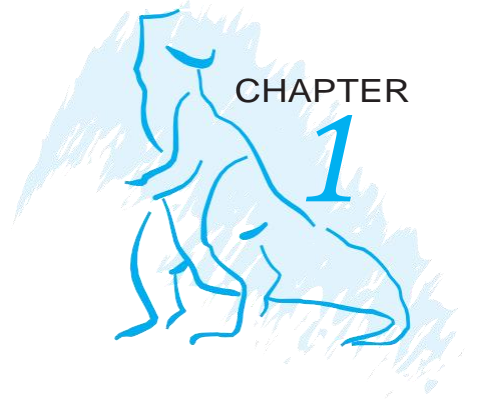# OPERATING SYSTEM

# *Introduction*

An **operating system** is a program that manages a computer's hardware. It also provides a basis for application programs and acts as an intermediary between the computer user and the computer hardware.

Before we can explore the details of computer system operation, we need to know something about system structure. We thus discuss the basic functions of system startup, I/O, and storage early in this chapter. We also describe the basic computer architecture that makes it possible to write a functional operating system.

Because an operating system is large and complex, it must be created piece by piece. Each of these pieces should be a well-delineated portion of the system, with carefully defined inputs, outputs, and functions. In this chapter, we provide a general overview of the major components of a contemporary computer system as well as the functions provided by the operating system. Additionally, we cover several other topics to help set the stage for the remainder of this text: data structures used in operating systems, computing environments, and open-source operating systems.

## 1.1 What Operating Systems Do?

We begin our discussion by looking at the operating system's role in the overall computer system. A computer system can be divided roughly into **four components**: the hardware, the operating system, the application programs, and the users (Figure 1.1).

The hardware —the central processing unit (CPU), the memory, and the input/output (I/O) devices —provides the basic computing resources for the system. The application programs —such as word processors, spreadsheets, compilers, and Web browsers— define the ways in which these resources are used to solve users' computing problems. **The operating system controls the hardware and coordinates its use among the various application programs for the various users.**
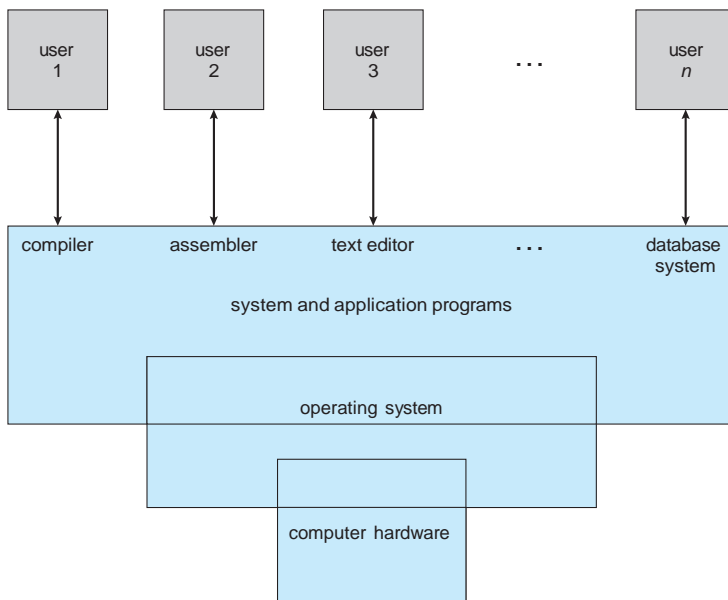
**Figure 1.1** Abstract view of the components of a computer

- OS is a resource allocator Manages all resources Decides between conflicting requests for efficient and fair resource use

- OS is a control program Controls execution of programs to prevent errors and improper use of the computer

- Kernel: The one program running at all times on the computer.

### 1.1.1 User View

The user's view of a computer depends on the type of interface being used. For personal computers (PCs), the system is designed for a single user, focusing on ease of use and performance, with little concern for resource sharing. In contrast, mainframes or minicomputers serve multiple users simultaneously, with an operating system optimized for efficient resource utilization. Workstations connected to networks balance individual usability with shared resource management, such as file and print servers. Mobile devices like smartphones and tablets, typically used for tasks like email and web browsing, offer touch-based interfaces and are increasingly replacing traditional computers. Lastly, embedded systems in devices like appliances and cars are designed to function with minimal user interaction.

**Users want convenience, ease of use and good performance**

-Don't care about resource utilization

• But shared computers such as mainframe or minicomputer must keep all users happy.

### 1.1.2 System View

From the system's perspective, the operating system (OS) is the key program that interacts with hardware, acting as a **resource allocator**. It manages resources like CPU time, memory, file storage, and I/O devices, deciding how to distribute them efficiently and fairly among programs and users, especially in systems with multiple users.

Additionally, the OS functions as a control program, overseeing the execution of user programs, preventing errors, and ensuring proper use of the system, with a focus on controlling I/O devices.

## 1.2 Computer-System Organization

Before we can explore the details of how computer systems operate, we need general knowledge of the structure of a computer system. In this section, we look at several parts of this structure. The section is mostly concerned with computer-system organization.

A modern general-purpose computer system consists of one or more CPUs and a number of device controllers connected through a common bus that provides access to shared memory (Figure 1.2).
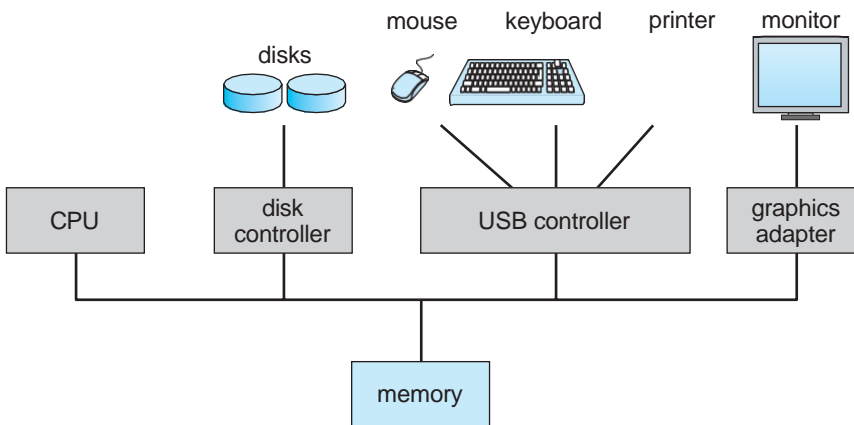


**Figure 1.2** A modern computer system.

### 1.2.1 Computer Startup

- The **bootstrap program** is loaded at power-up or reboot.
- It is **typically stored in ROM or EPROM**, commonly known as firmware.
- The bootstrap program **initializes all aspects of the system**.
- It then **loads the operating system kernel** and starts its execution.

This program is responsible for starting up the system and ensuring that the essential components are ready for the operating system to run.

### 1.2.2 Storage Structure

The basic unit of computer storage is the bit. A bit can contain one of two values, 0 and 1. All other storage in a computer is based on collections of bits. Given enough bits, it is amazing how many things a computer can represent: numbers, letters, images, movies, sounds, documents, and programs, to name a few. A byte is 8 bits, and on most computers it is the smallest convenient chunk of storage. For example, most computers don't have an instruction to move a bit but do have one to move a byte. A less common term is word, which is a given computer architecture's native unit of data. A word is made up of one or more bytes. For example, a computer that has 64-bit registers and 64-bit memory addressing typically has 64-bit (8-byte) words. A computer executes many operations in its native word size rather than a byte at a time. Computer storage, along with most computer throughput, is generally measured and manipulated in bytes and collections of bytes.

**A kilobyte, or KB, is 1,024 bytes**

**a megabyte, or MB, is $1,024^2$ bytes**

**a gigabyte, or GB, is 1,024$^3$ bytes**

 **a terabyte, or TB, is 1,024$^4$ bytes**

 **a petabyte, or PB, is 1,024$^5$ bytes**

Computer manufacturers often round off these numbers and say that a megabyte is 1 million bytes and a gigabyte is 1 billion bytes. Networking measurements are an exception to this general rule; they are given in bits (because networks move data a bit at a time).
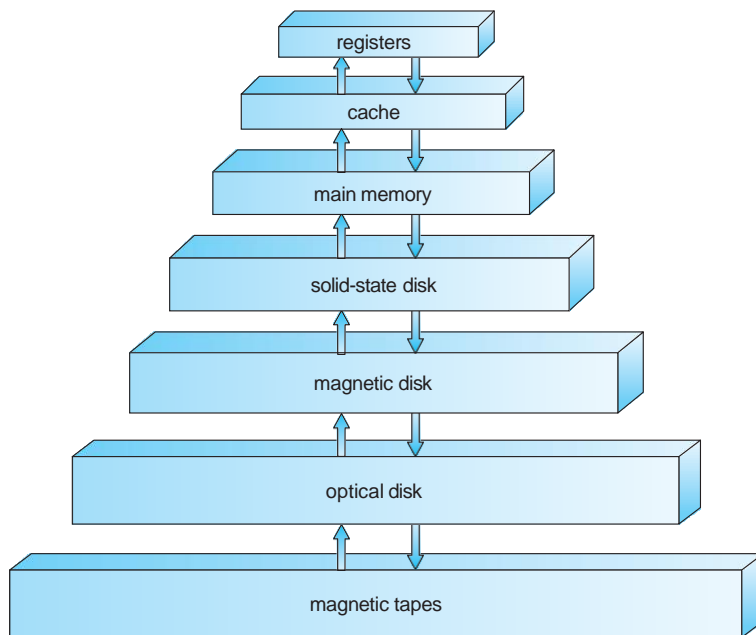


**Figure 1.4**  Storage-device hierarchy.

Main memory – only large storage media that the CPU can access directly

– Random access

– Typically volatile

• Secondary storage – extension of main memory that provides large nonvolatile storage capacity

• Hard disks – rigid metal or glass platters covered with magnetic recording material

– Disk surface is logically divided into tracks, which are subdivided into sectors

– The disk controller determines the logical interaction between the device and the computer

• Solid-state disks – faster than hard disks, nonvolatile

– Various technologies

– Becoming more popular

## 1.3 Computer-System Architecture

A computer system can be organized in a number of different ways, which we can categorize roughly according to the number of general-purpose processors used.
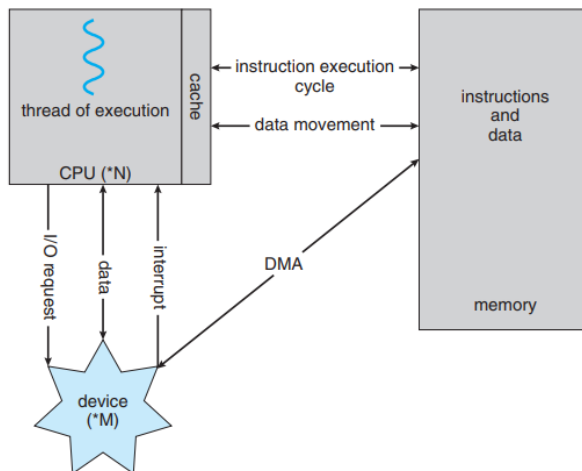


**Figure 1.5**  How a modern computer system works.

## 1.3. Single-processor systems

Single-processor systems feature one main CPU that executes general-purpose instructions and user processes. Though they may include special-purpose processors like disk or graphics controllers, these processors handle specific tasks with limited instruction sets and do not run user processes. They often assist the CPU by offloading tasks like disk scheduling or keyboard input processing. However, despite having additional specialized processors, the system remains classified as a single-processor system if there is only one general-purpose CPU handling the main workload.

## 1.3.2 Multiprocessor Systems

Within the past several years, **multiprocessor systems** (also known as **parallel systems** or **multicore systems**) have begun to dominate the landscape of computing. Such systems have two or more processors in close communication, sharing the computer bus and sometimes the clock, memory, and peripheral devices. Multiprocessor systems first appeared prominently appeared in servers and have since migrated to desktop and laptop systems. Recently, multiple processors have appeared on mobile devices such as smartphones and tablet computers.

Multiprocessor systems offer three main advantages:

1. **Increased throughput**: More processors mean more work can be done in less time, though overhead and resource contention limit the speed-up.
2. **Economy of scale**: They are more cost-efficient than multiple single-processor systems, as they can share peripherals, storage, and power supplies.
3. **Increased reliability**: If one processor fails, the others can take over its workload, ensuring the system slows down but doesn't completely fail, providing higher fault tolerance.

Multiprocessor systems can be classified into two types:

1. **Asymmetric multiprocessing**: One "boss" processor controls the system and assigns tasks to "worker" processors. Workers either follow instructions or have specific predefined tasks.
2. **Symmetric multiprocessing (SMP)**: All processors are peers, capable of performing all tasks. Each processor has its own local cache and shares memory, allowing efficient multitasking. Modern operating systems (Windows, Mac OS X, Linux) support SMP.

Additionally, multicore systems, with multiple cores on a single chip, offer better efficiency, reduced power consumption, and improved performance.
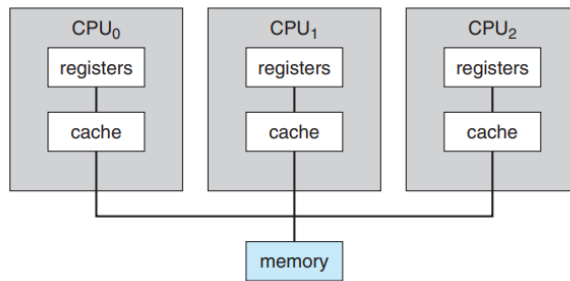


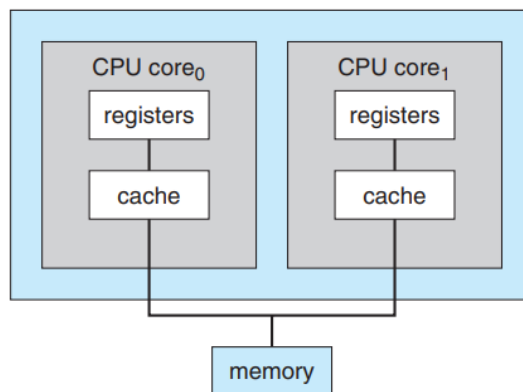**Figure 1.6** Symmetric multiprocessing architecture.



**Figure 1.7** A dual-core design with two cores placed on the same chip.
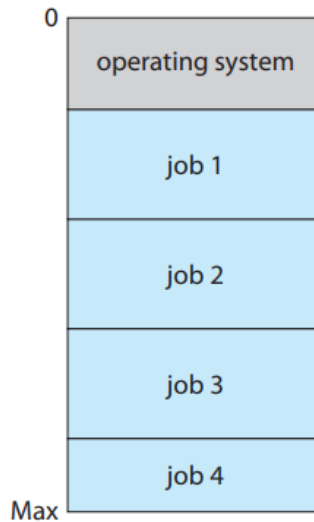
# 1.4 Operating-System Structure



**Figure 1.8** Memory layout for a multiprogramming system.

## 1.4.1 Operating-System Structure

An operating system provides the environment for executing programs and is structured to manage resources efficiently. A key feature is **multiprogramming**, where multiple programs run simultaneously, enhancing **CPU utilization**. Several jobs are kept in memory, with the remainder stored on disk in a **job pool**. The OS switches between jobs, ensuring that the CPU is always working, even when a program is waiting for I/O operations. This method ensures that the CPU is never idle, maximizing system efficiency.

 - Time-sharing systems, an extension of multiprogramming, allow multiple users to interact with programs simultaneously by rapidly switching between processes. This provides each user with the impression of sole access to the system. Time-sharing requires **CPU scheduling** and **memory management** to maintain efficiency, especially for interactive I/O tasks, which run at human speeds. **Virtual memory** allows running larger programs than available

physical memory. Additionally, these systems ensure file management, resource protection, and handle job synchronization to avoid **deadlocks** and ensure smooth operation across multiple processes.

## 1.5 Operating-System Operations

Modern operating systems are interrupt-driven, meaning they respond to events like I/O requests or errors through interrupts or traps. A trap is a software-generated interrupt caused by errors (e.g., division by zero) or specific requests for OS services. Each interrupt triggers a service routine to manage the event. The OS ensures errors in one program don't affect others, providing necessary protection. Without such protection, bugs could corrupt other processes or the OS itself, potentially causing system-wide issues or incorrect execution of multiple programs.

### 1.5.1 Dual-Mode and Multimode Operation

In order to ensure the proper execution of the operating system, we must be able to distinguish between the execution of operating-system code and user- defined code. The approach taken by most computer systems is to provide hardware support that allows us to differentiate among various modes of execution.
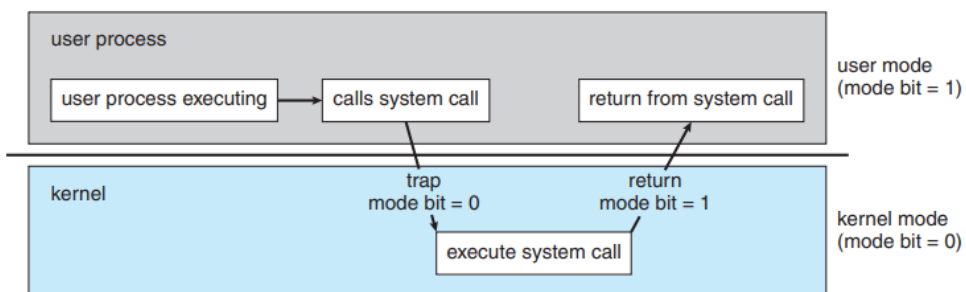
**Figure 1.9** Transition from user to kernel mode

At the very least, we need two separate *modes* of operation: **user mode** and **kernel mode** (also called **supervisor mode**, **system mode**, or **privileged mode**). A bit, called the **mode bit**, is added to the hardware of the computer to indicate the current mode: kernel (0) or user (1). With the mode bit, we can distinguish between a task that is executed on behalf of the operating system and one that is executed on behalf of the user. When the computer system is executing on behalf of a user application, the system is in user mode. However, when a user application requests a service from the operating system (via a system call), the system must transition from user to kernel mode to fulfill the request. This is shown in Figure 1.9. As we shall see, this architectural enhancement is useful for many other aspects of system operation as well.

At system boot time, the hardware starts in kernel mode. The operating system is then loaded and starts user applications in user mode. Whenever a trap or interrupt occurs, the hardware switches from user mode to kernel mode (that is, changes the state of the mode bit to 0). Thus, whenever the operating system gains control of the computer, it is in kernel mode. The system always switches to user mode (by setting the mode bit to 1) before passing control to a user program.

## 1.5 Storage Management

To make the computer system convenient for users, the operating system provides a uniform, logical view of information storage. The operating system abstracts from the physical properties of its storage devices to define a logical storage unit, the **file**. The operating system maps files onto physical media and accesses these files via the storage devices.

### 1.5.1 File-System Management

File management is a key function of an operating system, involving the storage and organization of data across different physical media like magnetic disks, optical disks, and tapes. Each medium has unique characteristics such as

speed, capacity, and access methods. Files, which can be programs or data, are organized into directories to enhance usability. Operating systems also manage access control, determining who can read, write, or modify files. This ensures efficient handling of storage devices and protection of shared resources.

## 1.5.2 Mass-Storage Management

As we have already seen, because main memory is too small to accommodate all data and programs, and because the data that it holds are lost when power is lost, the computer system must provide secondary storage to back up main memory. Most modern computer systems use disks as the principal on-line storage medium for both programs and data. Most programs— including compilers, assemblers, word processors, editors, and formatters— are stored on a disk until loaded into memory. They then use the disk as both the source and destination of their processing. Hence, the proper management of disk storage is of central importance to a computer system. The operating system is responsible for the following activities in connection with disk management:

- Free-space management

- Storage allocation

- Disk scheduling

Because secondary storage is used frequently, it must be used efficiently. The entire speed of operation of a computer may hinge on the speeds of the disk subsystem and the algorithms that manipulate that subsystem.

There are, however, many uses for storage that is slower and lower in cost (and sometimes of higher capacity) than secondary storage. Backups of disk data, storage of seldom-used data, and long-term archival storage are some examples. Magnetic tape drives and their tapes and CD and DVD drives

and platters are typical **tertiary storage** devices. The media (tapes and optical platters) vary between **WORM** (write-once, read-many-times) and **RW** (read– write) formats.

Tertiary storage is not crucial to system performance, but it still must be managed. Some operating systems take on this task, while others leave tertiary-storage management to application programs. Some of the functions that operating systems can provide include mounting and unmounting media in devices, allocating and freeing the devices for exclusive use by processes, and migrating data from secondary to tertiary storage.

| Level | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| Name | registers | cache | main memory | solid state disk | magnetic disk |
| Typical size | < 1 KB | < 16MB | < 64GB | < 1 TB | < 10 TB |
| Implementation technology | custom memory with multiple ports CMOS | on-chip or off-chip CMOS SRAM | CMOS SRAM | flash memory | magnetic disk |
| Access time (ns) | 0.25 - 0.5 | 0.5 - 25 | 80 - 250 | 25,000 - 50,000 | 5,000,000 |
| Bandwidth (MB/sec) | 20,000 - 100,000 | 5,000 - 10,000 | 1,000 - 5,000 | 500 | 20 - 150 |
| Managed by | compiler | hardware | operating system | operating system | operating system |
| Backed by | cache | main memory | disk | disk | disk or tape |

**Figure 1.11** Performance of various levels of storage.

## 1.6 Protection and Security

Protection and security in operating systems ensure that only authorized users can access important resources like files, memory, and the CPU. Protection keeps systems safe from errors or misuse by controlling who can use certain parts of the computer. Security helps defend against harmful things like viruses or hackers. Each user has a unique ID, and permissions can be set for individuals or groups. Sometimes, a user might need extra privileges temporarily,

like in UNIX, where certain programs can run with special permissions.

## Review Questions

1. **What is the primary role of an operating system in a computer system?**
2. **List the four main components of a computer system as described in the text.**
3. **How does an operating system manage resources in a computer?**
4. **What are the differences between the user view and system view of a computer?**

5. **Explain how the operating system acts as a resource allocator.**
6. **What is the purpose of the kernel in an operating system?**
7. **Describe the function of the bootstrap program during system startup.**

# *Operating - System Structures*

An operating system provides the environment within which programs are executed. Internally, operating systems vary greatly in their makeup, since they are organized along many different lines. The design of a new operating system is a major task. It is important that the goals of the system be well defined before the design begins. These goals form the basis for choices among various algorithms and strategies.

We can view an operating system from several vantage points. One view focuses on the services that the system provides; another, on the interface that it makes available to users and programmers; a third, on its components and their interconnections. In this chapter, we explore all three aspects of operating systems, showing the viewpoints of users, programmers, and operating system designers. We consider what services an operating system provides, how they are provided, how they are debugged, and what the various methodologies are for designing such systems. Finally, we describe how operating systems are created and how a computer starts its operating system.

# 2.1 Operating-System Services

An operating system provides an environment for the execution of programs. It provides certain services to programs and to the users of those programs.
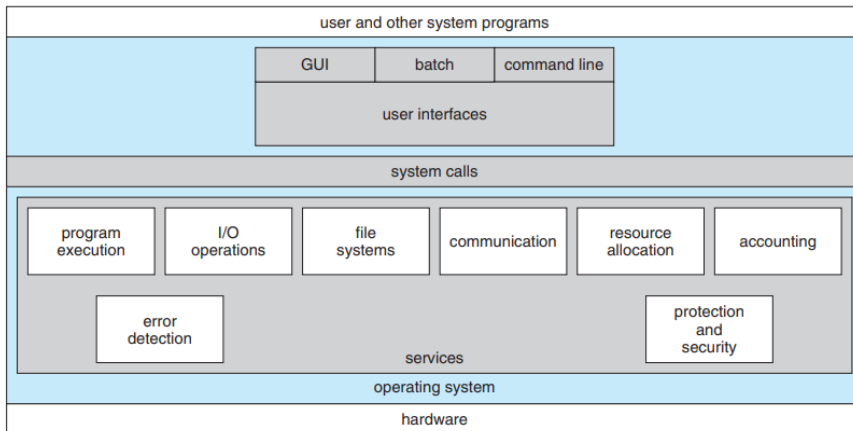


**Figure 2.1**   A view of operating system services.

Figure 2.1 shows one view of the various operating-system services and how they interrelate.

## Operating system services :

- **User interface**. Almost all operating systems have a **user interface** (**UI**). This interface can take several forms. One is a **command-line interface** (**CLI**), which uses text commands and a method for entering them (say, a keyboard for typing in commands in a specific format with specific options). Another is a **batch interface**, in which commands and directives to control those commands are entered into files, and those files are executed. Most commonly, a **graphical user interface** (**GUI**) is used. Here, the interface is a window system with a pointing device to direct I/O, choose from menus, and make selections and a keyboard to enter text. Some systems provide two or all three of these variations.

- **Program execution**. The system must be able to load a program into memory and to run that program. The program must be able to end its execution, either normally or abnormally (indicating error).

- **I/O operations**.A running program may require I/O, which may involve a file or an I/O device. For specific devices, special functions may be desired (such as recording to a CD or DVD drive or blanking a display screen). For efficiency and protection, users usually cannot control I/O devices directly. Therefore, the operating system must provide a means to do I/O.

- **File-system manipulation**. The file system is of particular interest. Obvi- ously, programs need to read and write files and directories. They also need to create and delete them by name, search for a given file, and list file information. Finally, some operating systems include permissions management to allow or deny access to files or directories based on file ownership. Many operating systems provide a variety of file systems, sometimes to allow personal choice and sometimes to provide specific features or performance characteristics.

- **Communications**. There are many circumstances in which one process needs to exchange information with another process. Such communication may occur between processes that are executing on the same computer or between processes that are executing on different computer systems tied together by a computer network. Communications may be implemented via **shared memory**, in which two or more processes read and write to a shared section of memory, or **message passing**, in which packets of information in predefined formats are moved between

processes by the operating system.

- **Error detection**. The operating system needs to be detecting and correcting errors constantly. Errors may occur in the CPU and memory hardware (such as a memory error or a power failure), in I/O devices (such as a parity error on disk, a connection failure on a network, or lack of paper in the printer), and in the user program (such as an arithmetic overflow, an attempt to access an illegal memory location, or a too-great use of CPU time). For each type of error, the operating system should take the appropriate action to ensure correct and consistent computing. Sometimes, it has no choice but to halt the system. At other times, it might terminate an error-causing process or return an error code to a process for the process to detect and possibly correct.

Another set of operating system functions exists not for helping the user but rather for ensuring the efficient operation of the system itself. Systems with multiple users can gain efficiency by sharing the computer resources among the users.

- **Resource allocation**. When there are multiple users or multiple jobs running at the same time, resources must be allocated to each of them. The operating system manages many different types of resources. Some (such as CPU cycles, main memory, and file storage) may have special allocation code, whereas others (such as I/O devices) may have much more general request and release code. For instance, in determining how best to use the CPU, operating systems have CPU-scheduling routines that take into account the speed of the CPU, the jobs that must be executed, the number of registers available, and other factors. There may also be routines to allocate printers, USB storage drives, and other peripheral devices.

- **Accounting**. We want to keep track of which users use how much and what kinds of computer resources. This record keeping may be used for accounting (so that users can be billed) or simply for accumulating usage statistics. Usage statistics may be a valuable tool for researchers who wish to reconfigure the system to improve computing services.

- **Protection and security**. The owners of information stored in a multiuser or networked computer system may want to control use of that information. When several separate processes execute concurrently, it should not be possible for one process to interfere with the others or with the operating system itself. Protection involves ensuring that all access to system resources is controlled. Security of the system from outsiders is also important.

## 2.2 User and Operating-System Interface

We mentioned earlier that there are several ways for users to

interface with the operating system. Here, we discuss two fundamental approaches. One provides a command-line interface, or **command interpreter**, that allows users to directly enter commands to be performed by the operating system. The other allows users to interface with the operating system via a graphical user interface, or GUI.

## 2.2.1 Command Interpreters

◎ Command interpreters, also known as shells, are programs responsible for interpreting and executing user commands.
◎ Some operating systems integrate the command interpreter into the kernel, while others, such as Windows and UNIX, treat it as a separate program that runs when a user logs in.
◎ UNIX and Linux systems offer multiple shells to choose from, such as Bourne, C shell, Korn, and others, allowing users to select one based on their preferences.

◎ The primary function of the shell is to execute user-specified commands like creating, deleting, or listing files, as seen in MS-DOS and UNIX shells.

Figure 2.2 shows the Bourne shell command interpreter being used on Solaris 10.

```
                              Terminal
 File  Edit  View  Terminal  Tabs  Help
fd0       0.0     0.0     0.0     0.0  0.0  0.0     0.0   0   0
sd0       0.0     0.2     0.0     0.2  0.0  0.0     0.4   0   0
sd1       0.0     0.0     0.0     0.0  0.0  0.0     0.0   0   0
                 extended device statistics
device    r/s     w/s     kr/s    kw/s wait actv  svc_t  %w  %b
fd0       0.0     0.0     0.0     0.0  0.0  0.0     0.0   0   0
sd0       0.6     0.0     38.4    0.0  0.0  0.0     8.2   0   0
sd1       0.0     0.0     0.0     0.0  0.0  0.0     0.0   0   0
(root@pbg-nv64-vm)-(11/pts)-(00:53 15-Jun-2007)-(global)
-(/var/tmp/system-contents/scripts)# swap -sh
total: 1.1G allocated + 190M reserved = 1.3G used, 1.6G available
(root@pbg-nv64-vm)-(12/pts)-(00:53 15-Jun-2007)-(global)
-(/var/tmp/system-contents/scripts)# uptime
 12:53am  up 9 min(s),  3 users,  load average: 33.29, 67.68, 36.81
(root@pbg-nv64-vm)-(13/pts)-(00:53 15-Jun-2007)-(global)
-(/var/tmp/system-contents/scripts)# w
  4:07pm  up 17 day(s), 15:24,  3 users,  load average: 0.09, 0.11, 8.66
User      tty           login@ idle   JCPU   PCPU  what
root      console       15Jun0718days      1            /usr/bin/ssh-agent -- /usr/bi
n/d
root      pts/3         15Jun07           18      4  w
root      pts/4         15Jun0718days              w
(root@pbg-nv64-vm)-(14/pts)-(16:07 02-Jul-2007)-(global)
-(/var/tmp/system-contents/scripts)#
```

**Figure 2.2**   The Bourne shell command interpreter in Solrais 10.

## 2.2.2 Graphical User Interfaces

◎ GUIs provide a user-friendly, graphical way to interact with the operating system, replacing text-based command interfaces.

◎ Users interact with GUIs through a mouse-based window-and-menu system, where icons represent files, programs, and system functions.

◎The first GUI appeared on the Xerox Alto in 1973, but GUIs became popular with the Apple Macintosh in the 1980s and later versions of Microsoft Windows.

◎ GUIs have evolved, particularly with the rise of mobile devices, where touchscreens have replaced mouse interactions, allowing users to swipe and press to interact with the system.

◎ UNIX systems traditionally used command-line interfaces, but GUIs like CDE, X-Windows, KDE, and GNOME have become common, especially in open-source environments like Linux.

Figure 2.3 illustrates the touchscreen of the Apple iPad.

Whereas earlier smartphones included a physical keyboard, most smartphones now simulate a keyboard on the touchscreen.



**Figure 2.3** The iPad touchscreen.

## 2.3 System Calls

**System calls** System calls provide a crucial interface between a user program and the operating system, enabling the program to request services like reading files or printing data. These calls are often written in high-level languages like C or C++, but certain low-level tasks may require assembly language.

**Example: File Copy Program**

To illustrate how system calls are used, consider writing a program that reads data from one file and copies it to another. This process involves several steps:

1. **Getting File Names:**

- The program needs input and output file names. In an interactive system, the program might prompt the user for these names through system calls, displaying a message on the screen and reading the user's input from the keyboard.
- In graphical systems, the file names could be selected using a mouse and a menu system, requiring many input/output (I/O) system calls.

2. **Opening Files:**

- The program opens the input file and creates the output file using system calls.
- If errors occur (e.g., the input file doesn't exist or is protected), the program needs additional system calls to handle the errors and possibly terminate abnormally.
- If the output file already exists, the program may use system calls to delete it, replace it, or ask the user what to do.

3. **Reading and Writing:**

- The program reads data from the input file and writes it to the output file. Each read and write operation is a system call.
- The program must handle potential errors, such as reaching the end of the file or encountering hardware failures during reading, and issues like disk space running out during writing.

4. **Closing Files and Exiting:**

After the file copy is complete, the program closes both files using system calls.
- Finally, it outputs a message and terminates normally using a final system call.

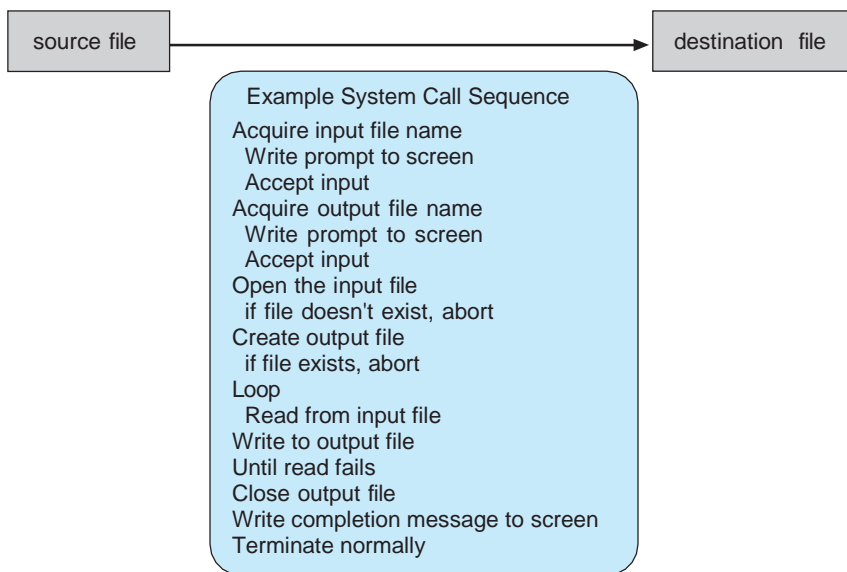 This system-call sequence is shown in Figure 2.5

**Figure 2.5** Example of how system calls are used.

Even simple programs rely heavily on system calls, with systems executing thousands of them per second. However, most programmers interact with an operating system through an **Application Programming Interface (API)** rather than directly invoking system calls.

1. **API vs. System Calls:**
    - APIs provide a higher-level, user-friendly interface for application programmers to interact with the system without dealing with the complexity of system calls.
    - APIs specify available functions, parameters, and expected return values, making programming easier.
2. **Common APIs:**
    - **Windows API**: Used for Windows-based systems.
    - **POSIX API**: Used for UNIX-based systems like Linux, Mac OS X, and others.

- o **Java API**: Used for programs running on the Java Virtual Machine (JVM).
3. **How APIs Work:**
    - o APIs call system functions indirectly. For example, the Windows `CreateProcess()` function calls the `NTCreateProcess()` system call in the Windows kernel.
4. **Benefits of Using APIs:**
    - o **Portability**: Programs written using an API can run on any system that supports the same API, though some architectural differences may still exist.
    - o **Simplification**: APIs abstract away the more complex, detailed system calls, making programming less error-prone and easier to manage.
5. **Correlation Between API and System Calls:**
    - o Many APIs, like those in POSIX and Windows, closely resemble the system calls in their respective operating systems, simplifying the development process.

## 2.4 Types of System Calls

System calls can be grouped roughly into six major categories: **process control**, **file manipulation**, **device manipulation**, **information maintenance**, **communications**, and **protection**. Most of these system calls support, or are supported by, concepts and functions that are discussed in later chapters. Figure 2.8 summarizes the types of system calls normally provided by an operating system. As mentioned, in this text, we normally refer to the system calls by generic names. Throughout the text, however, we provide examples of the actual counterparts to the system calls for Windows, UNIX, and Linux systems.

- Process control
  - end, abort
  - load, execute
  - create process, terminate process
  - get process attributes, set process attributes
  - wait for time
  - wait event, signal event
  - allocate and free memory
- File management
  - create file, delete file
  - open, close
  - read, write, reposition
  - get file attributes, set file attributes
- Device management
  - request device, release device
  - read, write, reposition
  - get device attributes, set device attributes
  - logically attach or detach devices
- Information maintenance
  - get time or date, set time or date
  - get system data, set system data
  - get process, file, or device attributes
  - set process, file, or device attributes
- Communications
  - create, delete communication connection
  - send, receive messages
  - transfer status information
  - attach or detach remote devices

Types of system calls.

## EXAMPLES OF WINDOWS AND UNIX SYSTEM CALLS

|  | Windows | Unix |
|---|---|---|
| **Process Control** | CreateProcess()<br>ExitProcess()<br>WaitForSingleObject() | fork()<br>exit()<br>wait() |
| **File Manipulation** | CreateFile()<br>ReadFile()<br>WriteFile()<br>CloseHandle() | open()<br>read()<br>write()<br>close() |
| **Device Manipulation** | SetConsoleMode()<br>ReadConsole()<br>WriteConsole() | ioctl()<br>read()<br>write() |
| **Information Maintenance** | GetCurrentProcessID()<br>SetTimer()<br>Sleep() | getpid()<br>alarm()<br>sleep() |
| **Communication** | CreatePipe()<br>CreateFileMapping()<br>MapViewOfFile() | pipe()<br>shm_open()<br>mmap() |
| **Protection** | SetFileSecurity()<br>InitlializeSecurityDescriptor() | chmod()<br>umask() |

## Review Questions

1) What are the differences between a command interpreter and a graphical user interface?

2) Explain the correlation between APIs and system calls.

3) What are the six major categories of system calls?

4) In what ways does an operating system manage resource allocation among multiple users or jobs?

5) What options do users have when choosing a shell?

# *Processes*

Early computers operated one program at a time, giving it complete control over system resources. Modern computer systems, however, support concurrent execution of multiple programs, necessitating greater control and compartmentalization, leading to the concept of a **process**— a program in execution and the fundamental unit of work in time-sharing systems.

As operating systems grow more complex, they are expected to handle not only user program execution but also various system tasks that should remain separate from the kernel. Thus, a system comprises a collection of processes, including both operating system processes (running system code) and user processes (running user code). These processes can execute simultaneously, allowing the operating system to switch the CPU between them, enhancing overall computer productivity. This chapter will explore the nature of processes and their functionality.

# 3.1 Process Concept

**Program vs. Process**

- A program is a passive entity such as the file that contains the list of instructions stored on a disk always referred to as an executable file.
- A program becomes a process when an executable file is loaded into the memory and then becomes an active entity.
- The fundamental task of any operating system is the process management.
- Processes include not only a text but also include a set of resources such as open files and pending signals. Processes also contain internal kernel data, processor state, an address space, and a data section.

**Process elements**

Segments of a process represent the following components:

**Text Section:** the program code. This is typically read-only, and might be shared by a number of processes.

**Data Section:** containing global variables.

**Heap:** containing memory dynamically allocated during run time.

**Stack:** containing temporary data.

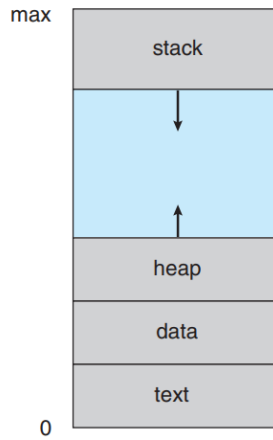- Function parameters, return addresses, local variables.

**Figure 3.1** Process in memory.

### 3.1.1 Process State

As a process executes, it changes **state**. The state of a process is defined in part by the current activity of that process. A process may be in one of the following states:

- **New**. The process is being created.

- **Running**. Instructions are being executed.

- **Waiting**. The process is waiting for some event to occur (such as an I/O completion or reception of a signal).

- **Ready**. The process is waiting to be assigned to a processor.

- **Terminated**. The process has finished execution.

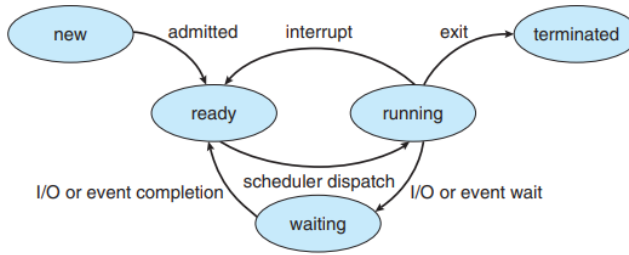The state diagram corresponding to these states is presented in Figure 3.2.

**Figure 3.2** Diagram of process state.

### 3.1.3 Process Control Block

Each process is represented in the operating system by a **process control block** (**PCB**) — also called a **task control block**.A PCB is shown in Figure 3.3. It contains many pieces of information associated with a specific process, including these:

- **Process state**. The state may be new, ready, running, waiting, halted, and so on.

- **Program counter**. The counter indicates the address of the next instruction to be executed for this process.

- **CPU registers**. The registers vary in number and type, depending on the computer architecture. They include accumulators, index registers, stack pointers, and general-purpose registers, plus any condition-code information. Along with the program counter, this state information must be saved when an interrupt occurs, to allow the process to be continued correctly afterward (Figure 3.4).

- **CPU-scheduling information**. This information includes a process priority, pointers to scheduling queues, and any other scheduling parameters.
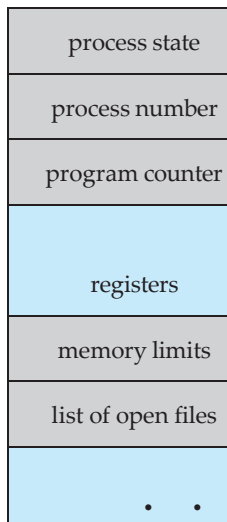
| process state |
| process number |
| program counter |
| registers |
| memory limits |
| list of open files |
| . . |

**Figure 3.3** Process control block (PCB).

- **Memory-management information**. This information may include such items as the value of the base and limit registers and the page tables, or the segment tables, depending on the memory system used by the operating system.
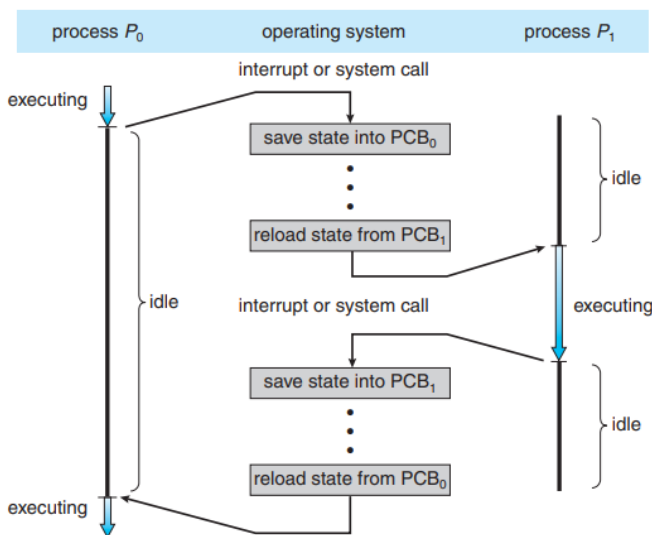


**Figure 3.4** Diagram showing CPU switch from process to process.

- **Accounting information**. This information includes the amount of CPU and real time used, time limits,

account numbers, job or process numbers, and so on.

- **I/O status information**. This information includes the list of I/O devices allocated to the process, a list of open files, and so on.

In brief, the PCB simply serves as the repository for any information that may vary from process to process.

## 3.2  Process Scheduling

The objective of multiprogramming is to have some process running at all times, to maximize CPU utilization. The objective of time sharing is to switch the CPU among processes so frequently that users can interact with each program while it is running. To meet these objectives, the **process scheduler** selects  an available process (possibly from a set of several available processes) for program execution on the CPU. For a single-processor system, there will never be more than one running process. If there are more processes, the rest will have to wait until the CPU is free and can be rescheduled.

### 3.2.1 Scheduling Queues

When processes enter a system, they are placed in a **job queue**, which includes all processes. Processes that are in **main memory** and ready to execute are listed in the **ready queue**, typically organized as a linked list with a header that points to the first and last Process Control Blocks (PCBs). Each PCB has a pointer to the next PCB in the queue.

The system also features **device queues** for processes waiting on I/O requests, such as to a shared disk. If a process requests I/O while the device is busy, it must wait in the respective device queue, which is unique to each device.

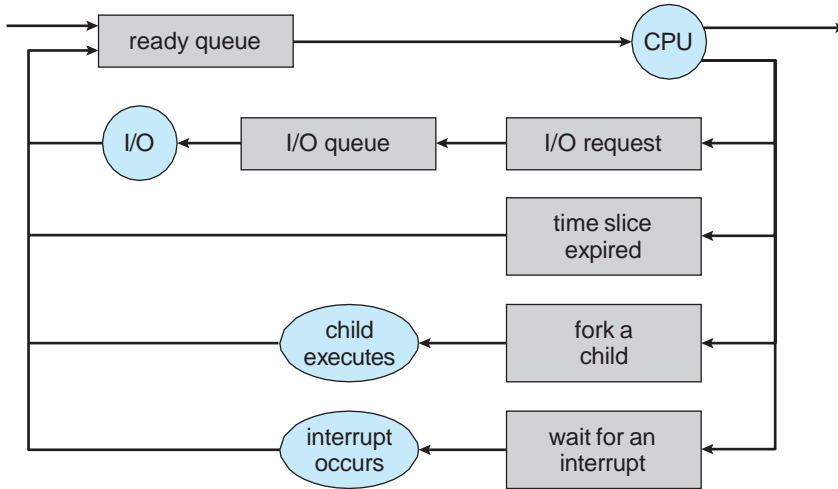Each device has its own device queue (Figure 3.5).

**Figure 3.6**  Queueing-diagram representation of process scheduling.

### 3.2.2 Schedulers

A process migrates among the various scheduling queues throughout its lifetime. The operating system must select, for scheduling purposes, processes from these queues in some fashion. The selection process is carried out by the appropriate **scheduler**.

Often, in a batch system, more processes are submitted than can be executed immediately. These processes are spooled to a mass-storage device (typically a disk), where they are kept for later execution. The **long-term scheduler**, or **job scheduler**, selects processes from this pool and loads them into memory for execution. The **short-term scheduler**, or **CPU scheduler**, selects from among the processes that are ready to execute and allocates the CPU to one of them.

38

**Short-term scheduler (or CPU scheduler)** – selects which process should be executed next and allocates CPU

– Sometimes the only scheduler in a system

– Short-term scheduler is invoked frequently (milliseconds) ⇒ (must be fast)

• **Long-term scheduler (or job scheduler)** – selects which processes should be brought into the ready queue

–Long-term scheduler is invoked infrequently (seconds, minutes) ⇒ (may be slow)

– The long-term scheduler controls the degree of multiprogramming

**Medium-term scheduler** can be added if degree of multiple programming needs to decrease

- Remove process from memory, store on disk, bring back in from disk to continue execution: swapping
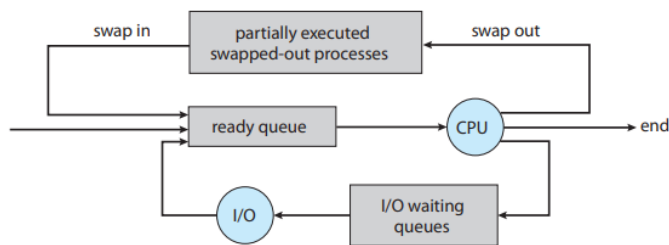


**Figure 3.7** Addition of medium-term scheduling to the queueing diagram.

## 3.3 Interprocess Communication

Processes executing concurrently in the operating system may be either independent processes or cooperating processes. A process is *independent* if it cannot affect or be affected by the other processes executing in the system. Any process that does not share data with any other process is independent. A process is *cooperating* if it can affect or be affected by the other processes executing in the system. Clearly, any process that shares data with other processes is a cooperating process.

There are several reasons for providing an environment that allows process cooperation:

- **Information sharing**. Since several users may be interested in the same piece of information (for instance, a shared file), we must provide an environment to allow concurrent access to such information.

- **Computation speedup**. If we want a particular task to run faster, we must break it into subtasks, each of which will be executing in parallel with the others. Notice that such a speedup can be achieved only if the computer has multiple processing cores.

- **Modularity**. We may want to construct the system in a modular fashion, dividing the system functions into separate processes or threads.

- **Convenience**. Even an individual user may work on many tasks at the same time. For instance, a user may be editing, listening to music, and compiling in parallel.

**Cooperating processes** require an **interprocess communication (IPC)** mechanism for data exchange, typically implemented through two fundamental models: **shared memory** and **message passing**.

1.  **Shared Memory**: In this model, processes share a designated region of memory, allowing them to read and write data directly to this space for communication. Once the shared memory is established, accesses occur as routine memory operations, minimizing kernel intervention.
2.  **Message Passing**: This model facilitates communication through messages exchanged between processes. It is particularly beneficial for transferring smaller data amounts, as it avoids conflicts, and is easier to implement in distributed systems.

Both models have their advantages. Shared memory tends to be faster due to less frequent kernel involvement, while message passing offers better performance in multi-core systems because it circumvents cache coherence issues that can arise with shared memory. As the number of processing cores increases, message passing may become the preferred IPC method. The two communications models are contrasted in Figure 3.8.
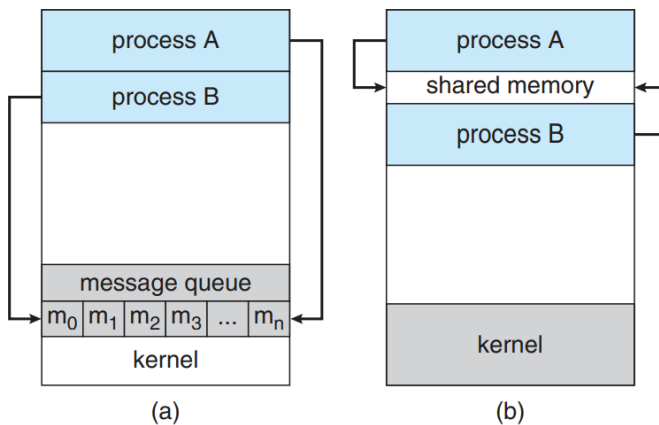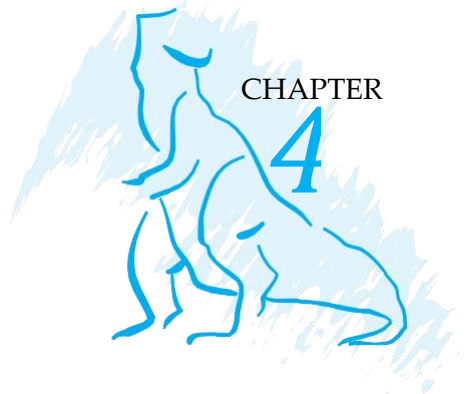
**Figure 3.8** Communications models. (a) Message passing. (b) Shared memory.

# Review Questions

**1) What is the difference between a program and a process?**
**2) Why do modern computer systems support the execution of multiple processes simultaneously?**
**3) What are the main components of a process in memory?**
**4) Describe the different states a process can be in during its lifecycle.**
**5) What is the role of the Process Control Block (PCB) in managing processes?**
**6) How does the short-term scheduler differ from the long-term scheduler in process management?**

**7) What are the different types of queues used in process scheduling, and what is their purpose?**

# *Threads*

The process model introduced in Chapter 3 assumed that a process was an executing program with a single thread of control. Virtually all modern operating systems, however, provide features enabling a process to contain multiple threads of control. In this chapter, we introduce many concepts associated with multithreaded computer systems, including a discussion of the APIs for the Pthreads, Windows, and Java thread libraries. We look at a number of issues related to multithreaded programming and its effect on the design of operating systems. Finally, we explore how the Windows and Linux operating systems support threads at the kernel level.

## 4.1 Overview

A thread is a basic unit of CPU utilization; it comprises a thread ID, a program counter, a register set, and a stack. It shares with other threads belonging to the same process its code section, data section, and other operating-system resources, such as open files and signals. A traditional (or *heavyweight*) process has a single thread of control. If a process has multiple threads of control, it can perform more than one task at a time. Figure 4.1 illustrates the difference between a traditional **single-threaded** process and a **multithreaded** process.

### 4.1.1Motivation

Modern software applications are predominantly multithreaded, with multiple threads running within a single process to perform various tasks concurrently. For example, a web browser may use different threads for displaying content and retrieving data, while a word processor might manage keystrokes, display graphics, and check grammar simultaneously. Multithreading is especially beneficial in multicore systems, allowing CPU-intensive tasks to be executed in parallel across different cores.

In certain scenarios, such as web servers handling multiple client requests, multithreading proves far more efficient than the traditional single-threaded approach. A single-threaded web server can only handle one client request at a time, causing delays for other clients. Previously, web servers created separate processes for each request, which is resource-intensive and slow. However, with multithreading, the server creates a new thread for each client request, reducing the overhead of process creation and enabling faster, more efficient handling of multiple concurrent requests. This is illustrated in Figure 4.2.
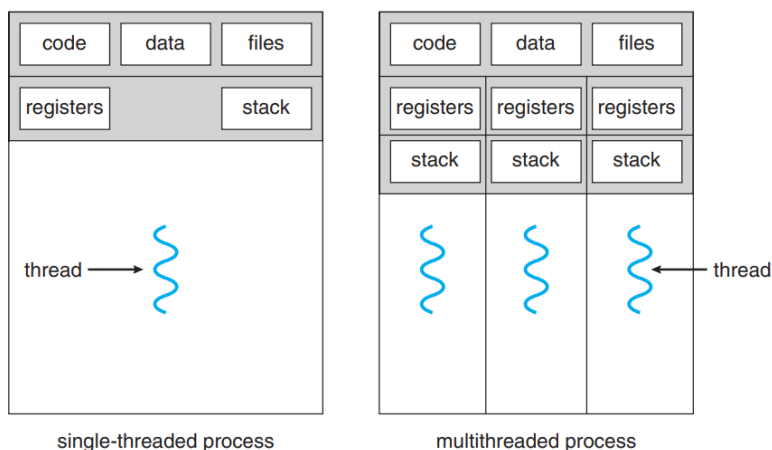


**Figure 4.1**  Single-threaded and multithreaded processes.

Threads also play a vital role in remote procedure call (RPC) systems. RPCs allow interprocess communication by providing a communication mechanism similar to ordinary function or procedure calls. Typically, RPC servers are multithreaded. When a server receives a message, it services the message using a separate thread. This allows the server to service several concurrent requests.



**Figure 4.2** Multithreaded server architecture.

Finally, most operating-system kernels are now multithreaded. Several threads operate in the kernel, and each thread performs a specific task, such as managing devices, managing memory, or interrupt handling. For example, Solaris has a set of threads in the kernel specifically for interrupt handling; Linux uses a kernel thread for managing the amount of free memory in the system.

## 4.1.2 Benefits

The benefits of multithreaded programming can be broken down into four major categories:

1. **Responsiveness**. Multithreading an interactive

application may allow a program to continue running even if part of it is blocked or is performing a lengthy operation, thereby increasing responsiveness to the user. This quality is especially useful in designing user interfaces. For instance, consider what happens when a user clicks a button that results in the performance of a time-consuming operation. A single-threaded application would be unresponsive to the user until the operation had completed. In contrast, if the time-consuming operation is performed in a separate thread, the application remains responsive to the user.

2. **Resource sharing**. Processes can only share resources through techniques such as shared memory and message passing. Such techniques must be explicitly arranged by the programmer. However, threads share the memory and the resources of the process to which they belong by default. The benefit of sharing code and data is that it allows an application to have several different threads of activity within the same address space.

3. **Economy**. Allocating memory and resources for process creation is costly. Because threads share the resources of the process to which they belong, it is more economical to create and context-switch threads. Empirically gauging the difference in overhead can be difficult, but in general it is significantly more time consuming to create and manage processes than threads. In Solaris, for example, creating a process is about thirty timesn slower than is creating a thread, and context switching is about five times slower.

4. **Scalability.** The benefits of multithreading can be even greater in a multiprocessor architecture, where threads may be running in parallel on different

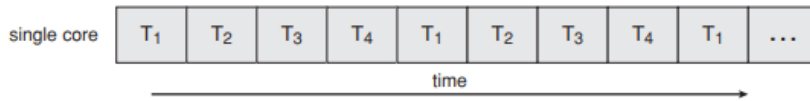processing cores. A single-threaded process can run on only one processor, regardless how many are available.



**Figure 4.3** Concurrent execution on a single-core system

## 4.2 Multicore Programming

 Earlier in the history of computer design, in response to the need for more computing performance, single-CPU systems evolved into multi-CPU systems. A more recent, similar trend in system design is to place multiple computing cores on a single chip. Each core appears as a separate processor to the operating system. Whether the cores appear across CPU chips or within CPU chips, we call these systems **multicore** or **multiprocessor** systems. Multithreaded programming provides a mechanism for more efficient use of these multiple computing cores and improved concurrency. Consider an application with four threads. On a system with a single computing core, concurrency merely means that the execution of the threads will be interleaved over time (Figure 4.3), because the processing core is capable of executing only one thread at a time. On a system with multiple cores, however, concurrency means that the threads can run in parallel, because the system can assign a separate thread to each core (Figure 4.4).

Notice the distinction between *parallelism* and *concurrency* in this discussion. A system is parallel if it can perform more than one task simultaneously. In contrast, a concurrent system supports more than one task by allowing all the tasks to make progress. Thus, it is possible to have concurrency without parallelism. Before the advent of SMP and multicore architectures, most com- puter systems had only a single processor. CPU schedulers were designed to provide the illusion of parallelism by rapidly switching between processes in

the system, thereby allowing each process to make progress. Such processes were running concurrently, but not in parallel. As systems have grown from tens of threads to thousands of threads, CPU designers have improved system performance by adding hardware to improve thread performance. Modern Intel CPUs frequently support two threads per core, while the Oracle T4 CPU supports eight threads per core. This support means that multiple threads can be loaded into the core for fast switching. Multicore computers will no doubt continue to increase in core counts and
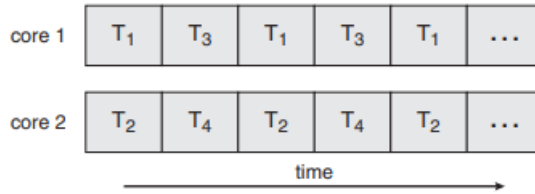
hardware thread support.



**Figure 4.4**Parallel execution on a multicore system.

## 4.2.1 Programming Challenges

The trend towards multicore systems continues to place pressure on system designers and application programmers to make better use of the multiple computing cores. Designers of operating systems must write scheduling algorithms that use multiple processing cores to allow the parallel execution shown in Figure 4.4. For application programmers, the challenge is to modify existing programs as well as design new programs that are multithreaded.

In general, five areas present challenges in programming for multicore systems:

1. **Identifying tasks**. This involves examining applications to find areas that can be divided into separate, concurrent tasks. Ideally, tasks are independent of one another and thus can run in parallel on individual cores.

2. **Balance**. While identifying tasks that can run in parallel, programmers must also ensure that the tasks perform equal work of equal value. In some instances, a certain task may not contribute as much value to the overall process as other tasks. Using a separate execution core to run that task may not be worth the

cost.

3. **Data splitting**. Just as applications are divided into separate tasks, the data accessed and manipulated by the tasks must be divided to run on separate cores.

4. **Data dependency**. The data accessed by the tasks must be examined for dependencies between two or more tasks. When one task depends on data from another, programmers must ensure that the execution of the tasks is synchronized to accommodate the data dependency.

5. **Testing and debugging**. When a program is running in parallel on multiple cores, many different execution paths are possible. Testing and debugging such concurrent programs is inherently more difficult than testing and debugging single-threaded applications.

## 4.2.2  Types of Parallelism

In general, there are two types of parallelism: data parallelism and task parallelism. **Data parallelism** focuses on distributing subsets of the same data across multiple computing cores and performing the same operation on each core. Consider, for example, summing the contents of an array of size $N$. On a single-core system, one thread would simply sum the elements $[0] \ldots [N-1]$. On a dual-core system, however, thread $A$, running on core 0, could sum the elements $[0] \ldots [N/2 - 1]$ while thread $B$, running on core 1, could sum the elements $[N/2] \ldots [N-1]$. The two threads would be running in parallel on separate computing cores.

**Task parallelism** involves distributing not data but tasks (threads) across multiple computing cores. Each thread is performing a unique operation. Different threads may be operating on the same data, or they may be operating on different data. Consider again our example above. In contrast to that situation, an example of task parallelism might involve two threads, each performing a unique statistical operation on the array of elements. The threads again are operating in parallel on separate computing cores, but each is performing a unique operation.

Fundamentally, then, data parallelism involves the distribution of data across multiple cores and task parallelism on the distribution of tasks across multiple cores. In practice, however, few applications strictly follow either data or task parallelism. In most instances, applications use a hybrid of these two strategies.

## 4.3 Multithreading Models

Our discussion so far has treated threads in a generic sense. However, support for threads may be provided either at the user level, for **user threads**, or by the kernel, for **kernel threads**. User threads are supported above the kernel and are managed without kernel support, whereas kernel threads are supported and managed directly by the operating system. Virtually all contemporary operating systems— including Windows, Linux, Mac OS X, and Solaris— support kernel threads.

Ultimately, a relationship must exist between user threads and kernel threads. In this section, we look at three common ways of establishing such a relationship: the many-to-one model, the one-to-one model, and the many-to- many models.

### 4.3.1 Many-to-One Model

The **many-to-one model** connects multiple user threads to a single kernel thread. It is efficient because thread management happens in user space. However, if one thread makes a system call that blocks, all threads in the process are blocked. Additionally, this model can't run threads in parallel on multicore systems. Due to these limitations, it's rarely used today.
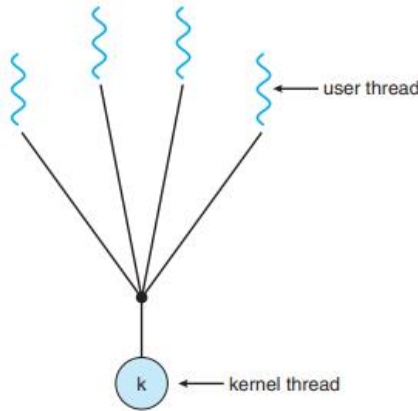
**Figure 4.5** Many-to-one model.

## 4.3.2 One-to-One Model

The one-to-one model (Figure 4.6) maps each user thread to a kernel thread. It provides more concurrency than the many-to-one model by allowing another thread to run when a thread makes a blocking system call. It also allows multiple threads to run in parallel on multiprocessors. The only drawback to this model is that creating a user thread requires creating the corresponding kernel thread. Because the overhead of creating kernel threads can burden the performance of an application, most implementations of this model restrict the number of threads supported by the system. Linux, along with the family of Windows operating systems, implement the one-to-one model.



**Figure 4.6** One-to-one model.

### 4.3.3 Many-to-Many Model

 The many-to-many model (Figure 4.7) multiplexes many user-level threads to a smaller or equal number of kernel threads. The number of kernel threads may be specific to either a particular application or a particular machine (an application may be allocated more kernel threads on a multiprocessor than on a single processor).

 Let's consider the effect of this design on concurrency. Whereas the many- to-one model allows the developer to create as many user threads as she wishes, it does not result in true concurrency, because the kernel can schedule only one thread at a time. The one-to-one model allows greater concurrency, but the developer has to be careful not to create too many threads within an application



**Figure 4.7** Many-to-many model.

**Figure 4.8** Two-level model.

The many-to-many model suffers from neither of these shortcomings: developers can create as many user threads as necessary, and the corresponding kernel threads can run in parallel on a multiprocessor. Also, when a thread performs a blocking system call, the kernel can schedule another thread for execution.

One variation on the many-to-many model still multiplexes many user- level threads to a smaller or equal number of kernel threads but also allows a user-level thread to be bound to a kernel thread. This variation is sometimes referred to as the **two-level model** (Figure 4.8). The Solaris operating system supported the two-level model in versions older than Solaris 9. However, beginning with Solaris 9, this system uses the one-to-one model.

## Review Questions

1) **What is a thread, and what are the basic components that define it?**
2) **How do threads in a process share resources?**
3) **What is the difference between a single-threaded and multithreaded process?**
4) **What are the main benefits of multithreading in modern software applications?**

5) **Why is multithreading particularly advantageous in multicore systems?**

# CPU Scheduling

CPU scheduling is the basis of multiprogrammed operating systems. By switching the CPU among processes, the operating system can make the computer more productive. In this chapter, we introduce basic CPU-scheduling concepts and present several CPU-scheduling algorithms. We also consider the problem of selecting an algorithm for a particular system.

In Chapter 4, we introduced threads to the process model. On operating systems that support them, it is kernel-level threads— not processes— that are in fact being scheduled by the operating system. However, the terms "process scheduling" and "thread scheduling" are often used interchangeably. In this chapter, we use *process scheduling* when discussing general scheduling concepts and *thread scheduling* to refer to thread-specific ideas.

## 5.1   Basic Concepts

In a single-processor system, only one process can run at a time. Others  must wait until the CPU is free and can be rescheduled. The objective of multiprogramming is to have some process running at all times, to maximize CPU  utilization. The idea is relatively simple. A process is executed until it must wait, typically for the completion of some I/O request. In a simple computer system, the CPU then just sits idle. All this waiting time is  wasted;  no  useful  work  is  accomplished.  With multiprogramming, we try to use this time productively. Several processes are kept in memory at one time.

**Figure 5.1**   Alternating sequence of CPU and I/O bursts.

When one process has to wait, the operating system takes the CPU away from that process and gives the CPU to another process. This pattern continues. Every time one process has to wait, another process can take over use of the CPU.

Scheduling of this kind is a fundamental operating-system function. Almost all computer resources are scheduled before use. The CPU is, of course, one of the primary computer resources. Thus, its scheduling is central to operating-system design.

### 5.1.1 CPU–I/O Burst Cycle

The success of CPU scheduling depends on an observed property of processes: process execution consists of a **cycle** of CPU execution and I/O wait. Processes alternate between these two states. Process execution begins with a **CPU burst**.

That is followed by an **I/O burst**, which is followed by another CPU burst, then another I/O burst, and so on. Eventually, the final CPU burst ends with a system request to terminate execution (Figure 5.1).

## 1.2 CPU Scheduler

Whenever the CPU becomes idle, the operating system must select one of the processes in the ready queue to be executed. The selection process is carried out by the **short-term scheduler**, or CPU scheduler. The scheduler selects a process from the processes in memory that are ready to execute and allocates the CPU to that process.

Note that the ready queue is not necessarily a first-in, first-out (FIFO) queue. As we shall see when we consider the various scheduling algorithms, a ready queue can be implemented as a FIFO queue, a priority queue, a tree, or simply an unordered linked list. Conceptually, however, all the processes in the ready queue are lined up waiting for a chance to run on the CPU. The records in the queues are generally process control blocks (PCBs) of the processes.

## 5.1.3 Preemptive Scheduling

CPU scheduling decisions occur under the following four conditions:

1. **Process moves from running to waiting state** (e.g., due to I/O requests or wait()).
2. **Process moves from running to ready state** (e.g., after an interrupt).
3. **Process moves from waiting to ready state** (e.g., when I/O completes).
4. **Process terminates**.

For conditions 1 and 4, scheduling is mandatory, as a new process must be selected for execution. For conditions 2 and

3, the operating system can choose whether to continue running the current process or to switch to another.

- **Non-preemptive (cooperative) scheduling**: Scheduling occurs only during conditions 1 and 4. Once a process gets the CPU, it retains control until it either finishes or enters the waiting state. Early systems like **Windows 3.x** used this model.
- **Preemptive scheduling**: Allows process switching under all four conditions. Modern systems like **Windows 95** and later versions, as well as **Mac OS X**, use preemptive scheduling.

Preemptive scheduling introduces potential issues, such as **race conditions** when processes share data. If a process is preempted while modifying data, another process could access inconsistent data. This challenge also extends to the operating system kernel, particularly during system calls or I/O handling.

To prevent inconsistency, some operating systems, like **UNIX**, wait until system calls are completed or I/O blocks occur before switching processes. However, this approach is unsuitable for **real-time computing** where tasks need to meet strict deadlines.

Additionally, **interrupts** must be handled promptly, and certain critical code sections disable interrupts temporarily to prevent simultaneous access by multiple processes. These sections are typically brief and infrequent, ensuring minimal impact on system performance.

### 5.1.4 Dispatcher

Another component involved in the CPU-scheduling function is the **dispatcher**. The dispatcher is the module that gives control of the CPU to the process selected by the short-term scheduler. This function involves the following:

- Switching context

- Switching to user mode

- Jumping to the proper location in the user program to restart that program

The dispatcher should be as fast as possible, since it is invoked during every process switch. The time it takes for the dispatcher to stop one process and start another running is known as the **dispatch latency**.

## 5.2 Scheduling Criteria

Different CPU-scheduling algorithms have different properties, and the choice of a particular algorithm may favor one class of processes over another. In choosing which algorithm to use in a particular situation, we must consider the properties of the various algorithms.

Many criteria have been suggested for comparing CPU-scheduling algo- rithms. Which characteristics are used for comparison can make a substantial difference in which algorithm is judged to be best. The criteria include the following:

- **CPU utilization**. We want to keep the CPU as busy as possible. Concep- tually, CPU utilization can range from 0 to 100 percent. In a real system, it should range from 40 percent (for a lightly loaded system) to 90 percent (for a heavily loaded system).

- **Throughput**. If the CPU is busy executing processes,

then work is being done. One measure of work is the number of processes that are completed per time unit, called throughput. For long processes, this rate may be one process per hour; for short transactions, it may be ten processes per second.

- **Turnaround time**. From the point of view of a particular process, the important criterion is how long it takes to execute that process. The interval from the time of submission of a process to the time of completion is the turnaround time. Turnaround time is the sum of the periods spent waiting to get into memory, waiting in the ready queue, executing on the CPU, and doing I/O.

- **Waiting time**. The CPU-scheduling algorithm does not affect the amount of time during which a process executes or does I/O. It affects only the amount of time that a process spends waiting in the ready queue. Waiting time is the sum of the periods spent waiting in the ready queue.

- **Response time**. In an interactive system, turnaround time may not be the best criterion. Often, a process can produce some output fairly early and can continue computing new results while previous results are being output to the user. Thus, another measure is the time from the submission of a request until the first response is produced. This measure, called response time, is the time it takes to start responding, not the time it takes to output the response. The turnaround time is generally limited by the speed of the output device.

It is desirable to maximize CPU utilization and throughput and to minimize turnaround time, waiting time, and response time. In most cases, we optimize the average

measure. However, under some circumstances, we prefer to optimize the minimum or maximum values rather than the average. For example, to guarantee that all users get good service, we may want to minimize the maximum response time.

Investigators have suggested that, for interactive systems (such as desktop systems), it is more important to minimize the variance in the response time than to minimize the average response time. A system with reasonable and predictable response time may be considered more desirable than a system that is faster on the average but is highly variable. However, little work has been done on CPU-scheduling algorithms that minimize variance.

As we discuss various CPU-scheduling algorithms in the following section, we illustrate their operation. An accurate illustration should involve many processes, each a sequence of several hundred CPU bursts and I/O bursts. For simplicity, though, we consider only one CPU burst (in milliseconds) per process in our examples. Our measure of comparison is the average waiting time.

## 5.3 Scheduling Algorithms

CPU scheduling deals with the problem of deciding which of the processes in the ready queue is to be allocated the CPU. There are many different CPU-scheduling algorithms. In this section, we describe several of them.

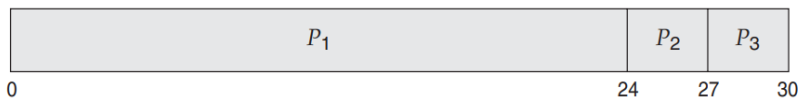### 5.3.1 First-Come, First-Served Scheduling

By far the simplest CPU-scheduling algorithm is the **first-come, first-served** (**FCFS**) scheduling algorithm. With this scheme, the process that requests the CPU first is allocated the CPU first. The implementation of the FCFS policy is easily managed with a FIFO queue. When a process enters the ready queue, its PCB is linked onto the tail of the queue. When the

CPU is free, it is allocated to the process at the head of the queue. The running process is then removed from the queue. The code for FCFS scheduling is simple to write and understand.

On the negative side, the average waiting time under the FCFS policy is often quite long. Consider the following set of processes that arrive at time 0, with the length of the CPU burst given in milliseconds:

| Process | Burst Time |
|---------|-----------|
| $P_1$ | 24 |
| $P_2$ | 3 |
| $P_3$ | 3 |

If the processes arrive in the order $P_1$, $P_2$, $P_3$, and are served in FCFS order, we get the result shown in the following **Gantt chart**, which is a bar chart that illustrates a particular schedule, including the start and finish times of each of the participating processes:

| $P_1$ | | | $P_2$ | $P_3$ |
|-------|---|---|-------|-------|
| 0 | | | 24 | 27 | 30 |

The waiting time is 0 milliseconds for process $P_1$, 24 milliseconds for process $P_2$, and 27 milliseconds for process $P_3$. Thus, the average waiting time is $(0 + 24 + 27)/3 = 17$ milliseconds. If the processes arrive in the order $P_2$, $P_3$, $P_1$, however, the results will be as shown in the following Gantt chart:

| $P_2$ | $P_3$ | $P_1$ |
|-------|-------|-------|
| 0 | 3 | 6 | 30 |

The average waiting time is now $(6 + 0 + 3)/3 = 3$ milliseconds. This reduction is substantial. Thus, the average waiting time under an FCFS policy is generally not minimal and may vary substantially if the processes' CPU burst times vary greatly.

In addition, consider the performance of FCFS scheduling in a dynamic situation. Assume we have one CPU-bound process and many I/O-bound processes. As the processes flow around the system, the following scenario may result. The CPU-bound process will get and hold the CPU. During this time, all the other processes will finish their I/O and will move into the ready queue, waiting for the CPU. While the processes wait in the ready queue, the I/O devices are idle. Eventually, the CPU-bound process finishes its CPU burst and moves to an I/O device. All the I/O-bound processes, which have short CPU bursts, execute quickly and move back to the I/O queues. At this point, the CPU sits idle. The CPU-bound process will then move back to the ready queue and be allocated the CPU. Again, all the I/O processes end up waiting in the ready queue until the CPU-bound process is done. There is a **convoy effect** as all the other processes wait for the one big process to get off the CPU. This effect results in lower CPU and device utilization than might be possible if the shorter processes were allowed to go first.

Note also that the FCFS scheduling algorithm is nonpreemptive. Once the CPU has been allocated to a process, that process keeps the CPU until it releases the CPU, either by terminating or by requesting I/O. The FCFS algorithm is thus particularly troublesome for time-sharing systems, where it is important that each user get a share of the CPU at regular intervals. It would be disastrous to allow one process to keep the CPU for an extended period.
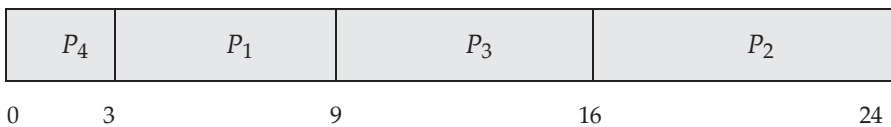
### 5.3.2 Shortest-Job-First Scheduling

A different approach to CPU scheduling is the **shortest-job-first** (**SJF**) scheduling algorithm. This algorithm associates with each process the length of the process's next CPU burst. When the CPU is available, it is assigned to the process that has the smallest next CPU burst. If the next CPU bursts of two processes are the same, FCFS scheduling is used to break the tie. Note that a more appropriate term for this scheduling method would be the *shortest-next- CPU-burst* algorithm, because scheduling depends on the length of the next CPU burst of a process, rather than its total length. We use the term SJF because most people and textbooks use this term to refer to this type of scheduling.

As an example of SJF scheduling, consider the following set of processes, with the length of the CPU burst given in milliseconds:

| Process | Burst Time |
|---------|------------|
| $P_1$   | 6          |
| $P_2$   | 8          |
| $P_3$   | 7          |
| $P_4$   | 3          |

Using SJF scheduling, we would schedule these processes according to the following Gantt chart:

| $P_4$ | $P_1$ | $P_3$ | $P_2$ |
|-------|-------|-------|-------|
| 0   3 |     9 |    16 |    24 |

The waiting time is 3 milliseconds for process $P_1$, 16 milliseconds for process $P_2$, 9 milliseconds for process $P_3$, and 0 milliseconds for

process $P_4$. Thus, the average waiting time is $(3 + 16 + 9 + 0)/4 = 7$ milliseconds. By comparison, if we were using the FCFS scheduling scheme, the average waiting time would be 10.25 milliseconds.

The SJF scheduling algorithm is provably optimal, in that it gives the minimum average waiting time for a given set of processes. Moving a short process before a long one decreases the waiting time of the short process more than it increases the waiting time of the long process. Consequently, the average waiting time decreases.
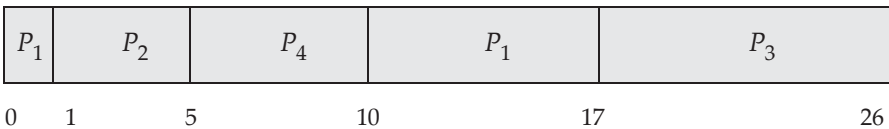
The real difficulty with the SJF algorithm is knowing the length of the next CPU request. For long-term (job) scheduling in a batch system, we can use the process time limit that a user specifies when he submits the job. In this situation, users are motivated to estimate the process time limit accurately, since a lower value may mean faster response but too low a value will cause a time-limit-exceeded error and require resubmission. SJF scheduling is used frequently in long-term scheduling.

Although the SJF algorithm is optimal, it cannot be implemented at the level of short-term CPU scheduling. With short-term scheduling, there is no way to know the length of the next CPU burst. One approach to this problem is to try to approximate SJF scheduling. We may not know the length of the next CPU burst, but we may be able to predict its value. We expect that the next CPU burst will be similar in length to the previous ones. By computing an approximation of the length of the next CPU burst, we can pick the process with the shortest predicted CPU burst.

As an example, consider the following four processes, with the length of the CPU burst given in milliseconds:

| Process | Arrival Time | Burst Time |
|---------|--------------|------------|
| $P_1$ | 0 | 8 |
| $P_2$ | 1 | 4 |
| $P_3$ | 2 | 9 |
| $P_4$ | 3 | 5 |

If the processes arrive at the ready queue at the times shown and need the indicated burst times, then the resulting preemptive SJF schedule is as depicted in the following Gantt chart:

| $P_1$ | $P_2$ | $P_4$ | $P_1$ | $P_3$ |
|-------|-------|-------|-------|-------|

0   1       5          10             17                      26

Process $P_1$ is started at time 0, since it is the only process in the queue. Process $P_2$ arrives at time 1. The remaining time for process $P_1$ (7 milliseconds) is larger than the time required by process $P_2$ (4 milliseconds), so process $P_1$ is preempted, and process $P_2$ is scheduled. The average waiting time for this example is $[(10 − 1) + (1 − 1) + (17 − 2) + (5 − 3)]/4 = 26/4 = 6.5$ milliseconds. Nonpreemptive SJF scheduling would result in an average waiting time of 7.75 milliseconds.

### 5.3.3 Priority Scheduling

The SJF algorithm is a special case of the general **priority-scheduling** algorithm. A priority is associated with each process, and the CPU is allocated to the process with the highest priority. Equal-priority processes are scheduled in FCFS order.
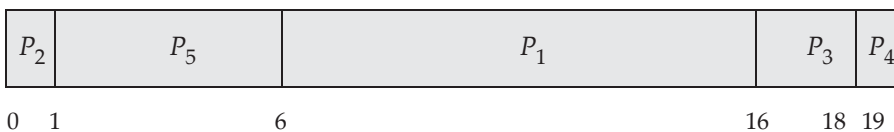
An SJF algorithm is simply a priority algorithm where the priority ($p$) is the inverse of the (predicted) next CPU burst. The larger the CPU burst, the lower the priority, and vice versa.

Note that we discuss scheduling in terms of **high** priority and **low** priority. Priorities are generally indicated by some fixed range of numbers, such as 0 to 7 or 0 to 4,095. However, there is no general agreement on whether 0 is the highest or lowest priority. Some systems use low numbers to represent low priority; others use low numbers for high priority. This difference can lead to confusion. In this text, we assume that low numbers represent high priority.

As an example, consider the following set of processes, assumed to have arrived at time 0 in the order $P_1$, $P_2$, ... , $P_5$, with the length of the CPU burst given in milliseconds:

| Process | Burst Time | Priority |
|---------|------------|----------|
| $P_1$ | 10 | 3 |
| $P_2$ | 1 | 1 |
| $P_3$ | 2 | 4 |
| $P_4$ | 1 | 5 |
| $P_5$ | 5 | 2 |

Using priority scheduling, we would schedule these processes according to the following Gantt chart:

| $P_2$ | $P_5$ | $P_1$ | $P_3$ | $P_4$ |
|---|---|---|---|---|

0    1                6                              16        18  19

The average waiting time is 8.2 milliseconds.

Priorities can be defined either internally or externally. Internally defined priorities use some measurable quantity or quantities to compute the priority of a process. For example, time limits, memory requirements, the number of open files, and the ratio of average I/O burst to average CPU burst have been used in computing priorities. External priorities are set by criteria outside  the operating system, such as the importance of the process, the type and amount of funds being paid for computer use, the department sponsoring the work, and other, often political, factors.

Priority scheduling can be either preemptive or nonpreemptive. When a process arrives at the ready queue, its priority is compared with the priority of the currently running process. A preemptive priority scheduling algorithm will preempt the CPU if the priority of the newly arrived process is higher than the priority of the currently running process. A nonpreemptive priority scheduling algorithm will simply put the new process at the head of the ready queue.

A major problem with priority scheduling algorithms is **indefinite block- ing**, or **starvation**. A process that is ready to run but waiting for the CPU can be considered blocked. A priority scheduling algorithm can leave some low- priority processes waiting indefinitely. In a heavily loaded computer system, a steady stream of higher-priority processes can prevent a low-priority process from ever getting the CPU. Generally, one of two things will happen. Either the process will eventually be run, or the computer system will eventually crash and lose all unfinished low-priority processes.

A solution to the problem of indefinite blockage of low-priority processes is **aging**. Aging involves gradually increasing the priority of processes that wait in the system for a long time. For example, if priorities range from 127 (low) to 0 (high), we could increase the priority of a waiting process

by 1 every 15 minutes. Eventually, even a process with an initial priority of 127 would have the highest priority in the system and would be executed. In fact, it would take no more than 32 hours for a priority-127 process to age to a priority-0 process.

### 5.3.4 Round-Robin Scheduling

The **round-robin** (**RR**) scheduling algorithm is designed especially for time- sharing systems. It is similar to FCFS scheduling, but preemption is added to enable the system to switch between processes. A small unit of time, called a **time quantum** or **time slice**, is defined. A time quantum is generally from 10 to 100 milliseconds in length. The ready queue is treated as a circular queue.

The CPU scheduler goes around the ready queue, allocating the CPU to each process for a time interval of up to 1 time quantum.

To implement RR scheduling, we again treat the ready queue as a FIFO queue of processes. New processes are added to the tail of the ready queue. The CPU scheduler picks the first process from the ready queue, sets a timer to interrupt after 1 time quantum, and dispatches the process.

One of two things will then happen. The process may have a CPU burst of less than 1 time quantum. In this case, the process itself will release the CPU voluntarily. The scheduler will then proceed to the next process in the ready queue. If the CPU burst of the currently running process is longer than 1 time quantum, the timer will go off and will cause an interrupt to the operating system. A context switch will be executed, and the process will be put at the tail of the ready queue. The CPU scheduler will then select the next process in the ready queue.

The average waiting time under the RR policy is often long. Consider the following set of processes that arrive at time 0, with the length of the CPU burst given in milliseconds:

| Process | Burst Time |
|---------|------------|
| $P_1$ | 24 |
| $P_2$ | 3 |
| $P_3$ | 3 |

If we use a time quantum of 4 milliseconds, then process $P_1$ gets the first 4 milliseconds. Since it requires another 20 milliseconds, it is preempted after the first time quantum, and the CPU is given to the next process in the queue, process $P_2$. Process $P_2$ does not need 4 milliseconds, so it quits before its time quantum expires. The CPU is then given to the next process, process $P_3$. Once each process has received 1 time quantum, the CPU is returned to process $P_1$ for an additional time quantum. The resulting RR schedule is as follows:

| $P_1$ | $P_2$ | $P_3$ | $P_1$ | $P_1$ | $P_1$ | $P_1$ | $P_1$ |
|-------|-------|-------|-------|-------|-------|-------|-------|
| 0   4 | 7 | 10 | 14 | 18 | 22 | 26 | 30 |

Let's calculate the average waiting time for this schedule. $P_1$ waits for 6 milliseconds (10 - 4), $P_2$ waits for 4 milliseconds, and $P_3$ waits for 7 milliseconds. Thus, the average waiting time is 17/3 = 5.66 milliseconds.

In the RR scheduling algorithm, no process is allocated the CPU for more than 1 time quantum in a row (unless it is the only runnable process). If a process's CPU burst exceeds 1 time quantum, that process is preempted and is put back in the ready queue. The RR scheduling algorithm is thus preemptive. If there are $n$ processes in the ready queue

and the time quantum is $q$, then each process gets $1/n$ of the CPU time in chunks of at most $q$ time units. Each process must wait no longer than $(n\text{-}1) \times q$ time units until its next time quantum. For example, with five processes and a time quantum of 20 milliseconds, each process will get up to 20 milliseconds every 100 milliseconds. The performance of the RR algorithm depends heavily on the size of the time quantum. At one extreme, if the time quantum is extremely large, the RR policy

 is the same as the FCFS policy. In contrast, if the time quantum is extremely small (say, 1 millisecond), the RR approach can result in a large number of context switches. Assume, for example, that we have only one process of 10 time units. If the quantum is 12 time units, the process finishes in less than 1 time quantum, with no overhead. If the quantum is 6 time units, however, the process requires 2 quanta, resulting in a context switch. If the time quantum is 1 time unit, then nine context switches will occur, slowing the execution of the process accordingly (Figure 5.4).
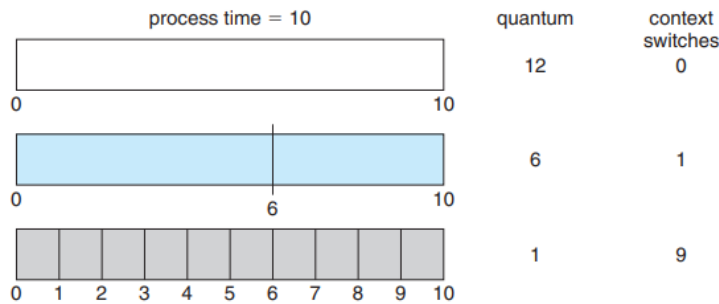
**Figure 5.4** How a smaller time quantum increases context switches.

Thus, we want the time quantum to be large with respect to the context- switch time. If the context-switch time is approximately 10 percent of the time quantum, then about 10 percent of the CPU time will be spent in context switching. In practice, most modern systems have time quanta ranging from 10 to 100 milliseconds. The time required for a context switch is typically less than 10 microseconds; thus, the context-switch time is a small fraction of the time quantum.

Turnaround time also depends on the size of the time quantum. the average turnaround time of a set of processes does not necessarily improve as the time-quantum size increases. In general, the average turnaround time can be improved if most processes finish their next CPU burst in a single time quantum. For example, given three processes of 10 time units each and a quantum of 1 time unit, the average turnaround time is 29. If the time quantum is 10, however, the average turnaround time drops to 20. If context-switch time is added in, the average turnaround time increases even more for a smaller time quantum, since more context switches are required.

Although the time quantum should be large compared with the context- switch time, it should not be too large. As we pointed out earlier, if the time quantum is too large, RR

scheduling degenerates to an FCFS policy. A rule of thumb is that 80 percent of the CPU bursts should be shorter than the time quantum.
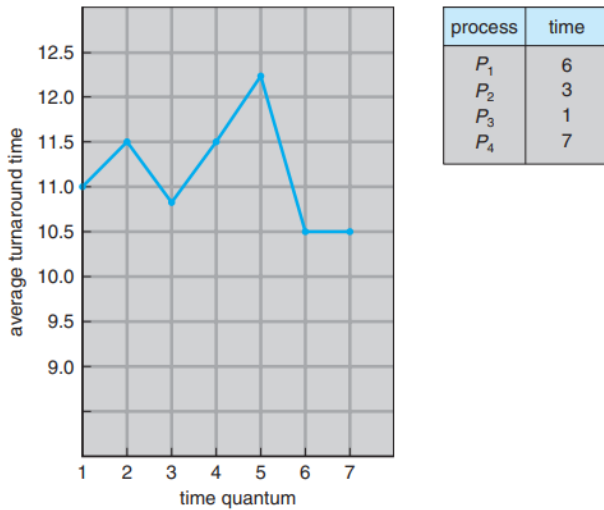


| process | time |
|---------|------|
| $P_1$ | 6 |
| $P_2$ | 3 |
| $P_3$ | 1 |
| $P_4$ | 7 |

**Figure 5.5**   How turnaround time varies with the time quantum.

## 5.3.5 Multilevel Queue Scheduling

Another class of scheduling algorithms has been created for situations in which processes are easily classified into different groups. For example, a common division is made between **foreground** (interactive) processes and **background** (batch) processes.

These two types of processes have different response-time requirements and so may have different scheduling needs. In addition, foreground processes may have priority (externally defined) over background processes.

A **multilevel queue** scheduling algorithm partitions the ready queue into several separate queues (Figure 5.6). The processes are permanently assigned to one queue, generally based on some property of the process, such as memory size, process priority, or process type. Each queue has its own

scheduling algorithm. For example, separate queues might be used for foreground and background processes. The foreground queue might be scheduled by an RR algorithm, while the background queue is scheduled by an FCFS algorithm.

In addition, there must be scheduling among the queues, which is com- monly implemented as fixed-priority preemptive scheduling. For example, the foreground queue may have absolute priority over the background queue.

Let's look at an example of a multilevel queue scheduling algorithm with five queues, listed below in order of priority:

1. System processes
2. Interactive processes
3. Interactive editing processes
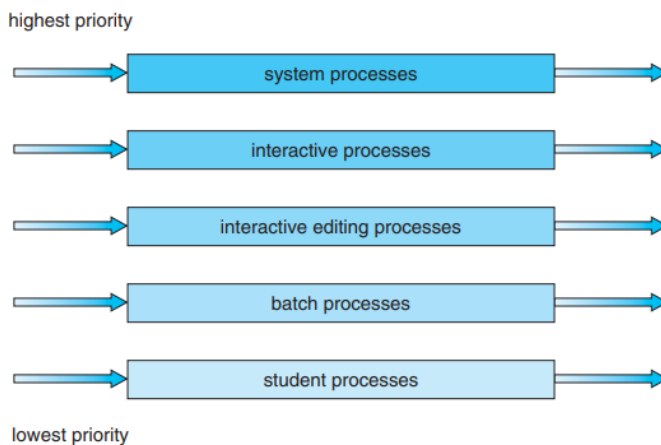4. Batch processes
5. Student processes



**Figure 5.6**  Multilevel queue scheduling.

Each queue has absolute priority over lower-priority queues.

No process in the batch queue, for example, could run unless the queues for system processes, interactive processes, and interactive editing processes were all empty. If an interactive editing process entered the ready queue while a batch process was running, the batch process would be preempted.

Another possibility is to time-slice among the queues. Here, each queue gets a certain portion of the CPU time, which it can then schedule among its various processes. For instance, in the foreground– background queue example, the foreground queue can be given 80 percent of the CPU time for RR scheduling among its processes, while the background queue receives 20 percent of the CPU to give to its processes on an FCFS basis.

### 5.3.6 Multilevel Feedback Queue Scheduling

Normally, when the multilevel queue scheduling algorithm is used, processes are permanently assigned to a queue when they enter the system. If there are separate queues for foreground and background processes, for example, processes do not move from one queue to the other, since processes do not change their foreground or background nature. This setup has the advantage of low scheduling overhead, but it is inflexible.

The **multilevel feedback queue** scheduling algorithm, in contrast, allows a process to move between queues. The idea is to separate processes according to the characteristics of their CPU bursts. If a process uses too much CPU time, it will be moved to a lower-priority queue. This scheme leaves I/O-bound and interactive processes in the higher-priority queues. In addition, a process that waits too long in a lower-priority queue may be moved to a higher-priority queue. This form of aging prevents starvation.

For example, consider a multilevel feedback queue scheduler with three queues, numbered from 0 to 2 (Figure 5.7). The scheduler first executes all processes in queue 0. Only when queue 0 is empty will it execute processes in queue 1. Similarly, processes in queue 2 will be executed only if queues 0 and 1 are empty. A process that arrives for queue 1 will preempt a process in queue 2. A process in queue 1 will in turn be preempted by a process arriving for queue 0.
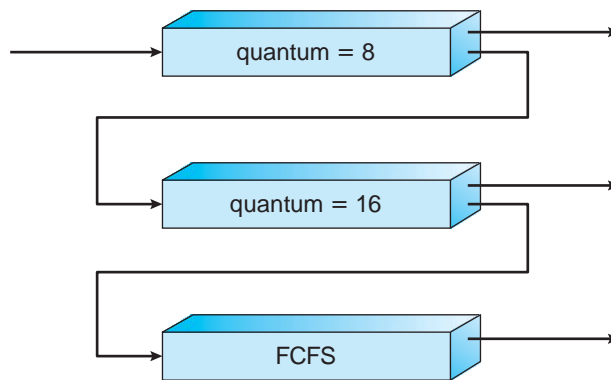
**Figure 5.7** Multilevel feedback queues.

A process entering the ready queue is put in queue 0. A process in queue 0 is given a time quantum of 8 milliseconds. If it does not finish within this time, it is moved to the tail of queue 1. If queue 0 is empty, the process at the head of queue 1 is given a quantum of 16 milliseconds. If it does not complete, it is preempted and is put into queue 2. Processes in queue 2 are run on an FCFS basis but are run only when queues 0 and 1 are empty.

This scheduling algorithm gives highest priority to any process with a CPU burst of 8 milliseconds or less. Such a process will quickly get the CPU, finish its CPU burst, and go off to its next I/O burst. Processes that need more than 8 but less than 24 milliseconds are also served quickly, although with lower priority than shorter processes. Long processes automatically sink to queue 2 and are served in FCFS order with any CPU cycles left over from queues 0 and 1.

In general, a multilevel feedback queue scheduler is defined by the following parameters:

- The number of queues
- The scheduling algorithm for each queue
- The method used to determine when to upgrade a

process to a higher- priority queue

- The method used to determine when to demote a process to a lower- priority queue

- The method used to determine which queue a process will enter when that process needs service

The definition of a multilevel feedback queue scheduler makes it the most general CPU-scheduling algorithm. It can be configured to match a specific system under design. Unfortunately, it is also the most complex algorithm, since defining the best scheduler requires some means by which to select values for all the parameters.

## Review Questions

1) **Discuss how the following pairs of scheduling criteria conflict in certain settings.**

a. CPU utilization and response time

b. Average turnaround time and maximum waiting time

c. I/O device utilization and CPU utilization

2) **Which of the following scheduling algorithms could result in starvation?**

a. **First-come, first-served**

b. **Shortest job first**

c. **Round robin**

d.**Priorit**

# Deadlocks

In a multiprogramming environment, several processes may compete for a finite number of resources. A process requests resources; if the resources are not available at that time, the process enters a waiting state. Sometimes, a waiting process is never again able to change state, because the resources it has requested are held by other waiting processes. This situation is called a **deadlock**.

Perhaps the best illustration of a deadlock can be drawn from a law passed by the Kansas legislature early in the 20th century. It said, in part: "When two trains approach each other at a crossing, both shall come to a full stop and neither shall start up again until the other has gone."

In this chapter, we describe methods that an operating system can use to prevent or deal with deadlocks. Although some applications can identify programs that may deadlock, operating systems typically do not provide deadlock-prevention facilities, and it remains the responsibility of program- mers to ensure that they design deadlock-free programs. Deadlock problems can only become more common, given current trends, including larger num- bers of processes, multithreaded programs, many more resources within a system, and an emphasis on long-lived file and

database servers rather than batch systems.

# 6.1 Deadlock Characterization

In a deadlock, processes never finish executing, and system resources are tied up, preventing other jobs from starting. Before we discuss the various methods for dealing with the deadlock problem, we look more closely at features that characterize deadlocks.

## 6.1.1 Necessary Conditions

A deadlock situation can arise if the following four conditions hold simultane- ously in a system:

1. **Mutual exclusion**. At least one resource must be held in a nonsharable mode; that is, only one process at a time can use the resource. If another process requests that resource, the requesting process must be delayed until the resource has been released.

2. **Hold and wait**. A process must be holding at least one resource and waiting to acquire additional resources that are currently being held by other processes.

3. **No preemption**. Resources cannot be preempted; that is, a resource can be released only voluntarily by the process holding it, after that process has completed its task.

4. **Circular wait**. A set $\{P_0, P_1, ..., P_n\}$ of waiting processes must exist such that $P_0$ is waiting for a resource held by $P_1$, $P_1$ is waiting for a resource held by $P_2$, ..., $P_{n-1}$ is waiting for a resource held by $P_n$, and $P_n$ is waiting for a resource held by $P_0$.

We emphasize that all four conditions must hold for a deadlock to occur. The circular-wait condition implies the

hold-and-wait condition, so the four conditions are not completely independent.

## 6.2.2 Resource-Allocation Graph

Deadlocks can be described more precisely in terms of a directed graph called a **system resource-allocation graph**. This graph consists of a set of vertices $V$ and a set of edges $E$. The set of vertices $V$ is partitioned into two different types of nodes: $P = \{P_1, P_2, ..., P_n\}$, the set consisting of all the active processes in the system, and $R = \{R_1, R_2, ..., R_m\}$, the set consisting of all resource types in the system.

A directed edge from process $P_i$ to resource type $R_j$ is denoted by $P_i \rightarrow R_j$; it signifies that process $P_i$ has requested an instance of resource type $R_j$ and is currently waiting for that resource. A directed edge from resource type $R_j$ to process $P_i$ is denoted by $R_j \rightarrow P_i$; it signifies that an instance of resource type $R_j$ has been allocated to process $P_i$. A directed edge $P_i \rightarrow R_j$ is called a **request edge**; a directed edge $R_j \rightarrow P_i$ is called an **assignment edge**.

Pictorially, we represent each process $P_i$ as a circle and each resource type $R_j$ as a rectangle. Since resource type $R_j$ may have more than one instance, we represent each such instance as a dot within the rectangle. Note that a request edge points to only the rectangle $R_j$, whereas an assignment edge must also designate one of the dots in the rectangle.

When process $P_i$ requests an instance of resource type $R_j$, a request edge is inserted in the resource-allocation graph. When this request can be fulfilled, the request edge is *instantaneously* transformed to an assignment edge. When the process no longer needs access to the resource, it releases the resource. As a result, the assignment edge is deleted.

The resource-allocation graph shown in Figure 6.1 depicts the following situation.

- The sets $P$, $R$, and $E$:
- $P = \{P_1, P_2, P_3\}$
- $R = \{R_1, R_2, R_3, R_4\}$
- $E = \{P_1 \rightarrow R_1, P_2 \rightarrow R_3, R_1 \rightarrow P_2, R_2 \rightarrow P_2, R_2 \rightarrow P_1, R_3 \rightarrow P_3\}$

- Resource instances:

- One instance of resource type $R_1$

- Two instances of resource type $R_2$

- One instance of resource type $R_3$

- Three instances of resource type $R_4$

➢ Process states:

- Process $P_1$ is holding an instance of resource type $R_2$ and is waiting for an instance of resource type $R_1$.

- Process $P_2$ is holding an instance of $R_1$ and an instance of $R_2$ and is waiting for an instance of $R_3$.

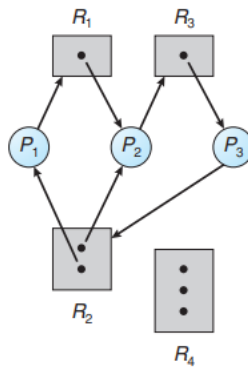- Process $P_3$ is holding an instance of $R_3$.



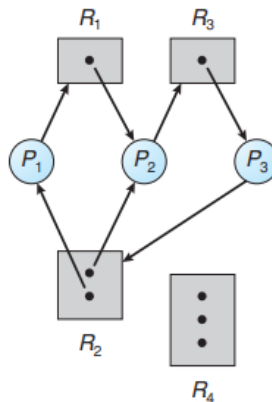**Figure 6.1** Resource-allocation graph.



**Figure 6.2** Resource-allocation graph with a deadlock.

Given the definition of a resource-allocation graph, it can be shown that, if the graph contains no cycles, then no process in the system is deadlocked. If the graph does contain a cycle, then a deadlock may exist.

If each resource type has exactly one instance, then a cycle implies that a deadlock has occurred. If the cycle involves only a set of resource types, each of which has only a single instance, then a deadlock has occurred. Each process involved in the cycle is deadlocked. In this case, a cycle in the graph is both a necessary and a sufficient condition for the existence of deadlock.

If each resource type has several instances, then a cycle does not necessarily imply that a deadlock has occurred. In this case, a cycle in the graph is a necessary but not a sufficient condition for the existence of deadlock.

To illustrate this concept, we return to the resource-allocation graph depicted in Figure 6.1. Suppose that process $P_3$ requests an instance of resource type R2. Since no resource instance is currently available, we add a request edge

P3 → R2 to the graph (Figure 6.2). At this point, two minimal cycles exist in the system: P1 → R1 → P2 → R3 → P3 → R2 → P1 P2 → R3 → P3 → R2 → P2 Processes P1, P2, and P3 are deadlocked. Process P2 is waiting for the resource R3, which is held by process P3. Process P3 is waiting for either process P1 or process P2 to release resource R2. In addition, process P1 is waiting for process P2 to release resource R1. Now consider the resource-allocation graph in Figure 6.3. In this example, we also have a cycle:

P1 → R1 → P3 → R2 → P1

However, there is no deadlock. Observe that process P4 may release its instance of resource type R2. That resource can then be allocated to P3, breaking the cycle. In summary, if a resource-allocation graph does not have a cycle, then the

system is not in a deadlocked state. If there is a cycle, then the system may or may not be in a deadlocked state. This observation is important when we deal with the deadlock problem.



**Figure 6.3** Resource-allocation graph with a cycle but no deadlock.

However, there is no deadlock. Observe that process $P_4$ may release its instance of resource type $R_2$. That resource can then be allocated to $P_3$, breaking the cycle. In summary, if a resource-allocation graph does not have a cycle, then the system is *not* in a deadlocked state. If there is a cycle, then the system may or may not be in a deadlocked state. This observation is important when we deal with the deadlock problem.
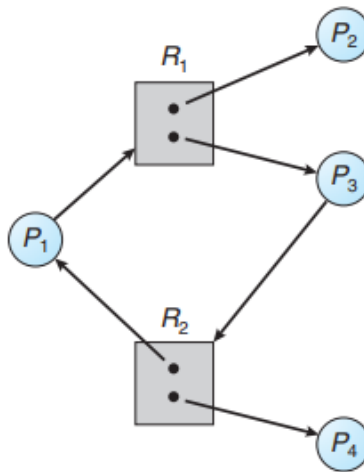
## 6.3 Methods for Handling Deadlocks

Generally speaking, we can deal with the deadlock problem in one of three ways:

- We can use a protocol to prevent or avoid deadlocks,

- ensuring that the system will *never* enter a deadlocked state.

- We can allow the system to enter a deadlocked state, detect it, and recover.

- We can ignore the problem altogether and pretend that deadlocks never occur in the system.

The third solution is the one used by most operating systems, including Linux and Windows. It is then up to the application developer to write programs that handle deadlocks.

To ensure that deadlocks never occur, the system can use either a deadlock- prevention or a deadlock-avoidance scheme. **Deadlock prevention** provides a set of methods to ensure that at least one of the necessary conditions cannot hold. These methods prevent deadlocks by constraining how requests for resources can be made.

**Deadlock avoidance** requires that the operating system be given additional information in advance concerning which resources a process will request and use during its lifetime. With this additional knowledge, the operating system can decide for each request whether or not the process should wait. To decide whether the current request can be satisfied or must be delayed, the system must consider the resources currently available, the resources currently allocated to each process, and the future requests and releases of each process.

If a system does not employ either a deadlock-prevention or a deadlock- avoidance algorithm, then a deadlock situation may arise. In this environment, the system can provide an algorithm that examines the state of the system to determine whether a deadlock has occurred and an algorithm to recover from the deadlock.

In the absence of algorithms to detect and recover from

deadlocks, we may arrive at a situation in which the system is in a deadlocked state yet has no way of recognizing what has happened. In this case, the undetected deadlock will cause the system's performance to deteriorate, because resources are being held by processes that cannot run and because more and more processes, as they make requests for resources, will enter a deadlocked state. Eventually, the system will stop functioning and will need to be restarted manually.

Although this method may not seem to be a viable approach to the deadlock problem, it is nevertheless used in most operating systems, as mentioned earlier. Expense is one important consideration. Ignoring the possibility of deadlocks is cheaper than the other approaches. Since in many systems, deadlocks occur infrequently (say, once per year), the extra expense of the other methods may not seem worthwhile. In addition, methods used to recover from other conditions may be put to use to recover from deadlock. In some circumstances, a system is in a frozen state but not in a deadlocked state. We see this situation, for example, with a real-time process running at the highest priority (or any process running on a nonpreemptive scheduler) and never returning control to the operating system. The system must have manual recovery methods for such conditions and may simply use those techniques for deadlock recovery.

## 6.4 Deadlock Prevention

for a deadlock to occur, each of the four necessary conditions must hold. By ensuring that at least one of these conditions cannot hold, we can *prevent* the occurrence of a deadlock. We elaborate on this approach by examining each of the four necessary conditions separately.

### 6.4.1 Mutual Exclusion

The mutual exclusion condition must hold. That is, at least one resource must be nonsharable. Sharable resources, in contrast,

do not require mutually exclusive access and thus cannot be involved in a deadlock. Read-only files are a good example of a sharable resource. If several processes attempt to open a read-only file at the same time, they can be granted simultaneous access to the file. A process never needs to wait for a sharable resource. In general, however, we cannot prevent deadlocks by denying the mutual-exclusion condition, because some resources are intrinsically nonsharable. For example, a mutex lock cannot be simultaneously shared by several processes.

## 6.4.2 Hold and Wait

To ensure that the hold-and-wait condition never occurs in the system, we must guarantee that, whenever a process requests a resource, it does not hold any other resources. One protocol that we can use requires each process to request and be allocated all its resources before it begins execution. We can implement this provision by requiring that system calls requesting resources for a process precede all other system calls.

An alternative protocol allows a process to request resources only when it has none. A process may request some resources and use them. Before it can request any additional resources, it must release all the resources that it is currently allocated.

To illustrate the difference between these two protocols, we consider a process that copies data from a DVD drive to a file on disk, sorts the file, and then prints the results to a printer. If all resources must be requested at the beginning of the process, then the process must initially request the DVD drive, disk file, and printer. It will hold the printer for its entire execution, even though it needs the printer only at the end. The second method allows the process to request initially only the DVD drive and disk file. It copies from the DVD drive to the disk and then releases both the DVD drive and the disk file. The process must then request the disk file and the printer. After copying the disk file to the printer, it

releases these two resources and terminates.

Both these protocols have two main disadvantages. First, resource utiliza- tion may be low, since resources may be allocated but unused for a long period. In the example given, for instance, we can release the DVD drive and disk file, and then request the disk file and printer, only if we can be sure that our data will remain on the disk file. Otherwise, we must request all resources at the beginning for both protocols.

Second, starvation is possible. A process that needs several popular resources may have to wait indefinitely, because at least one of the resources that it needs is always allocated to some other process.

### 6.4.3 No Preemption

The third necessary condition for deadlocks is that there be no preemption of resources that have already been allocated. To ensure that this condition does not hold, we can use the following protocol. If a process is holding some resources and requests another resource that cannot be immediately allocated to it (that is, the process must wait), then all resources the process is currently holding are preempted. In other words, these resources are implicitly released. The preempted resources are added to the list of resources for which the process is waiting. The process will be restarted only when it can regain its old resources, as well as the new ones that it is requesting.

Alternatively, if a process requests some resources, we first check whether they are available. If they are, we allocate them. If they are not, we check whether they are allocated to some other process that is waiting for additional resources. If so, we preempt the desired resources from the waiting process and allocate them to the requesting process. If the resources are neither available nor held by a waiting process, the

requesting process must wait. While it is waiting, some of its resources may be preempted, but only if another process requests them. A process can be restarted only when it is allocated the new resources it is requesting and recovers any resources that were preempted while it was waiting.

This protocol is often applied to resources whose state can be easily saved and restored later, such as CPU registers and memory space. It cannot generally be applied to such resources as mutex locks and semaphores.

### 6.4.4 Circular Wait

The fourth and final condition for deadlocks is the circular-wait condition. One way to ensure that this condition never holds is to impose a total ordering of all resource types and to require that each process requests resources in an increasing order of enumeration.

## 6.5 Deadlock Avoidance

Deadlock-prevention algorithms, prevent deadlocks by limiting how requests can be made. The limits ensure that at least one of the necessary conditions for deadlock cannot occur. Possible side effects of preventing deadlocks by this method, however, are low device utilization and reduced system throughput.

An alternative method for avoiding deadlocks is to require additional information about how resources are to be requested. For example, in a system with one tape drive and one printer, the system might need to know that process $P$ will request first the tape drive and then the printer before releasing both resources, whereas process $Q$ will request first the printer and then the tape drive. With this knowledge of the complete sequence of requests and releases for each process, the system can decide for each request whether or not the process should wait in order to avoid a possible future deadlock. Each request requires that in making this decision

the system consider the resources currently available, the resources currently allocated to each process, and the future requests and releases of each process.

The various algorithms that use this approach differ in the amount and type of information required. The simplest and most useful model requires that each process declare the *maximum number* of resources of each type that it may need. Given this a priori information, it is possible to construct an algorithm that ensures that the system will never enter a deadlocked state. A deadlock-avoidance algorithm dynamically examines the resource-allocation state to ensure that a circular-wait condition can never exist. The resource-allocation *state* is defined by the number of available and allocated resources and the maximum demands of the processes.

## Review Questions

1) Suppose that a system is in an unsafe state. Show that it is possible for the processes to complete their execution without entering a deadlocked state.

2) Is it possible to have a deadlock involving only one single-threaded process? Explain your answer.

# REFRENCES

*1- Silberschatz, A., Galvin, P. B., & Gagne, G. (2012). Operating system concepts. [Sl].*

*2-* R. Rojas and U. Hashagen, The First Computers— History and Architectures, MIT Press (2000).