists	274
Data Compression	275
Text Compression	275
Other Compression Schemes	276
Conclusion	277

	Key Terms	277
	Interesting Searches	279
	Exercises	280
<b>Chapter 9</b>	<b>Network Organization Concepts</b>	<b>283</b>
	Basic Terminology	284
	Network Topologies	286
	Star	287
	Ring	287
	Bus	289
	Tree	290
	Hybrid	291
	Network Types	292
	Local Area Network	292
	Metropolitan Area Network	293
	Wide Area Network	293
	Wireless Local Area Network	293
	Software Design Issues	295
	Addressing Conventions	295
	Routing Strategies	296
	Connection Models	298
	Conflict Resolution	301
	Transport Protocol Standards	305
	OSI Reference Model	305
	TCP/IP Model	309
	Conclusion	311
	Key Terms	311
	Interesting Searches	313
	Exercises	314
<b>Chapter 10</b>	<b>Management of Network Functions</b>	<b>317</b>
	History of Networks	318
	Comparison of Network and Distributed Operating Systems	318
	DO/S Development	321
	Memory Management	321
	Process Management	323
	Device Management	328
	File Management	330
	Network Management	334

NOS Development	336
Important NOS Features	337
Major NOS Functions	338
Conclusion	339
Key Terms	339
Interesting Searches	340
Exercises	340
<b>Chapter 11   Security and Ethics</b>	<b>343</b>
Role of the Operating System in Security	344
System Survivability	344
Levels of Protection	345
Backup and Recovery	346
Security Breaches	347
Unintentional Intrusions	347
Intentional Attacks	348
System Protection	354
Antivirus Software	355
Firewalls	356
Authentication	357
Encryption	359
Password Management	361
Password Construction	361
Password Alternatives	363
Social Engineering	365
Ethics	366
Conclusion	367
Key Terms	367
Interesting Searches	370
Exercises	370
<b>Chapter 12   System Management</b>	<b>373</b>
Evaluating an Operating System	374
Cooperation Among Components	374
Role of Memory Management	375
Role of Processor Management	375
Role of Device Management	376
Role of File Management	378
Role of Network Management	379

Measuring System Performance	380
Measurement Tools	380
Feedback Loops	383
Patch Management	385
Patching Fundamentals	386
Software Options	388
Timing the Patch Cycle	388
System Monitoring	388
Accounting	391
Conclusion	392
Key Terms	393
Interesting Searches	394
Exercises	394

## **Part Two    Operating Systems in Practice    397**

<b>Chapter 13</b>	<b>UNIX Operating System</b>	<b>401</b>
	Overview	402
	History	402
	The Evolution of UNIX	404
	Design Goals	405
	Memory Management	406
	Process Management	408
	Process Table Versus User Table	409
	Synchronization	411
	Device Management	414
	Device Classifications	414
	Device Drivers	416
	File Management	417
	File Naming Conventions	418
	Directory Listings	419
	Data Structures	422
	User Command Interface	423
	Script Files	425
	Redirection	426
	Pipes	427
	Filters	428
	Additional Commands	429

Conclusion	431
Key Terms	432
Interesting Searches	433
Exercises	433
<b>Chapter 14 MS-DOS Operating System</b>	<b>435</b>
History	436
Design Goals	438
Memory Management	440
Main Memory Allocation	442
Memory Block Allocation	443
Processor Management	444
Process Management	444
Interrupt Handlers	445
Device Management	446
File Management	447
Filename Conventions	447
Managing Files	448
User Interface	452
Batch Files	453
Redirection	454
Filters	455
Pipes	455
Additional Commands	456
Conclusion	458
Key Terms	458
Interesting Searches	459
Exercises	460
<b>Chapter 15 Windows Operating Systems</b>	<b>463</b>
Windows Development	464
Early Windows Products	464
Operating Systems for Home and Professional Users	465
Operating Systems for Networks	466
Design Goals	467
Extensibility	467
Portability	468
Reliability	468
Compatibility	469
Performance	470

Memory Management	470
User-Mode Features	471
Virtual Memory Implementation	472
Processor Management	474
Device Management	476
File Management	480
Network Management	483
Directory Services	484
Security Management	485
Security Basics	486
Security Terminology	486
User Interface	488
Conclusion	493
Key Terms	494
Interesting Searches	495
Exercises	495
<b>Chapter 16 Linux Operating System</b>	<b>499</b>
Overview	500
History	500
Design Goals	502
Memory Management	503
Processor Management	506
Organization of Table of Processes	506
Process Synchronization	507
Process Management	507
Device Management	508
Device Classifications	509
Device Drivers	509
Device Classes	510
File Management	511
Data Structures	511
Filename Conventions	512
Directory Listings	513
User Interface	515
Command-Driven Interfaces	516
Graphical User Interfaces	516
System Monitor	517
Service Settings	517
System Logs	518
Keyboard Shortcuts	518

System Management	519
Conclusion	520
Key Terms	520
Interesting Searches	521
Exercises	522

Appendix

<i>Appendix A</i> ACM Code of Ethics and Professional Conduct	525
---------------------------------------------------------------	-----

Glossary	529
----------	-----

Bibliography	559
--------------	-----

Index	563
-------	-----



# Preface

This book explains a very technical subject in a not-so-technical manner, putting the concepts of operating systems into a format that students can quickly grasp.

For those new to the subject, this text demonstrates what operating systems are, what they do, how they do it, how their performance can be evaluated, and how they compare with each other. Throughout the text we describe the overall function and tell readers where to find more detailed information, if they so desire.

For those with more technical backgrounds, this text introduces the subject concisely, describing the complexities of operating systems without going into intricate detail. One might say this book leaves off where other operating system textbooks begin.

To do so, we've made some assumptions about our audiences. First, we assume the readers have some familiarity with computing systems. Second, we assume they have a working knowledge of an operating system and how it interacts with them. We recommend (although we don't require) that readers be familiar with at least one operating system. In a few places, we found it necessary to include examples using Java or pseudocode to illustrate the inner workings of the operating systems; but, for readers who are unfamiliar with computer languages, we've added a prose description to each example that explains the events in more familiar terms.

## Organization and Features

---

This book is structured to explain the functions of an operating system regardless of the hardware that houses it. The organization addresses a recurring problem with textbooks about technologies that continue to change—that is, the constant advances in evolving subject matter can make textbooks immediately outdated. To address this problem, we've divided the material into two parts: first, the concepts—which do not change quickly—and second, the specifics of operating systems—which change dramatically over the course of years and even months. Our goal is to give readers the ability to apply the topics intelligently, realizing that, although a command, or series of commands, used by one operating system may be different from another, their goals are the same and the functions of the operating systems are also the same.

Although it is more difficult to understand how operating systems work than to memorize the details of a single operating system, understanding general operating system

concepts is a longer-lasting achievement. Such understanding also pays off in the long run because it allows one to adapt as technology changes—as, inevitably, it does. Therefore, the purpose of this book is to give computer users a solid background in the basics of operating systems, their functions and goals, and how they interact and interrelate.

Part One, the first 12 chapters, describes the theory of operating systems. It concentrates on each of the “managers” in turn and shows how they work together. Then it introduces network organization concepts, security, ethics, and management of network functions. Part Two examines actual operating systems, how they apply the theories presented in Part One, and how they compare with each other.

Chapter 1 gives a brief introduction to the subject. The meat of the text begins in Chapters 2 and 3 with memory management because it is the simplest component of the operating system to explain and has historically been tied to the advances from one operating system to the next. We explain the role of the Processor Manager in Chapters 4, 5, and 6, first discussing simple systems and then expanding the discussion to include multiprocessing systems. By the time we reach device management in Chapter 7 and file management in Chapter 8, readers will have been introduced to the four main managers found in every operating system. Chapters 9 and 10 introduce basic concepts related to networking, and Chapters 11 and 12 discuss security, ethics, and some of the tradeoffs that designers consider when attempting to satisfy the needs of their user population.

Each chapter includes learning objectives, key terms, and research topics. For technically oriented readers, the exercises at the end of each chapter include problems for advanced students. Please note that some advanced exercises assume knowledge of matters not presented in the book, but they’re good for those who enjoy a challenge. We expect some readers from a more general background will cheerfully pass them by.

In an attempt to bring the concepts closer to home, throughout the book we’ve added real-life examples to illustrate abstract concepts. However, let no one confuse our conversational style with our considerable respect for the subject matter. The subject of operating systems is a complex one and it cannot be covered completely in these few pages. Therefore, this textbook does not attempt to give an in-depth treatise of operating systems theory and applications. This is the overall view.

Part Two introduces four operating systems in the order of their first release: UNIX, MS-DOS, Windows, and Linux. Here, each chapter discusses how one operating system applies the concepts discussed in Part One and how it compares with the others. Again, we must stress that this is a general discussion—an in-depth examination of an operating system would require details based on its current standard version, which can’t be done here. We strongly suggest that readers use our discussion as a guide, a base to work from, when comparing the pros and cons of a specific operating system and supplement our work with research that’s as current as possible.

The text concludes with several reference aids. Terms that are important within a chapter are listed at its conclusion as key terms. The extensive end-of-book Glossary

includes brief definitions for hundreds of terms used in these pages. The Bibliography can guide the reader to basic research on the subject. Finally, the Appendix features the ACM Code of Ethics.

Not included in this text is a discussion of databases and data structures, except as examples of process synchronization problems, because they only tangentially relate to operating systems and are frequently the subject of other courses. We suggest that readers begin by learning the basics as presented in the following pages before pursuing these complex subjects.

## Changes to the Sixth Edition

---

This edition has been thoroughly updated and features many improvements over the fifth edition:

- New references to Macintosh OS X, which is based on UNIX
- Numerous new homework exercises in every chapter
- Updated references to the expanding influence of wireless technology
- More networking information throughout the text
- Continuing emphasis on system security and patch management
- More discussion describing the management of multiple processors
- Updated detail in the chapters that discuss UNIX, Windows, and Linux
- New research topics and student exercises for the chapters on UNIX, MS-DOS, Windows, and Linux

Other changes throughout the text are editorial clarifications, expanded captions, and improved illustrations.

## A Note for Instructors

---

The following supplements are available when this text is used in a classroom setting:

**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, and Discussion Topics.

**Distance Learning.** Course Technology is proud to present online test banks in WebCT and Blackboard to provide the most complete and dynamic learning experience possible. Instructors are encouraged to make the most of the course, both online and offline. For more information on how to access the online test bank, contact your local Course Technology sales representative.

**PowerPoint Presentations.** This book comes with Microsoft PowerPoint slides for each chapter. These are included as a teaching aid for classroom presentations, either to make available to students on the network for chapter review, or to be printed for classroom distribution. Instructors can add their own slides for additional topics that they introduce to the class.

**Solutions.** Selected solutions to Review Questions and Exercises are provided on the Instructor Resources CD-ROM and may also be found on the Cengage Course Technology Web site at [www.cengage.com/coursestechnology](http://www.cengage.com/coursestechnology). The solutions are password protected.

**Order of Presentation.** We have built this text with a modular construction to accommodate several presentation options, depending on the instructor's preference. For example, the syllabus can follow the chapters as listed in Chapter 1 through Chapter 12 to present the core concepts that all operating systems have in common. Using this path, students will learn about the management of memory, processors, devices, files, and networks, in that order. An alternative path might begin with Chapter 1, move next to processor management in Chapters 4 through 6, then to memory management in Chapters 2 and 3, touch on systems security and management in Chapters 11 and 12, and finally move to device and file management in Chapters 7 and 8. Because networking is often the subject of another course, instructors may choose to bypass Chapters 9 and 10, or include them for a more thorough treatment of operating systems.

We hope you find our discussion of ethics helpful in Chapter 11, which is included in response to requests by university adopters of the text who want to discuss this subject in their lectures.

In Part Two, we examine details about four specific operating systems in an attempt to show how the concepts in the first 12 chapters are applied by a specific operating system. In each case, the chapter is structured in a similar manner as the chapters in Part One. That is, they discuss the management of memory, processors, files, devices, networks, and systems. In addition, each includes an introduction to one or more user interfaces for that operating system. With this edition, we added exercises and research topics to each of these chapters to help students explore issues discussed in the preceding pages.

For the first time, we included references to the Macintosh OS X operating system in the UNIX chapter.

We continue to include MS-DOS in spite of its age because faculty reviewers and adopters have specifically requested it, presumably so students can learn the basics of this command-driven interface using a Windows emulator.

If you have suggestions for inclusion in this text, please send them along. Although we are squeezed for space, we are pleased to consider all possibilities.

## Acknowledgments

---

Our gratitude goes to all of our friends and colleagues, who were so generous with their encouragement, advice, and support. Special thanks go to Robert Kleinmann, Eleanor Irwin, Charles R. Woratschek, Terri Lennox, and Roger Flynn for their assistance.

Special thanks also to those at Course Technology, Brooks/Cole, and PWS Publishing who made significant contributions to all six editions of this text, especially Alyssa Pratt, Kallie Swanson, Mike Sugarman, and Mary Thomas Stone. In addition, the following individuals made key contributions to this edition: Jennifer Feltri, Content Project Manager, and Sreejith Govindan, Integra.

We deeply appreciate the comments of the reviewers who helped us refine this edition:

Proposal Reviewers:

Nisheeth Agrawal: Calhoun Community College

Brian Arthur: Mary Baldwin College

Margaret Moore: University of Phoenix

Chapter Reviewers:

Kent Einspahr: Concordia University

Gary Heisler: Lansing Community College

Paul Hemler: Hampden-Sydney College

And to the many students and instructors who have sent helpful comments and suggestions since publication of the first edition in 1991, we thank you. Please keep them coming.

*Ann McIver McHoes*, mchoesa@duq.edu

*Ida M. Flynn*

*This page intentionally left blank*

## Part One

# Operating Systems Concepts

*“So work the honey-bees,  
Creatures that by a rule in nature teach  
The act of order to a peopled kingdom.”*

—William Shakespeare (1564–1616; in Henry V)

All operating systems have certain core items in common: each must manage memory, processing capability, devices and peripherals, files, and networks. In Part One of this text we present an overview of these operating systems essentials.

- Chapter 1 introduces the subject.
- Chapters 2–3 discuss main memory management.
- Chapters 4–6 cover processor management.
- Chapter 7 concentrates on device management.
- Chapter 8 is devoted to file management.
- Chapters 9–10 briefly review networks.
- Chapter 11 discusses system security issues.
- Chapter 12 explores system management and the interaction of the operating system’s components.

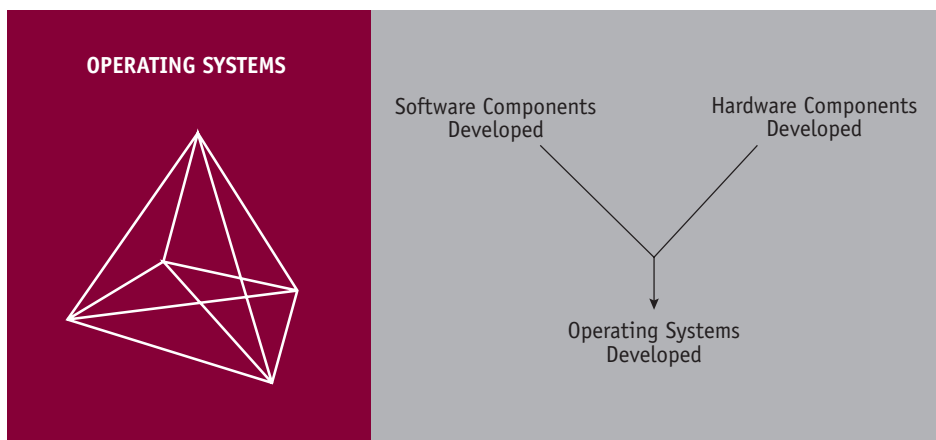
Then, in Part Two of the text (Chapters 13–16), we look at specific operating systems and how they apply the theory presented here in Part One.

Throughout our discussion of this very technical subject, we try to include definitions of terms that might be unfamiliar to you. However, it isn't always possible to describe a function and define the technical terms while keeping the explanation clear. Therefore, we've put the key terms with definitions at the end of each chapter, and at the end of the text is an extensive glossary for your reference. Items listed in the Key Terms are shown in boldface the first time they appear.

Throughout the book we keep our descriptions and examples as simple as possible to introduce you to the system's complexities without getting bogged down in technical detail. Therefore, be aware that for almost every topic explained in the following pages, there's much more information that can be studied. Our goal is to introduce you to the subject, and to encourage you to pursue your interest using other texts or primary sources if you need more detail.



# Introducing Operating Systems



*“I think there is a world market for maybe five computers.”*

—Thomas J. Watson (1874–1956; chairman of IBM 1949–1956)

---

## Learning Objectives

After completing this chapter, you should be able to describe:

- Innovations in operating system development
  - The basic role of an operating system
  - The major operating system software subsystem managers and their functions
  - The types of machine hardware on which operating systems run
  - The differences among batch, interactive, real-time, hybrid, and embedded operating systems
  - Multiprocessing and its impact on the evolution of operating system software
  - Virtualization and core architecture trends in new operating systems
-

## Introduction

---

To understand an operating system is to understand the workings of an entire computer system, because the operating system manages each and every piece of hardware and software. This text explores what operating systems are, how they work, what they do, and why.

This chapter briefly describes how simple operating systems work and how, in general, they've evolved. The following chapters explore each component in more depth and show how its function relates to the other parts of the operating system. In other words, you see how the pieces work harmoniously to keep the computer system working smoothly.

## What Is an Operating System?

---

A computer system consists of **software** (programs) and **hardware** (the physical machine and its electronic components). The **operating system** software is the chief piece of software, the portion of the computing system that manages all of the hardware and all of the other software. To be specific, it controls every file, every device, every section of main memory, and every nanosecond of processing time. It controls who can use the system and how. In short, it's the boss.

Therefore, each time the user sends a command, the operating system must make sure that the command is executed; or, if it's not executed, it must arrange for the user to get a message explaining the error. Remember: This doesn't necessarily mean that the operating system executes the command or sends the error message—but it does control the parts of the system that do.

## Operating System Software

---

The pyramid shown in Figure 1.1 is an abstract representation of an operating system and demonstrates how its major components work together.

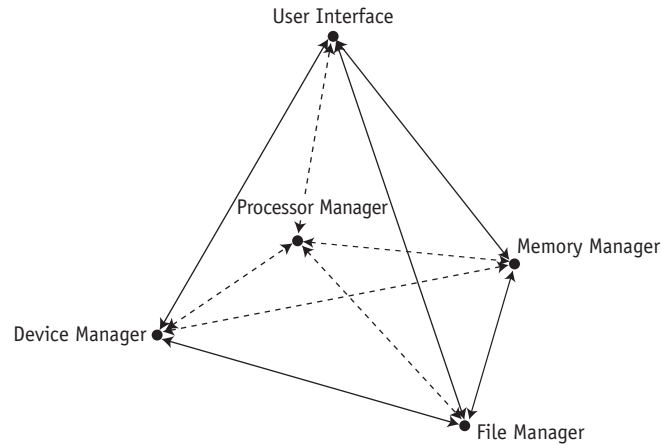
At the base of the pyramid are the four essential managers of every operating system: the **Memory Manager**, the **Processor Manager**, the **Device Manager**, and the **File Manager**. In fact, these managers are the basis of all operating systems and each is discussed in detail throughout the first part of this book. Each manager works closely with the other managers and performs its unique role regardless of which specific operating system is being discussed. At the top of the pyramid is the User Interface, from which users issue commands to the operating system. This is the component that's unique to each operating system—sometimes even between different versions of the same operating system.



Unless we mention networking or the Internet, our discussions apply to the most basic elements of operating systems. Chapters 9 and 10 are dedicated to networking.

**(figure 1.1)**

*This model of a non-networked operating system shows four subsystem managers supporting the User Interface.*



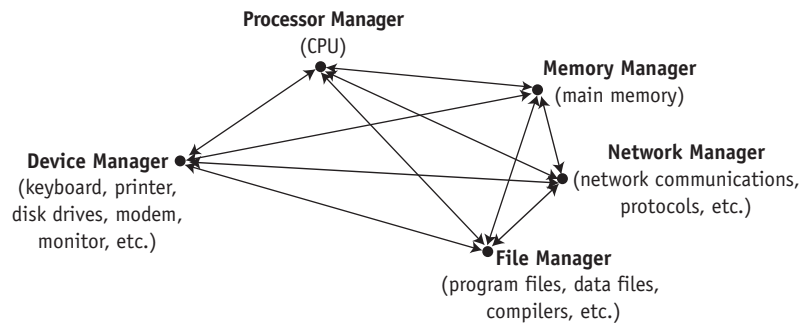
A **network** was not always an integral part of operating systems; early systems were self-contained with all network capability added on top of existing operating systems. Now most operating systems routinely incorporate a **Network Manager**. The base of a pyramid for a networked operating system is shown in Figure 1.2.

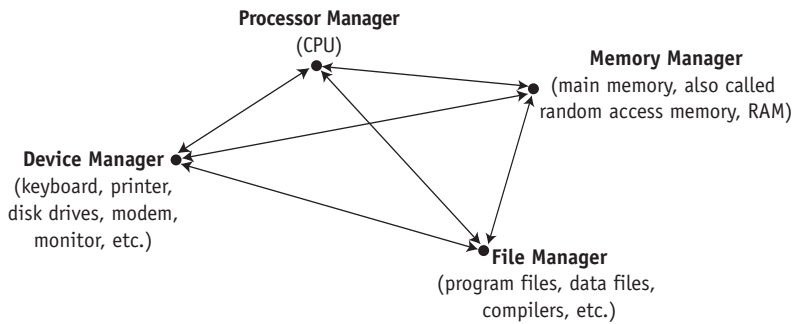
Regardless of the size or configuration of the system, each of the subsystem managers, shown in Figure 1.3, must perform the following tasks:

- Monitor its resources continuously
- Enforce the policies that determine who gets what, when, and how much
- Allocate the resource when appropriate
- Deallocate the resource when appropriate

**(figure 1.2)**

*Networked systems have a Network Manager that assumes responsibility for networking tasks while working harmoniously with every other manager.*





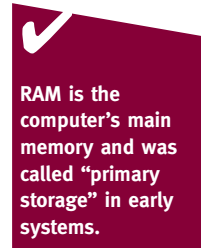
(figure 1.3)

*Each subsystem manager at the base of the pyramid takes responsibility for its own tasks while working harmoniously with every other manager.*

## Main Memory Management

The Memory Manager (the subject of Chapters 2–3) is in charge of main memory, also known as RAM, short for Random Access Memory. The Memory Manager checks the validity of each request for memory space and, if it is a legal request, it allocates a portion of memory that isn't already in use. In a multiuser environment, the Memory Manager sets up a table to keep track of who is using which section of memory. Finally, when the time comes to reclaim the memory, the Memory Manager deallocates memory.

A primary responsibility of the Memory Manager is to protect the space in main memory occupied by the operating system itself—it can't allow any part of it to be accidentally or intentionally altered.



## Processor Management

The Processor Manager (the subject of Chapters 4–6) decides how to allocate the central processing unit (CPU). An important function of the Processor Manager is to keep track of the status of each process. A process is defined here as an instance of execution of a program.

The Processor Manager monitors whether the CPU is executing a process or waiting for a READ or WRITE command to finish execution. Because it handles the processes' transitions from one state of execution to another, it can be compared to a traffic controller. Once the Processor Manager allocates the processor, it sets up the necessary registers and tables and, when the job is finished or the maximum amount of time has expired, it reclaims the processor.

Think of it this way: The Processor Manager has two levels of responsibility. One is to handle jobs as they enter the system and the other is to manage each process within those jobs. The first part is handled by the Job Scheduler, the high-level portion of the Processor Manager, which accepts or rejects the incoming jobs. The second part is

handled by the Process Scheduler, the low-level portion of the Processor Manager, which is responsible for deciding which process gets the CPU and for how long.

## Device Management

The Device Manager (the subject of Chapter 7) monitors every device, channel, and control unit. Its job is to choose the most efficient way to allocate all of the system's devices, printers, ports, disk drives, and so forth, based on a scheduling policy chosen by the system's designers.

The Device Manager does this by allocating each resource, starting its operation, and, finally, deallocating the device, making it available to the next process or job.

## File Management

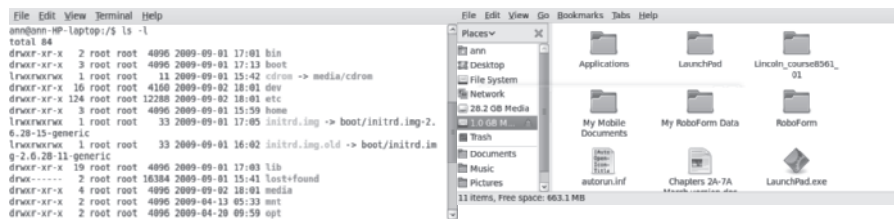
The File Manager (the subject of Chapter 8) keeps track of every file in the system, including data files, program files, compilers, and applications. By using predetermined access policies, it enforces restrictions on who has access to which files. The File Manager also controls what users are allowed to do with files once they access them. For example, a user might have read-only access, read-and-write access, or the authority to create and delete files. Managing access control is a key part of file management. Finally, the File Manager allocates the necessary resources and later deallocates them.

## Network Management

Operating systems with Internet or networking capability have a fifth essential manager called the Network Manager (the subject of Chapters 9–10) that provides a convenient way for users to share resources while controlling users' access to them. These resources include hardware (such as CPUs, memory areas, printers, tape drives, modems, and disk drives) and software (such as compilers, application programs, and data files).

## User Interface

The user interface is the portion of the operating system that users interact with directly. In the old days, the user interface consisted of commands typed on a keyboard and displayed on a monitor, as shown in Figure 1.4. Now most systems allow users to choose a menu option from a list. The user interface, desktops, and formats vary widely from one operating system to another, as shown in Chapters 13–16 in Part Two of this text.



(figure 1.4)

Two user interfaces from Linux: a command-driven interface (left) and a menu-driven interface (right).

## Cooperation Issues

However, it is not enough for each manager to perform its individual tasks. It must also be able to work harmoniously with every other manager. Here is a simplified example. Let's say someone chooses an option from a menu to execute a program. The following major steps must occur in sequence:

1. The Device Manager must receive the electrical impulses from the mouse or keyboard, form the command, and send the command to the User Interface, where the Processor Manager validates the command.
2. The Processor Manager then sends an acknowledgment message to be displayed on the monitor so the user realizes the command has been sent.
3. When the Processor Manager receives the command, it determines whether the program must be retrieved from storage or is already in memory, and then notifies the appropriate manager.
4. If the program is in storage, the File Manager must calculate its exact location on the disk and pass this information to the Device Manager, which retrieves the program and sends it to the Memory Manager.
5. The Memory Manager then finds space for it and records its exact location in memory. Once the program is in memory, the Memory Manager must track its location in memory (even if it's moved) as well as its progress as it's executed by the Processor Manager.
6. When the program has finished executing, it must send a finished message to the Processor Manager so that the processor can be assigned to the next program waiting in line.
7. Finally, the Processor Manager must forward the finished message to the Device Manager, so that it can notify the user and refresh the screen.

Although this is a vastly oversimplified demonstration of a complex operation, it illustrates some of the incredible precision required for the operating system to work smoothly. So although we'll be discussing each manager in isolation for much of this text, no single manager could perform its tasks without the active cooperation of every other part.

## A Brief History of Machine Hardware

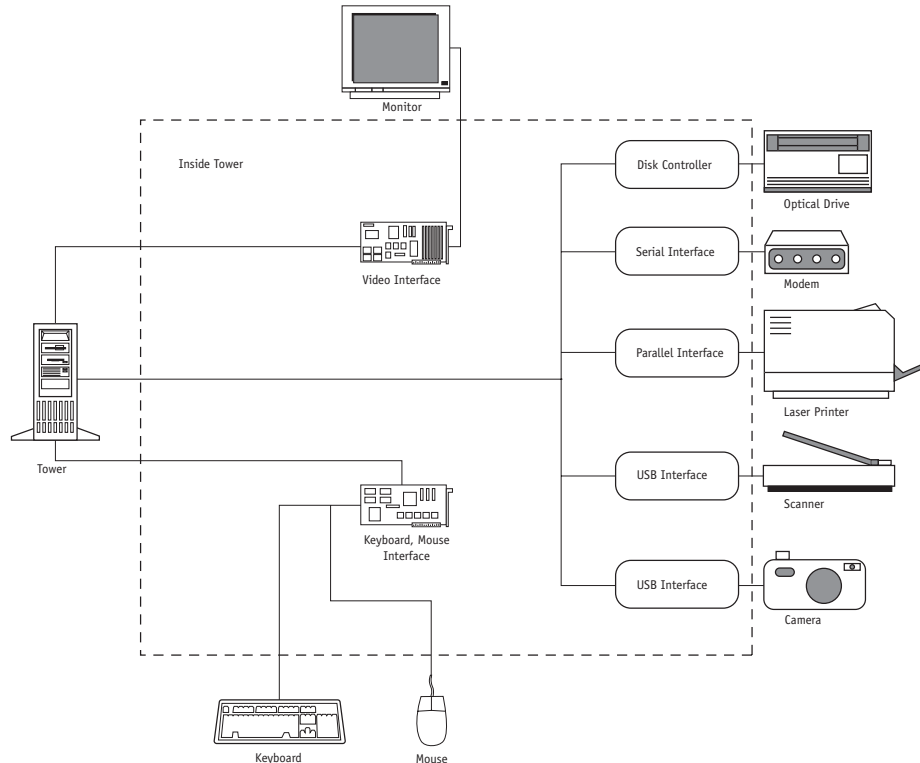
To appreciate the role of the operating system (which is software), we need to discuss the essential aspects of the computer system's hardware, the physical machine and its electronic components, including memory chips, input/output devices, **storage** devices, and the central processing unit (CPU).

- **Main memory** (random access memory, RAM) is where the data and instructions must reside to be processed.
- **I/O devices**, short for input/output devices, include every peripheral unit in the system such as printers, disk drives, CD/DVD drives, flash memory, keyboards, and so on.
- The **central processing unit** (CPU) is the brains with the circuitry (sometimes called the chip) to control the interpretation and execution of instructions. In essence, it controls the operation of the entire computer system, as illustrated in Figure 1.5. All storage references, data manipulations, and I/O operations are initiated or performed by the CPU.

Until the mid-1970s, computers were classified by capacity and price. A **mainframe** was a large machine—in size and in internal memory capacity. The IBM 360, introduced in

(figure 1.5)

*A logical view of a typical computer system hardware configuration. The tower holds the central processing unit, the arithmetic and logic unit, registers, cache, and main memory, as well as controllers and interfaces shown within the dotted lines.*



1964, is a classic example of an early mainframe. The IBM 360 model 30 required an air-conditioned room about 18 feet square to house the CPU, the operator's console, a printer, a card reader, and a keypunch machine. The CPU was 5 feet high and 6 feet wide, had an internal memory of 64K (considered large at that time), and a price tag of \$200,000 in 1964 dollars. Because of its size and price at the time, its applications were generally limited to large computer centers belonging to the federal government, universities, and very large businesses.

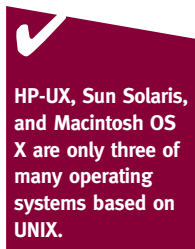
The **minicomputer** was developed to meet the needs of smaller institutions, those with only a few dozen users. One of the early minicomputers was marketed by Digital Equipment Corporation to satisfy the needs of large schools and small colleges that began offering computer science courses in the early 1970s. (The price of its PDP-8 was less than \$18,000.) Minicomputers are smaller in size and memory capacity and cheaper than mainframes. Today, computers that fall between microcomputers and mainframes in capacity are often called midrange computers.

The **supercomputer** was developed primarily for government applications needing massive and fast number-crunching ability to carry out military operations and weather forecasting. Business and industry became interested in the technology when the massive computers became faster and less expensive. A Cray supercomputer is a typical example with six to thousands of processors performing up to 2.4 trillion floating point operations per second (2.4 teraflops). Supercomputers are used for a wide range of tasks from scientific research to customer support and product development. They're often used to perform the intricate calculations required to create animated motion pictures. And they help oil companies in their search for oil by analyzing massive amounts of data (Stair, 1999).

The **microcomputer** was developed to offer inexpensive computation capability to individual users in the late 1970s. Early models featured a revolutionary amount of memory: 64K. Their physical size was smaller than the minicomputers of that time, though larger than the microcomputers of today. Eventually, microcomputers grew to accommodate software with larger capacity and greater speed. The distinguishing characteristic of the first microcomputer was its single-user status.

Powerful microcomputers developed for use by commercial, educational, and government enterprises are called **workstations**. Typically, workstations are networked together and are used to support engineering and technical users who perform massive mathematical computations or computer-aided design (CAD), or use other applications requiring very powerful CPUs, large amounts of main memory, and extremely high-resolution graphic displays to meet their needs.

**Servers** are powerful computers that provide specialized services to other computers on client/server networks. Examples can include print servers, Internet servers, e-mail servers, etc. Each performs critical network tasks. For instance, a file server, usually a





powerful computer with substantial file storage capacity (such as a large collection of hard drives), manages file storage and retrieval for other computers, called clients, on the network.

(table 1.1)

A brief list of platforms and sample operating systems listed in alphabetical order.

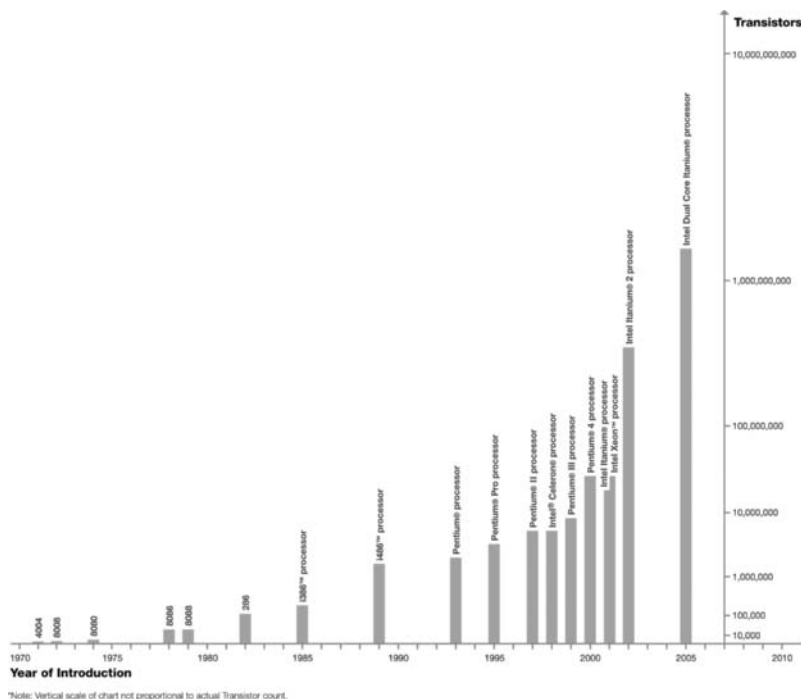
Platform	Operating System
Microcomputers	Linux, UNIX (includes Mac), Windows
Mainframe computers	IBM z/390, Linux, UNIX
Supercomputers	IRIX, Linux, UNICOS
Workstations, servers	Linux, UNIX, Windows
Networks	Linux, NetWare, UNIX, Windows
Personal digital assistants	BlackBerry, Linux, Palm OS, Windows Mobile

Some typical operating systems for a wide variety of platforms are shown in Table 1.1. Since the mid-1970s, rapid advances in computer technology have blurred the distinguishing characteristics of early machines: physical size, cost, and memory capacity. The most powerful mainframes today have multiple processors coordinated by the Processor Manager. Simple mainframes still have a large main memory, but now they’re available in desk-sized cabinets.

Networking is an integral part of modern computer systems because it can connect workstations, servers, and peripheral devices into integrated computing systems. Networking capability has become a standard feature in many computing devices: personal organizers, personal digital assistants (PDAs), cell phones, and handheld Web browsers.

At one time, computers were classified by memory capacity; now they’re distinguished by processor capacity. We must emphasize that these are relative categories and what is large today will become medium-sized and then small sometime in the near future.

In 1965, Intel executive Gordon Moore observed that each new processor chip contained roughly twice as much capacity as its predecessor, and each chip was released within 18–24 months of the previous chip. He predicted that the trend would cause computing power to rise exponentially over relatively brief periods of time. Now known as Moore’s Law, shown in Figure 1.6, the trend has continued and is still remarkably accurate. The Intel 4004 chip in 1971 had 2,300 transistors while the Pentium II chip 20 years later had 7.5 million, and the Pentium 4 Extreme Edition processor introduced in 2004 had 178 million transistors. Moore’s Law is often used by industry observers to make their chip capacity forecasts.



(figure 1.6)

*Demonstration of Moore's Law. Gordon Moore's 1965 prediction has held up for more than three decades.*

*Copyright © 2005 Intel Corporation*

## Types of Operating Systems

Operating systems for computers large and small fall into five categories distinguished by response time and how data is entered into the system: batch, interactive, real-time, hybrid, and embedded systems.

**Batch systems** date from the earliest computers, when they relied on stacks of punched cards or reels of magnetic tape for input. Jobs were entered by assembling the cards into a deck and running the entire deck of cards through a card reader as a group—a batch. The efficiency of a batch system is measured in **throughput**—the number of jobs completed in a given amount of time (for example, 550 jobs per hour).

**Interactive systems** give a faster turnaround than batch systems but are slower than the real-time systems we talk about next. They were introduced to satisfy the demands of users who needed fast turnaround when debugging their programs. The operating system required the development of time-sharing software, which would allow each user to interact directly with the computer system via commands entered from a type-writer-like terminal. The operating system provides immediate feedback to the user and response time can be measured in fractions of a second.

**Real-time systems** are used in time-critical environments where reliability is key and data must be processed within a strict time limit. The time limit need not be ultra-fast

**(figure 1.7)**

*The state-of-the-art computer interface box for the Apollo spacecraft in 1968. The guidance computer had few moving parts and no vacuum tubes, making it both rugged and compact.*

*Courtesy of NASA*



(though it often is), but system response time must meet the deadline or risk significant consequences. These systems also need to provide contingencies to fail gracefully—that is, preserve as much of the system’s capabilities and data as possible to facilitate recovery. For example, real-time systems are used for space flights (as shown in Figure 1.7), airport traffic control, fly-by-wire aircraft, critical industrial processes, certain medical equipment, and telephone switching, to name a few.

There are two types of real-time systems depending on the consequences of missing the deadline:

- Hard real-time systems risk total system failure if the predicted time deadline is missed.
- Soft real-time systems suffer performance degradation, but not total system failure, as a consequence of a missed deadline.

Although it’s theoretically possible to convert a general-purpose operating system into a real-time system by merely establishing a deadline, the unpredictability of these systems can’t provide the guaranteed response times that real-time performance requires (Dougherty, 1995). Therefore, most embedded systems and real-time environments require operating systems that are specially designed to meet real-time needs.

**Hybrid systems** are a combination of batch and interactive. They appear to be interactive because individual users can access the system and get fast responses, but such a system actually accepts and runs batch programs in the background when the interactive load is light. A hybrid system takes advantage of the free time between high-demand usage of the system and low-demand times. Many large computer systems are hybrids.

**Embedded systems** are computers placed inside other products to add features and capabilities. For example, you find embedded computers in household appliances, automobiles, digital music players, elevators, and pacemakers. In the case of automobiles, embedded computers can help with engine performance, braking, and navigation. For example, several projects are under way to implement “smart roads,” which would alert drivers in cars equipped with embedded computers to choose alternate routes when traffic becomes congested.

Operating systems for embedded computers are very different from those for general computer systems. Each one is designed to perform a set of specific programs, which are not interchangeable among systems. This permits the designers to make the operating system more efficient and take advantage of the computer’s limited resources, such as memory, to their maximum.

Before a general-purpose operating system, such as Linux, UNIX, or Windows, can be used in an embedded system, the system designers must select which components, from the entire operating system, are needed in that particular environment. The final version of this operating system will include only the necessary elements; any unneeded features or functions will be dropped. Therefore, operating systems with a small **kernel** (the core portion of the software) and other functions that can be mixed and matched to meet the embedded system requirements will have potential in this market.

## Brief History of Operating System Development

The evolution of operating system software parallels the evolution of the computer hardware it was designed to control. Here’s a very brief overview of this evolution.

### 1940S

The first generation of computers (1940–1955) was a time of vacuum tube technology and computers the size of classrooms. Each computer was unique in structure and purpose. There was little need for standard operating system software because each computer’s use was restricted to a few professionals working on mathematical, scientific, or military applications, all of whom were familiar with the idiosyncrasies of their hardware.

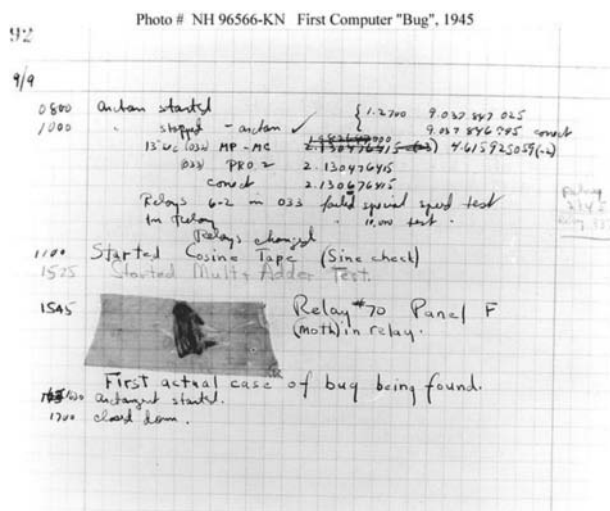
A typical program would include every instruction needed by the computer to perform the tasks requested. It would give explicit directions to the card reader (when to begin, how to interpret the data on the cards, when to end), the CPU (how and where to store the instructions in memory, what to calculate, where to find the data, where to send



One example of a software product to help developers build an embedded system is Windows Automotive.

(figure 1.8)

Dr. Grace Hopper's research journal from her work on Harvard's Mark I computer in 1945 included the remains of the first computer "bug," a moth that had become trapped in the computer's relays causing the system to crash. Today's use of the term "bug" stems from that first moth.



the output), and the output device (when to begin, how to print out the finished product, how to format the page, and when to end).

The machines were operated by the programmers from the main console—it was a hands-on process. In fact, to debug a program, the programmer would stop the processor, read the contents of each register, make the corrections in memory locations, and then resume operation. The first bug was a moth trapped in a Harvard computer that caused it to fail, as shown in Figure 1.8.

To run programs, the programmers would have to reserve the machine for the length of time they estimated it would take the computer to execute the program. As a result, the machine was poorly utilized. The CPU processed data and made calculations for only a fraction of the available time and, in fact, the entire system sat idle between reservations.

In time, computer hardware and software became more standard and the execution of a program required fewer steps and less knowledge of the internal workings of the computer. Compilers and assemblers were developed to translate into binary code the English-like commands of the evolving high-level languages.

Rudimentary operating systems started to take shape with the creation of macros, library programs, standard subroutines, and utility programs. And they included device driver subroutines—prewritten programs that standardized the way input and output devices were used.

These early programs were at a significant disadvantage because they were designed to use their resources conservatively at the expense of understandability. That meant

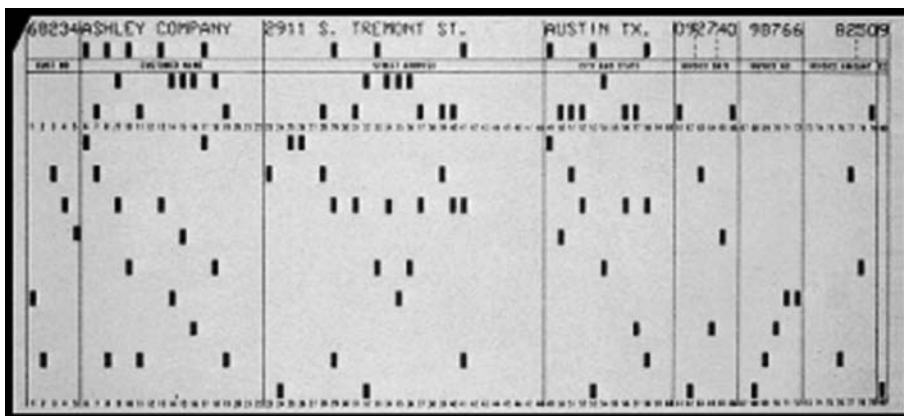
that many programs used convoluted logic that only the original programmer could understand, so it was nearly impossible for anyone else to debug or change the program later on.

## 1950s

Second-generation computers (1955–1965) were developed to meet the needs of new markets—government and business researchers. The business environment placed much more importance on the cost effectiveness of the system. Computers were still very expensive, especially when compared to other office equipment (the IBM 7094 was priced at \$200,000). Therefore, throughput had to be maximized to make such an investment worthwhile for business use, which meant dramatically increasing the usage of the system.

Two improvements were widely adopted: Computer operators were hired to facilitate each machine's operation, and job scheduling was instituted. Job scheduling is a productivity improvement scheme that groups together programs with similar requirements. For example, several FORTRAN programs would be run together while the FORTRAN compiler was still resident in memory. Or all of the jobs using the card reader for input might be run together, and those using the tape drive would be run later. Some operators found that a mix of I/O device requirements was the most efficient combination. That is, by mixing tape-input programs with card-input programs, the tapes could be mounted or rewound while the card reader was busy. A typical punch card is shown in Figure 1.9.

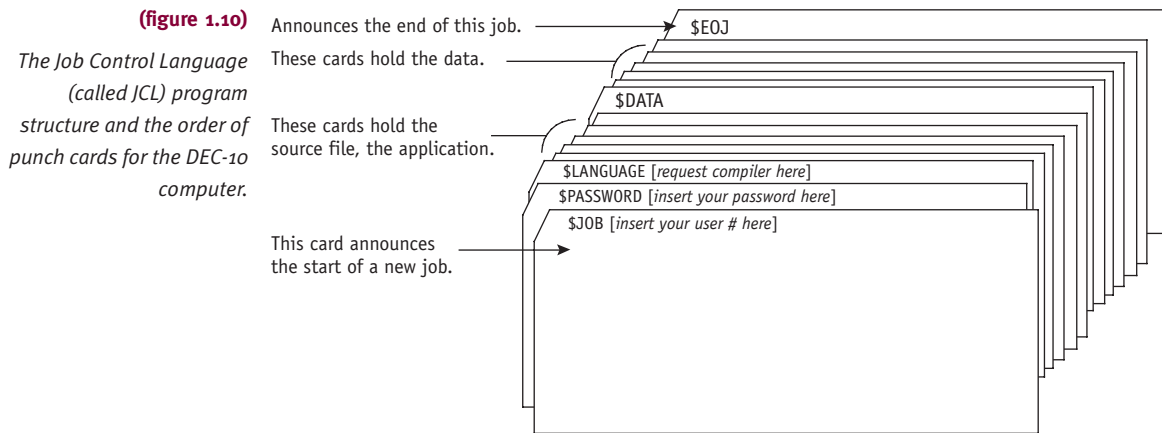
Job scheduling introduced the need for control cards, which defined the exact nature of each program and its requirements, illustrated in Figure 1.10. This was one of the first uses of a job control language, which helped the operating system coordinate and manage the system resources by identifying the users and their jobs and specifying the resources required to execute each job.



(figure 1.9)

Each letter or number printed along the top of the punch card is represented by a unique combination of holes beneath it.

From ibm.com



But even with batching techniques, the faster second-generation computers allowed expensive time lags between the CPU and the I/O devices. For example, a job with 1600 cards could take 79 seconds to be read by the card reader and only 5 seconds of CPU time to assemble or compile. That meant the CPU was idle 94 percent of the time and busy only 6 percent of the time it was dedicated to that job—an inefficiency that resulted in poor overall system use.

Eventually, several factors helped improve the performance of the CPU:

- First, the speeds of I/O devices such as drums, tape drives, and disks gradually increased.
- Second, to use more of the available storage area in these devices, records were grouped into blocks before they were retrieved or stored. (This is called blocking, meaning that several logical records are grouped within one physical record, and is discussed in detail in Chapter 7.)
- Third, to reduce the discrepancy in speed between the I/O and the CPU, an interface called the control unit was placed between them to act as a buffer. A buffer is an interim storage area that works as a temporary holding place. As the slow input device reads one record, the control unit places each character of the record into the buffer. When the buffer is full, the entire record is quickly transmitted to the CPU. The process is just the opposite for output devices: The CPU places the entire record into the buffer, which is then passed on by the control unit at the slower rate required by the output device.

The buffers of this generation were conceptually similar to those now used routinely by Internet browsers to make video and audio playback smoother, as shown in Figure 1.11.

If a control unit had more than one buffer, the I/O process could be made even faster. For example, if the control unit had two buffers, the second buffer could be loaded while the first buffer was transmitting its contents to or from the CPU. Ideally, by

the time the first was transmitted, the second was ready to go, and so on. In this way, input or output time was cut in half.



(figure 1.11)

Three typical browser buffering progress indicators.

- Fourth, in addition to buffering, an early form of spooling was developed by moving offline the operations of card reading, printing, and “punching.” For example, incoming jobs would be transferred from card decks to reels of magnetic tape offline. Then they could be read into the CPU from the tape at a speed much faster than that of the card reader. The spooler worked the same way as a buffer but, in this example, it was a separate offline device while a buffer was part of the main computer hardware.

Also during the second generation, techniques were developed to manage program libraries, create and maintain each data direct access address, and create and check file labels. Timer interrupts were developed to allow job sharing and to prevent infinite loops on programs that were mistakenly instructed to execute a single series of commands forever. Because a fixed amount of execution time was allocated to each program when it entered the system, and was monitored by the operating system, programs that were still running when the time expired were terminated.

During the second generation, programs were still assigned to the processor one at a time. The next step toward better use of the system’s resources was the move to shared processing.

## 1960s

Third-generation computers date from the mid-1960s. They were designed with faster CPUs, but their speed still caused problems when they interacted with printers and other I/O devices that ran at slower speeds. The solution was **multiprogramming**, which introduced the concept of loading many programs at one time and sharing the attention of a single CPU.

The first multiprogramming systems allowed each program to be serviced in turn, one after another. The most common mechanism for implementing multiprogramming was the introduction of the concept of the interrupt, whereby the CPU was notified of events needing operating system services. For example, when a program issued a print command (called an input/output command or an I/O command), it generated an interrupt requesting the services of the I/O processor and the CPU was released to begin execution of the next job. This was called *passive multiprogramming* because



the operating system didn't control the interrupts but waited for each job to end an execution sequence. It was less than ideal because if a job was CPU-bound (meaning that it performed a great deal of nonstop CPU processing before issuing an interrupt), it could tie up the CPU for a long time while all other jobs had to wait.

To counteract this effect, the operating system was soon given a more active role with the advent of *active multiprogramming*, which allowed each program to use only a preset slice of CPU time, which is discussed in Chapter 4. When time expired, the job was interrupted and another job was allowed to begin execution. The interrupted job had to wait until it was allowed to resume execution later. The idea of time slicing soon became common in many time-sharing systems.

Program scheduling, which was begun with second-generation systems, continued at this time but was complicated by the fact that main memory was occupied by many jobs. To solve this problem, the jobs were sorted into groups and then loaded into memory according to a preset rotation formula. The sorting was often determined by priority or memory requirements—whichever was found to be the most efficient use of the available resources. In addition to scheduling jobs, handling interrupts, and allocating memory, the operating systems also had to resolve conflicts whenever two jobs requested the same device at the same time, something we will explore in Chapter 5.

Even though there was progress in processor management, few major advances were made in data management.

## 1970s

After the third generation, during the late 1970s, computers had faster CPUs, creating an even greater disparity between their rapid processing speed and slower I/O access time. The first Cray supercomputer was released in 1976. Multiprogramming schemes to increase CPU use were limited by the physical capacity of the main memory, which was a limited resource and very expensive.

A solution to this physical limitation was the development of virtual memory, which took advantage of the fact that the CPU could process only one instruction at a time. With virtual memory, the entire program didn't need to reside in memory before execution could begin. A system with virtual memory would divide the programs into parts and keep them in secondary storage, bringing each part into memory only as it was needed. (Programmers of second-generation computers had used this concept with the roll in/roll out programming method, also called overlay, to execute programs that exceeded the physical memory of those computers.)

At this time there was also growing attention to the need for data resource conservation. Database management software became a popular tool because it organized data in an integrated manner, minimized redundancy, and simplified updating and

access of data. A number of query systems were introduced that allowed even the novice user to retrieve specific pieces of the database. These queries were usually made via a terminal, which in turn mandated a growth in terminal support and data communication software.

Programmers soon became more removed from the intricacies of the computer, and application programs started using English-like words, modular structures, and standard operations. This trend toward the use of standards improved program management because program maintenance became faster and easier.

## 1980s

Development in the 1980s dramatically improved the cost/performance ratio of computer components. Hardware was more flexible, with logical functions built on easily replaceable circuit boards. And because it was less costly to create these circuit boards, more operating system functions were made part of the hardware itself, giving rise to a new concept—**firmware**, a word used to indicate that a program is permanently held in read-only memory (ROM), as opposed to being held in secondary storage. The job of the programmer, as it had been defined in previous years, changed dramatically because many programming functions were being carried out by the system's software, hence making the programmer's task simpler and less hardware dependent.

Eventually the industry moved to **multiprocessing** (having more than one processor), and more complex languages were designed to coordinate the activities of the multiple processors servicing a single job. As a result, it became possible to execute programs in parallel, and eventually operating systems for computers of every size were routinely expected to accommodate multiprocessing.

The evolution of personal computers and high-speed communications sparked the move to networked systems and distributed processing, enabling users in remote locations to share hardware and software resources. These systems required a new kind of operating system—one capable of managing multiple sets of subsystem managers, as well as hardware that might reside half a world away.

With network operating systems, users generally became aware of the existence of many networked resources, could log in to remote locations, and could manipulate files on networked computers distributed over a wide geographical area. Network operating systems were similar to single-processor operating systems in that each machine ran its own local operating system and had its own users. The difference was in the addition of a network interface controller with low-level software to drive the local operating system, as well as programs to allow remote login and remote file access. Still, even with these additions, the basic structure of the network operating system was quite close to that of a standalone system.

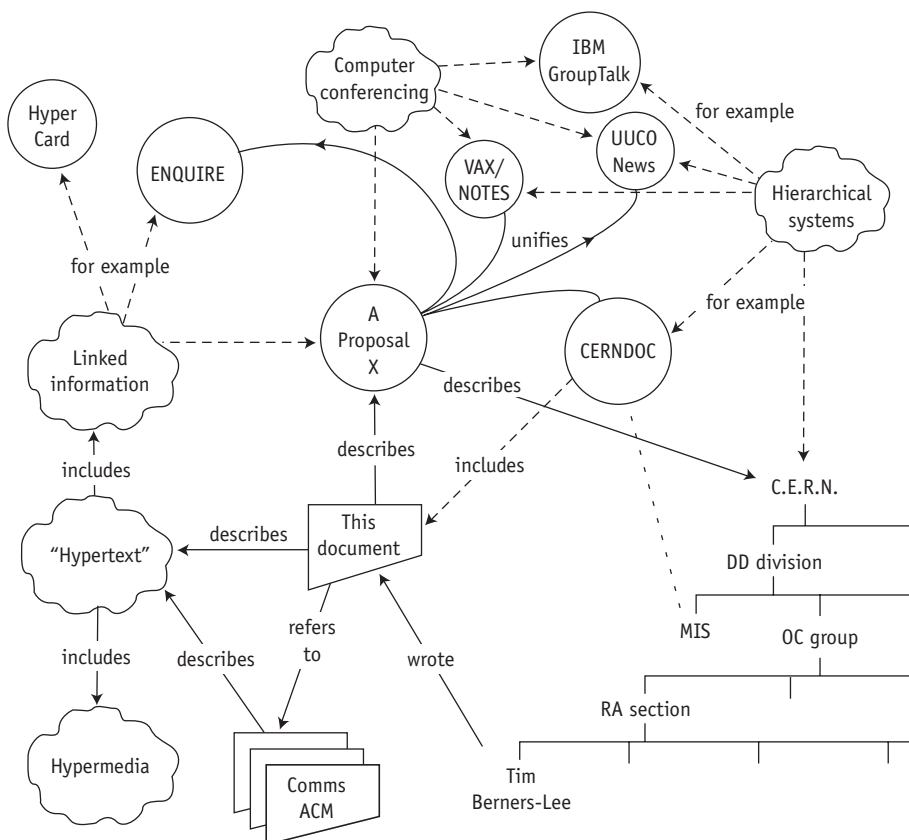
On the other hand, with distributed operating systems, users could think they were working with a typical uniprocessor system when in fact they were connected to a cluster of many processors working closely together. With these systems, users didn't need to know which processor was running their applications or which devices were storing their files. These details were all handled transparently by the operating system—something that required more than just adding a few lines of code to a uniprocessor operating system. The disadvantage of such a complex operating system was the requirement for more complex processor-scheduling algorithms. In addition, communications delays within the network sometimes meant that scheduling algorithms had to operate with incomplete or outdated information.

## 1990s

The overwhelming demand for Internet capability in the mid-1990s sparked the proliferation of networking capability. The World Wide Web, conceived in a paper, shown in Figure 1.12, by Tim Berners-Lee made the Internet accessible by computer users

(figure 1.12)

*Illustration from the first page of the 1989 proposal by Tim Berners-Lee describing his revolutionary "linked information system." Based on this research, he designed the first World Wide Web server and browser, making it available to the general public in 1991.*



worldwide, not just the researchers who had come to depend on it for global communications. Web accessibility and e-mail became standard features of almost every operating system. However, increased networking also sparked increased demand for tighter security to protect hardware and software.

The decade also introduced a proliferation of multimedia applications demanding additional power, flexibility, and device compatibility for most operating systems. A typical multimedia computer houses devices to perform audio, video, and graphic creation and editing. Those functions can require many specialized devices such as a microphone, digital piano, Musical Instrument Digital Interface (MIDI), digital camera, digital video disc (DVD) drive, optical disc (CD) drives, speakers, additional monitors, projection devices, color printers, and high-speed Internet connections. These computers also require specialized hardware (such as controllers, cards, busses) and software to make them work together properly.

Multimedia applications need large amounts of storage capability that must be managed gracefully by the operating system. For example, each second of a 30-frame-per-minute full-screen video requires 27MB of storage unless the data is compressed in some way. To meet the demand for compressed video, special-purpose chips and video boards have been developed by hardware companies.

What's the effect of these technological advances on the operating system? Each advance requires a parallel advance in the software's management capabilities.

## 2000S

The new century emphasized the need for operating systems to offer improved flexibility, reliability, and speed. To meet the need for computers that could accommodate multiple operating systems running at the same time and sharing resources, the concept of virtual machines, shown in Figure 1.13, was developed and became commercially viable.

**Virtualization** is the creation of partitions on a single server, with each partition supporting a different operating system. In other words, it turns a single physical server into multiple virtual servers, often with multiple operating systems. Virtualization requires the operating system to have an intermediate manager to oversee each operating system's access to the server's physical resources. For example, with virtualization, a single processor can run 64 independent operating systems on workstations using a processor capable of allowing 64 separate threads (instruction sequences) to run at the same time.

(figure 1.13)

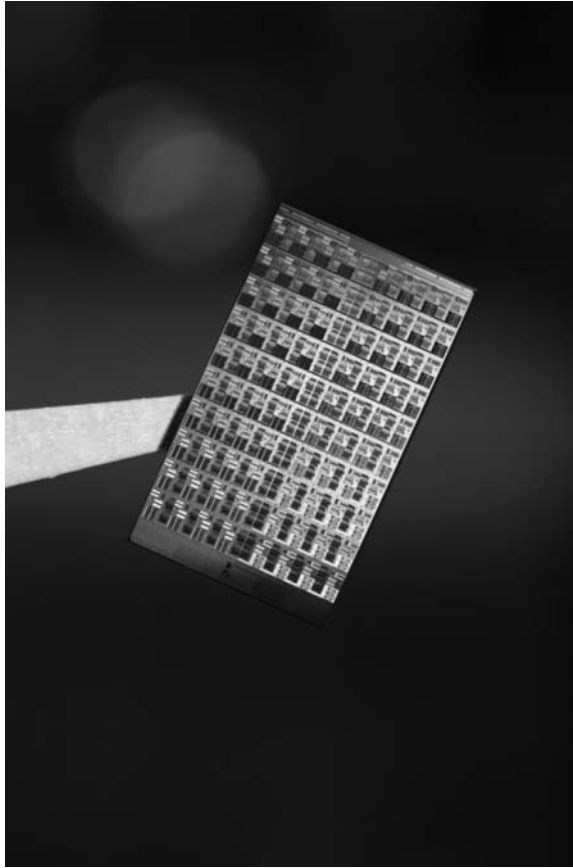
With virtualization, different operating systems can run on a single computer.

Courtesy of Parallels, Inc.



Processing speed has enjoyed a similar advancement with the development of multi-core processors, shown in Figure 1.14. Until recent years, the silicon wafer that forms the base of the computer chip circuitry held only a single CPU. However, with the introduction of dual-core processors, a single chip can hold multiple processor **cores**. Thus, a dual-core chip allows two sets of calculations to run at the same time, which sometimes leads to faster completion of the job. It's as if the user has two separate computers, and two processors, cooperating on a single task. As of this writing, designers have created chips that can hold 80 simple cores.

Does this hardware innovation affect the operating system software? Absolutely, because it must now manage the work of these multiple processors and be able to schedule and manage the processing of their multiple tasks. We'll explore some of the complexities of this in Chapter 6.



**(Figure 1.14)**

*A single piece of silicon can hold 80 cores, which (to put it in simplest terms) can perform 80 calculations at one time.*

*Courtesy of Intel Corporation*

## Threads

Multi-core technology helps the operating system handle threads, multiple actions that can be executed at the same time. First, an explanation: The Processor Manager is responsible for processing each job submitted by a user. Jobs are made up of processes (sometimes called tasks in other textbooks), and processes consist of multiple threads.

A process has two characteristics:

- It requires space in main memory where it resides during its execution; although, from time to time, it requires other resources such as data files or I/O devices.
- It passes through several states (such as running, waiting, ready) from its initial arrival into the computer system to its completion.

Multiprogramming and virtual memory dictate that processes be swapped between main memory and secondary storage during their execution. With conventional processes (also known as heavyweight processes), this swapping results in a lot of

overhead. That's because each time a swap takes place, all process information must be saved to preserve the process's integrity.

A **thread** (or lightweight process) can be defined as a unit smaller than a process, which can be scheduled and executed. Using this technique, the heavyweight process, which owns the resources, becomes a more passive element, while a thread becomes the element that uses the CPU and is scheduled for execution. Manipulating threads is less time consuming than manipulating processes, which are more complex. Some operating systems support multiple processes with a single thread, while others support multiple processes with multiple threads.



Web browsers routinely use multithreading to allow users to explore multiple areas of interest on the Internet at the same time.

Multithreading allows applications to manage a separate process with several threads of control. Web browsers use multithreading routinely. For instance, one thread can retrieve images while another sends and retrieves e-mail. Multithreading is also used to increase responsiveness in a time-sharing system to increase resource sharing and decrease overhead.

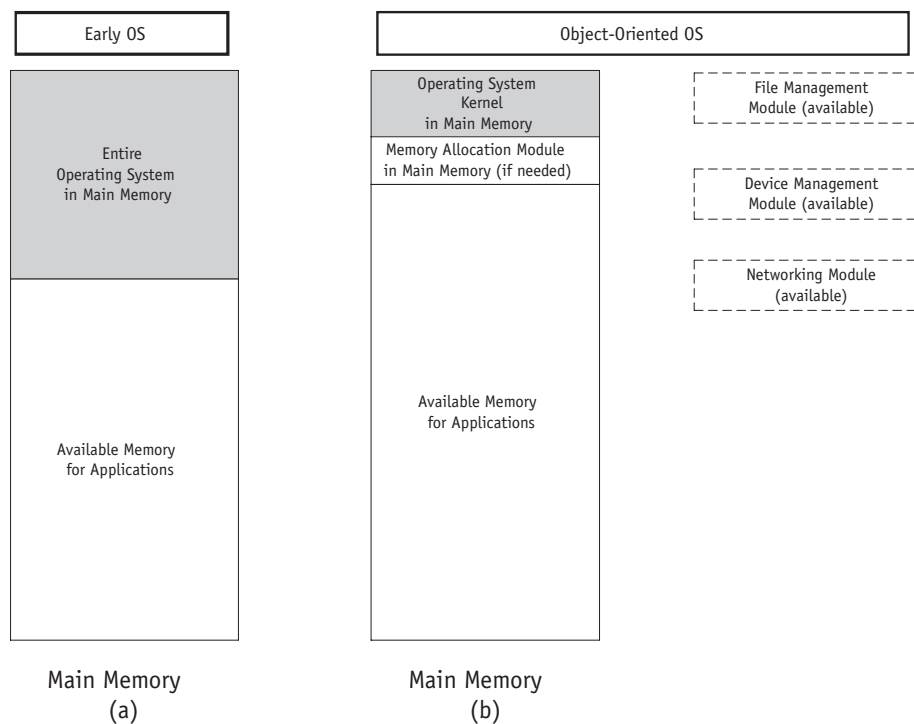
## Object-Oriented Design

An important area of research that resulted in substantial efficiencies was that of the system architecture of operating systems—the way their components are programmed and organized, specifically the use of **object-oriented** design and the reorganization of the operating system's nucleus, the kernel. The kernel is the part of the operating system that resides in memory at all times, performs the most essential operating system tasks, and is protected by hardware from user tampering.

The first operating systems were designed as a comprehensive single unit, as shown in Figure 1.15 (a). They stored all required elements of the operating system in memory such as memory allocation, process scheduling, device allocation, and file management. This type of architecture made it cumbersome and time consuming for programmers to add new components to the operating system, or to modify existing ones.

Most recently, the part of the operating system that resides in memory has been limited to a few essential functions, such as process scheduling and memory allocation, while all other functions, such as device allocation, are provided by special modules, which are treated as regular applications, as shown in Figure 1.15 (b). This approach makes it easier to add new components or modify existing ones.

Object-oriented design was the driving force behind this new organization. Objects are self-contained modules (units of software) that provide models of the real world and can be reused in different applications. By working on objects, programmers can modify and customize pieces of an operating system without disrupting the integrity of the remainder of the system. In addition, using a modular, object-oriented approach can

**(figure 1.15)**

*Early operating systems (a) loaded in their entirety into main memory.*

*Object-oriented operating systems (b) load only the critical elements into main memory and call other objects as needed.*

make software development groups more productive than was possible with procedural structured programming.

## Conclusion

In this chapter, we looked at the overall function of operating systems and how they have evolved to run increasingly complex computers and computer systems; but like any complex subject, there's much more detail to explore. As we'll see in the remainder of this text, there are many ways to perform every task and it's up to the designer of the operating system to choose the policies that best match the system's environment.

In the following chapters, we'll explore in detail how each portion of the operating system works, as well as its features, functions, benefits, and costs. We'll begin with the part of the operating system that's the heart of every computer: the module that manages main memory.



## Key Terms

**batch system:** a type of system developed for the earliest computers that used punched cards or tape for input, which were entered in a batch.

**central processing unit (CPU):** the component with the circuitry, the “chips,” to control the interpretation and execution of instructions.

**core:** the processing part of a CPU chip made up of the control unit and the arithmetic logic unit (ALU).

**Device Manager:** the section of the operating system responsible for controlling the use of devices. It monitors every device, channel, and control unit and chooses the most efficient way to allocate all of the system’s devices.

**embedded system:** a dedicated computer system, often small and fast, that resides in a larger physical system such as jet aircraft or ships.

**File Manager:** the section of the operating system responsible for controlling the use of files.

**firmware:** software instructions or data that are stored in a fixed or “firm” way, usually implemented on read-only memory (ROM).

**hardware:** the physical machine and its components, including main memory, I/O devices, I/O channels, direct access storage devices, and the central processing unit.

**hybrid system:** a computer system that supports both batch and interactive processes.

**interactive system:** a system that allows each user to interact directly with the operating system via commands entered from a keyboard.

**kernel:** the primary part of the operating system that remains in random access memory (RAM) and is charged with performing the system’s most essential tasks, such as managing main memory and disk access.

**main memory:** the memory unit that works directly with the CPU and in which the data and instructions must reside in order to be processed. Also called primary storage or internal memory.

**mainframe:** the historical name given to a large computer system characterized by its large size, high cost, and high performance.

**Memory Manager:** the section of the operating system responsible for controlling the use of memory. It checks the validity of each request for memory space and, if it’s a legal request, allocates the amount needed to execute the job.

**microcomputer:** a small computer equipped with all the hardware and software necessary to perform one or more tasks.

**minicomputer:** a small to medium-sized computer system, also called a midrange computer.

**multiprocessing:** when two or more CPUs share the same main memory, most I/O devices, and the same control program routines. They service the same job stream and execute distinct processing programs concurrently.

**multiprogramming:** a technique that allows a single processor to process several programs residing simultaneously in main memory and interleaving their execution by overlapping I/O requests with CPU requests.

**network:** a system of interconnected computer systems and peripheral devices that exchange information with one another.

**Network Manager:** the section of the operating system responsible for controlling access to and the use of networked resources.

**object-oriented:** a programming philosophy whereby programs consist of self-contained, reusable modules called objects, each of which supports a specific function, but which are categorized into classes of objects that share the same function.

**operating system:** the software that manages all the resources of a computer system.

**Processor Manager:** a composite of two submanagers, the Job Scheduler and the Process Scheduler, which decides how to allocate the CPU.

**real-time system:** a computing system used in time-critical environments that require guaranteed response times, such as navigation systems, rapid transit systems, and industrial control systems.

**server:** a node that provides to clients various network services, such as file retrieval, printing, or database access services.

**software:** a collection of programs used to perform certain tasks. Software falls into three main categories: operating system programs, compilers and assemblers, and application programs.

**storage:** a place where data is stored in the computer system. Primary storage is main memory and secondary storage is nonvolatile media.

**supercomputer:** the fastest, most sophisticated computers made, used for complex calculations.

**thread:** a portion of a program that can run independently of other portions. Multithreaded application programs can have several threads running at one time with the same or different priorities.

**throughput:** a composite measure of a system's efficiency that counts the number of jobs served in a given unit of time.