# OPERATING SYSTEM – [OS]    UNIT-I

**Introduction:**

Operating System is software that works as an interface between a user and the computer hardware. The primary objective of an operating system is to make computer system convenient to use and to utilize computer hardware in an efficient manner. The operating system performs the basic tasks such as receiving input from the keyboard, processing instructions and sending output to the screen.

The Software is the Non-Touchable Parts of the Computer, and Software's are those which are used for Performing an Operation So that Software's are just used for Making an Application but hardware's are those which are used for Performing an Operation. Operating system is software that is required in order to run application programs and utilities. It works as a bridge to perform better interaction between application programs and hardware of the computer. Various types of operating systems' are UNIX, MS-DOS, MS-Windows - 98/XP/Vista, WindowsNT/2000, OS/2 and Mac OS.

Operating system manages overall activities of a computer and the input/output devices attached to the computer. It is the first software you see when you turn on the computer, and the last software you see when the computer is turned off. It is the software that enables all the programs you use. At the simplest level, an operating system does two things:

The first, it manages the hardware and software resources of the computer system. These resources include the processor, memory, disk space, etc. The second, it provides a stable, consistent way for applications to deal with the hardware without having-to know all the details of the hardware.

The first task is very important i.e. managing the hardware and software resources, as various processes compete to each other for getting the CPU time and memory space to complete the task. In this regard; the operating system acts as a manager to allocate the available resources to 'satisfy the requirements of each process.

The second task i.e. providing a consistent application interface is especially important. A consistent application program interface (API) allows a user (or S/W developer) to write an application program on any computer and to run this program on another computer, even if the hardware configuration is different like as amount of memory, type of CPU or storage disk. It shields the user of the machine from the low-level details of the machine's operation and provides frequently needed facilities. When you turn on the computer, the operating system program is loaded into the main memory. This program is called the kernel. Once initialized, the system program is prepared to run the user programs and permits them to use the hardware efficiently. Windows 98/XP is an excellent example that supports different types of hardware configurations from thousands of vendors and accommodates thousands of different I/O devices like printers, disk drives, scanners and cameras.

Operating systems may be classified based on if multiple tasks can be performed simultaneously, and if the system can be used by multiple users. It can be termed as single-user or multiuser OS, and single-tasking or multi-tasking OS.A multi-user system must be multi-tasking. MS-DOS and Windows 3x are examples of single user operating system. Whereas UNIX is an example of multi-user and multitasking operating system.

**Classification of Operating systems:**

- **Multi-user**            : Allows two or more users to run programs at the same time. Some operating systems permit hundreds or even thousands of concurrent users.
- **Multiprocessing**       : Supports running a program on more than one CPU.
- **Multitasking**          : Allows more than one program to run concurrently.
- **Multithreading**        : Allows different parts of a single program to run concurrently.
- **Real time**             : Responds to input instantly. General-purpose operating systems, such as DOS and UNIX, are not real-time.

For Example if we want to Perform Some Paintings on the Screen, then we must use the Application Software as Paint and Hardware as a Mouse for Drawing an Object. But how the System knows what to do when Mouse Moves on the Screen and When the Mouse Draws a Line on the System so that Operating System is Necessary which Interact between or which Communicates with the Hardware and the Software.

**Characteristics or Functions of OS:**

For Better understanding you can see the Working of the Operating System. So we can say that the Operating System has the Following **Characteristics**:

- It boots the computer.
- Operating System is a Collection of Programs those are Responsible for the Execution of other Programs.
- Operating System is that which Responsible is for Controlling all the Input and Output Devices those are connected to the System.
- Operating System is that which Responsible is for Running all the Application Software's.
- Operating System is that which Provides Scheduling to the Various Processes Means Allocates the Memory to various Process those Wants to Execute.
- Operating System is that which provides the Communication between the user and the System.
- Operating System is Stored into the BIOS Means in the Basic Input and Output System means when a user Starts his System then this will Read all the instructions those are Necessary for Executing the System Means for Running the Operating System, Operating System Must be

Loaded into the Computer For this, this will use the Floppy or Hard Disks Which Stores the Operating System.

- It provides file management which refers to the way that the operating system manipulates, stores, retrieves and saves data.
- Error Handling is done by the operating system. It takes preventive measures whenever required to avoid errors.

**Most Popular Desktop Operating Systems:**

The three most popular types of operating systems for personal and business computing include Linux, Windows and Mac.

**Windows -** Microsoft Windows is a family of operating systems for personal and business computers. Windows dominates the personal computer world, offering a graphical user interface (GUI), virtual memory management, multitasking, and support for many peripheral devices.

**Mac -** Mac OS is the official name of the Apple Macintosh operating system. Mac OS features a graphical user interface (GUI) that utilizes windows, icons, and all applications that run on a Macintosh computer have a similar user interface.

**Linux -** Linux is a freely distributed open source operating system that runs on a number of hardware platforms. The Linux kernel was developed mainly by Linus Torvalds and it is based on Unix.

In the same way that a desktop OS controls your desktop or laptop computer, a mobile operating system is the software platform on top of which other programs can run on mobile devices, however, these systems are designed specifically to run on mobile devices such as mobile phones, smartphones, PDAs, tablet computers and other handhelds. The mobile OS is responsible for determining the functions and features available on your device, such as thumb wheel, keyboards, WAP, synchronization with applications, email, text messaging and more. The mobile OS will also determine which third-party applications (mobile apps) can be used on your device.

**What is an Operating System?**

A program that acts as an intermediary between a user of a computer and the computer hardware Operating system goals:

- Execute user programs and make solving user problems easier
- Make the computer system convenient to use
- Use the computer hardware in an efficient manner

**Operating System Definition**

- 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

- No universally accepted definition

- Everything a vendor ships when you order an operating system" is good approximation. But varies wildly

===================================================================

## Definition:

The operating system is the "one program running at all times on the computer"- called the **kernel**. (Along with the kernel, there are two other types of programs:

- **System programs**- which are associated with the operating system but are not necessarily part of the kernel, and

- **Application programs**- which include all programs not associated with the operation of the system.)
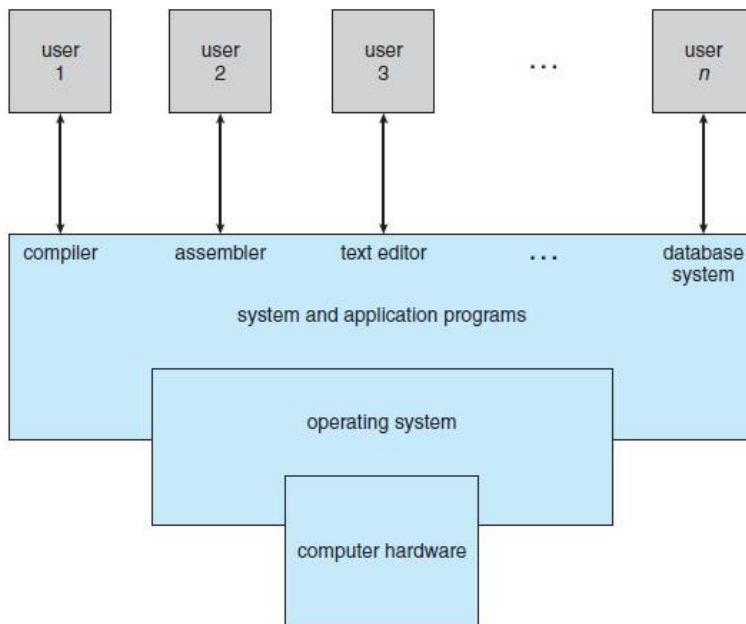
===================================================================

**Why should I study Operating Systems?**

- Need to understand interaction between the hardware and software

- Need to understand basic principles in the design of computer systems

- Efficient resource management, security, flexibility

- Because it enables you to do things that are difficult/impossible otherwise.

**OS challenges**

•Reliability
    -Does the system do what it was designed to do?
•Availability
    -What portion of the time is the system working?
    -Mean Time To Failure (MTTF), Mean Time to Repair
•Security
    -Can the system be compromised by an attacker?
•Privacy
    -Data is accessible only to authorized users
•Performance
    •Latency/response time -How long does an operation take to complete?
    •Throughput -How many operations can be done per unit of time?
    •Overhead -How much extra work is done by the OS?
    •Fairness -How equal is the performance received by different users?
    •Predictability -How consistent is the performance over time?

# WHAT OPERATING SYSTEMS DO

A computer system can be divided roughly into four components: the *hardware,* the *operating system,* the *application programs,* and the *users.*



Abstract view of the components of a computer system.

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.

We can also view a computer system as consisting of hardware, software, and data.

## Two Views of Operating System

1. User's View
2. System View

## User View:

The user's view of the computer varies according to the interface being used. Most computer users sit in front of a PC, consisting of a monitor, keyboard, mouse, and system unit. Such a system is designed for one user to monopolize its resources. The goal is to maximize the work (or play) that the user is performing. In this case, the operating system is designed mostly for **ease of use**, with some attention paid to performance and none paid to **resource utilization**—how various hardware and software resources are shared. Performance is, of course, important to the user; but such systems are optimized for the single-user experience rather than the requirements of multiple users.

In other cases, a user sits at a terminal connected to a **mainframe** or a **minicomputer**. Other users are accessing the same computer through other terminals. These users share resources and may exchange information. The operating system in such cases is designed to maximize resource utilization to assure

that all available CPU time, memory, and I/O are used efficiently and that no individual user takes more than her fair share.

In other cases, users sit at **workstations** connected to networks of other workstations and **servers**.

Recently, many varieties of mobile computers, such as smart phones and tablets, have come into fashion. The user interface for mobile computers generally features a **touch screen**, where the user interacts with the system by pressing and swiping fingers across the screen rather than using a physical keyboard and mouse.

**System View:**

From the computer's point of view, the operating system is the program most intimately involved with the hardware. In this context, we can view an operating system as a **resource allocator**.
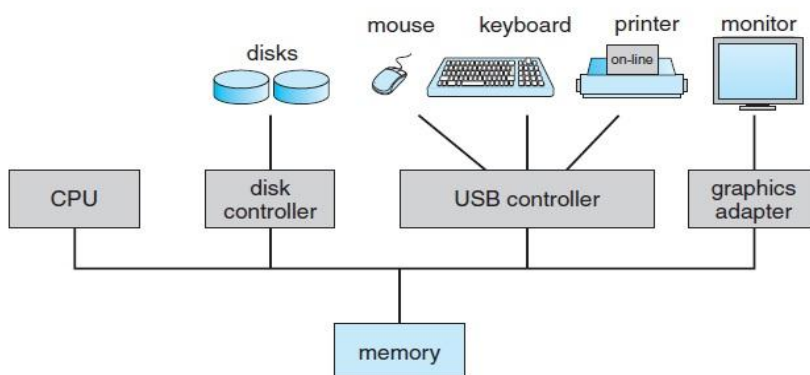
A computer system has many resources that may be required to solve a problem: CPU time, memory space, file-storage space, I/O devices, and so on. The operating system acts as the manager of these resources.

A slightly different view of an operating system emphasizes the need to control the various I/O devices and user programs. An operating system is a control program. A **control program** manages the execution of user programs to prevent errors and improper use of the computer. It is especially concerned with the operation and control of I/O devices.

## 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, so you can skim or skip it if you already understand the concepts.

**Computer-System Operation:**



A modern computer system.

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). Each device controller is in charge of a specific type of device (for

example, disk drives, audio devices, or video displays). The CPU and the device controllers can execute in parallel, competing for memory cycles. To ensure orderly access to the shared memory, a memory controller synchronizes access to the memory.

For a computer to start running—for instance, when it is powered up or rebooted—it needs to have an initial program to run.

This initial program, or **bootstrap program**, tends to be simple.

Typically, it is stored within the computer hardware in read-only memory (**ROM**) or electrically erasable programmable read-only memory (**EEPROM**), known by the general term **firmware**.

It initializes all aspects of the system, from CPU registers to device controllers to memory contents.

The bootstrap program must know how to load the operating system and how to start executing that system.

To accomplish this goal, the bootstrap program must locate the operating-system kernel and load it into memory.

Once the kernel is loaded and executing, it can start providing services to the system and its users. Some services are provided outside of the kernel, by system programs that are loaded into memory at boot time to become **system processes**, or **system daemons** that run the entire time the kernel is running. On UNIX, the first system process is "init," and it starts many other daemons. Once this phase is complete, the system is fully booted, and the system waits for some event to occur.
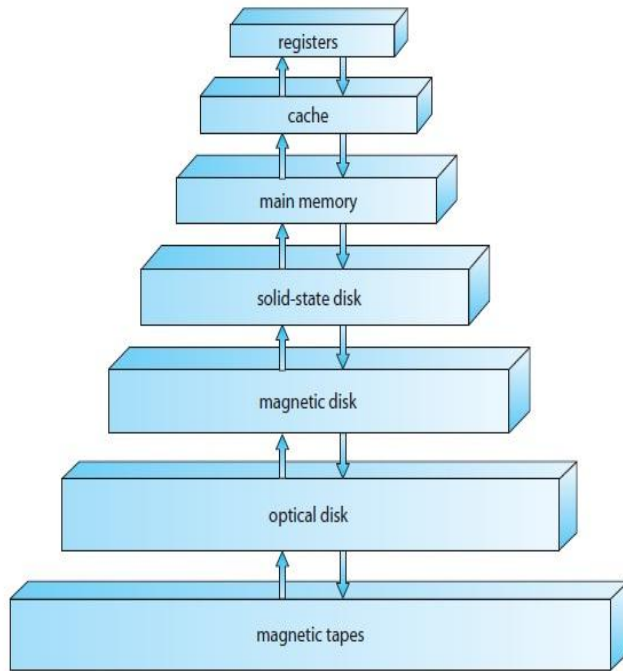
The occurrence of an event is usually signaled by an **interrupt** from either the hardware or the software. Hardware may trigger an interrupt at any time by sending a signal to the CPU, usually by way of the system bus. Software may trigger an interrupt by executing a special operation called a **system call** (also called a **monitor call**).

**Storage Structure:**

The CPU can load instructions only from memory, so any programs to run must be stored there. General-purpose computers run most of their programs from rewritable memory, called main memory (also called **random-access memory**, or **RAM**). Main memory commonly is implemented in a semiconductor technology called **dynamic random-access memory (DRAM)**. Ideally, we want the programs and data to reside in main memory permanently.

This arrangement usually is not possible for the following two reasons:

**1.** Main memory is usually too small to store all needed programs and data permanently.

**2.** Main memory is a **volatile** storage device that loses its contents when power is turned off or
   otherwise lost.

Storage-device hierarchy

Thus, most computer systems provide **secondary storage** as an extension of main memory. The main requirement for secondary storage is that it be able to hold large quantities of data permanently. The most common secondary-storage device is a **magnetic disk**, which provides storage for both programs and data.

The wide variety of storage systems can be organized in a hierarchy (Figure) according to speed and cost. The higher levels are expensive, but they are fast. As we move down the hierarchy, the cost per bit generally decreases, whereas the access time generally increases.

## Storage Definition and notation

- 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 1,0242 bytes;
    a **gigabyte**, or **GB**, is 1,0243 bytes;
    a **terabyte**, or **TB**, is 1,0244 bytes; and
    a **petabyte**, or **PB**, is 1,0245 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).

**I/O Structure:**

Storage is only one of many types of I/O devices within a computer. A large portion of operating system code is dedicated to managing I/O, both because of its importance to the reliability and performance of a system and because of the varying nature of the devices. Next, we provide an overview of I/O.

A general-purpose computer system consists of CPUs and multiple device controllers that are connected through a common bus. Each device controller is in charge of a specific type of device. Depending on the controller, more than one device may be attached. For instance, seven or more devices can be attached to the **small computer-systems interface (SCSI)** controller.

A device controller maintains some local buffer storage and a set of special-purpose registers. The device controller is responsible for moving the data between the peripheral devices that it controls and its local buffer storage. Typically, operating systems have a **device driver** for each device controller. This device driver understands the device controller and provides the rest of the operating system with a uniform interface to the device.

This form of interrupt-driven I/O is fine for moving small amounts of data but can produce high overhead when used for bulk data movement such as disk I/O. To solve this problem, **direct memory access (DMA)** is used.

## 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.

**Single-Processor Systems**

Until recently, most computer systems used a single processor. On a single processor system, there is one main CPU capable of executing a general-purpose instruction set, including instructions from user processes.
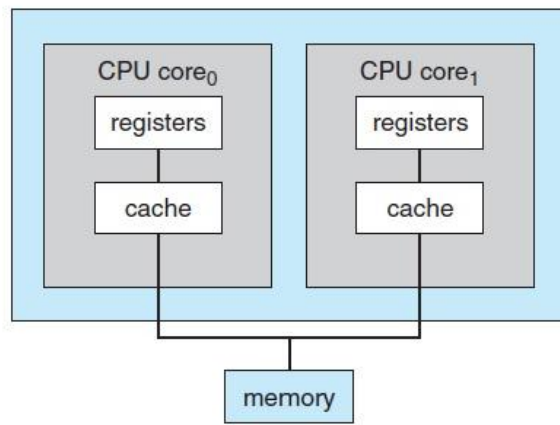
Almost all single processor systems have other special-purpose processors as well. They may come in the form of device-specific processors, such as disk, keyboard, and graphics controllers; or, on mainframes, they may come in the form of more general-purpose processors, such as I/O processors that move data rapidly among the components of the system.

All of these special-purpose processors run a limited instruction set and do not run user processes. Sometimes, they are managed by the operating system, in that the operating system sends them information about their next task and monitors their status.

For example, a disk-controller microprocessor receives a sequence of requests from the main CPU and implements its own disk queue and scheduling algorithm. This arrangement relieves the main CPU of the overhead of disk scheduling. PCs contain a microprocessor in the keyboard to convert the keystrokes into codes to be sent to the CPU.

In other systems or circumstances, special-purpose processors are low-level components built into the hardware. The operating system cannot communicate with these processors; they do their jobs autonomously. The use of special-purpose microprocessors is common and does not turn a single-processor system into a multiprocessor. If there is only one general-purpose CPU, then the system is a single-processor system.

**Multiprocessor Systems**



A dual-core design with two cores placed on the same chip.

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.
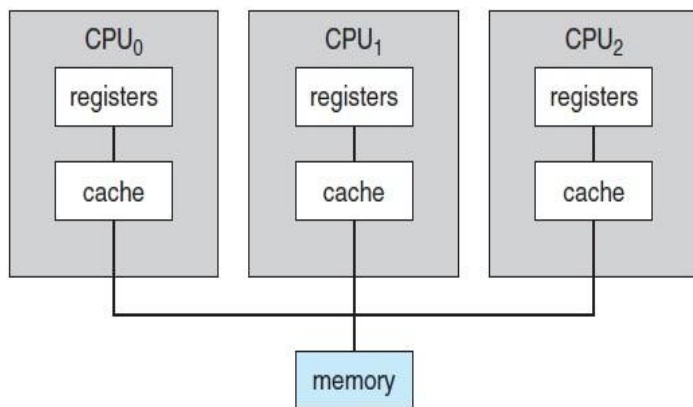
Multiprocessor systems have three main advantages:

**1. Increased throughput**. By increasing the number of processors, we expect to get more work done in less time. The speed-up ratio with $N$ processors is not $N$, however; rather, it is less than $N$. When multiple processors cooperate on a task, a certain amount of overhead is incurred in keeping all the parts working correctly. This overhead, plus contention for shared resources, lowers the expected gain from additional processors. Similarly, $N$ programmers working closely together do not produce $N$ times the amount of work a single programmer would produce.

**2. Economy of scale**. Multiprocessor systems can cost less than equivalent multiple single-processor systems, because they can share peripherals, mass storage, and power supplies. If several programs operate on the same set of data, it is cheaper to store those data on one disk and to have all the processors share them than to have many computers with local disks and many copies of the data.

**3. Increased reliability**. If functions can be distributed properly among several processors, then the failure of one processor will not halt the system, only slow it down. If we have ten processors and one fails, then each of the remaining nine processors can pick up a share of the work of the failed processor. Thus, the entire system runs only 10 percent slower, rather than failing altogether.

The multiple-processor systems in use today are of two types.

Some systems use **asymmetric multiprocessing**, in which each processor is assigned a specific task.

A *boss* processor controls the system; the other processors either look to the boss for instruction or have predefined tasks. This scheme defines a boss–worker relationship. The boss processor schedules and allocates work to the worker processors.

The most common systems use **symmetric multiprocessing (SMP)**, in which each processor performs all tasks within the operating system. SMP means that all processors are peers; no boss–worker relationship exists between processors.

Symmetric multiprocessing architecture.

**Clustered Systems:**

- Like parallel systems, clustered systems gather together multiple CPUs to accomplish computational work.
- Clustered systems differ from parallel systems, however, in that they are composed of two or more individual systems coupled together.
- The definition of the term clustered is **not concrete;** the general accepted definition is that clustered computers share storage and is closely linked via LAN networking.
- Clustering is usually performed to provide **high availability**.
- A layer of cluster software runs on the cluster nodes. Each node can monitor one or more of the others. If the monitored machine fails, the monitoring machine can take ownership of its storage, and restart the application(s) that were running on the failed machine. The failed machine can remain down, but the users and clients of the application would only see a brief interruption of service.

- **Asymmetric Clustering -** In this, one machine is in hot standby mode while the other is running the applications. The hot standby host (machine) does nothing but monitor the active server. If that server fails, the hot standby host becomes the active server.
- **Symmetric Clustering -** In this, two or more hosts are running applications, and they are monitoring each other. This mode is obviously more efficient, as it uses all of the available hardware.
- **Parallel Clustering -** Parallel clusters allow multiple hosts to access the same data on the shared storage. Because most operating systems lack support for this simultaneous data access by multiple hosts, parallel clusters are usually accomplished by special versions of software and special releases of applications.

Clustered technology is rapidly changing. Clustered system use and features should expand greatly as **Storage Area Networks (SANs)**. SANs allow easy attachment of multiple hosts to multiple storage units. Current clusters are usually limited to two or four hosts due to the complexity of connecting the hosts to shared storage.

## OPERATING-SYSTEM STRUCTURE

Now that we have discussed basic computer-system organization and architecture, we are ready to talk about operating systems.

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. There are, however, many commonalities, which we consider in this section.

One of the most important aspects of operating systems is the ability to **Multiprogram.** A single program cannot, in general, keep either the CPU or the I/O devices busy at all times. Single users frequently have multiple programs running.

**Multiprogramming** increases CPU utilization by organizing jobs (code and data) so that the CPU always has one to execute. The idea is as follows: The operating system keeps several jobs in memory simultaneously.

This idea is common in other life situations. A lawyer does not work for only one client at a time, for example. While one case is waiting to go to trial or have papers typed, the lawyer can work on another case. If he has enough clients, the lawyer will never be idle for lack of work. (Idle lawyers tend to become politicians, so there is a certain social value in keeping lawyers busy.) Multiprogrammed systems provide an environment in which the various system resources (for example, CPU, memory, and peripheral devices) are utilized effectively, but they do not provide for user interaction with the computer system.

**Time sharing** (or **multitasking**) is a logical extension of multiprogramming. In time-sharing systems, the CPU executes multiple jobs by switching among them, but the switches occur so frequently that the users can interact with each program while it is running.

Time sharing requires an **interactive** computer system, which provides direct communication between the user and the system. The user gives instructions to the operating system or to a program directly, using a input device such as a keyboard, mouse, touch pad, or touch screen, and waits for immediate results on an output device. Accordingly, the **response time** should be short—typically less than one second.

Each user has at least one separate program in memory. A program loaded into memory and executing is called a **process**. When a process executes, it typically executes for only a short time before it either finishes or needs to perform I/O. I/O may be interactive; that is, output goes to a display for the user, and input comes from a user keyboard, mouse, or other device. Since interactive I/O typically runs at "people speeds," it may take a long time to complete. Input, for example, may be bounded by the user's typing speed; seven characters per second is fast for people but incredibly slow for computers. Rather than let the CPU sit idle as this interactive input takes place, the operating system will rapidly switch the CPU to the program of some other user.

Time sharing and multiprogramming require that several jobs be kept simultaneously in memory. If several jobs are ready to be brought into memory, and if there is not enough room for all of them, then the system must choose among them. Making this decision involves **job scheduling.**

When the operating system selects a job from the job pool, it loads that job into memory for execution. Having several programs in memory at the same time requires some form of memory management. In addition, if several jobs are ready to run at the same time, the system must choose which job will run first. Making this decision is **CPU scheduling**.

Finally, running multiple jobs concurrently requires that their ability to affect one another be limited in all phases of the operating system; including process scheduling, disk storage and memory management.

In a time-sharing system, the operating system must ensure reasonable response time. This goal is sometimes accomplished through **swapping**, whereby processes are swapped in and out of main memory to the disk.

A more common method for ensuring reasonable response time is **virtual memory**, a technique that allows the execution of a process that is not completely in memory. The main advantage of the virtual-memory scheme is that it enables users to run programs that are larger than actual **physical memory**. Further, it abstracts main memory into a large, uniform array of storage, separating **logical memory** as viewed by the user from physical memory. This arrangement frees programmers from concern over memory-storage limitations.

# OPERATING-SYSTEM OPERATIONS

In modern operating systems are **interrupt driven**. If there are no processes to execute, no I/O devices to service, and no users to whom to respond, an operating system will sit quietly, waiting for something to happen. Events are almost always signaled by the occurrence of an interrupt or a trap.
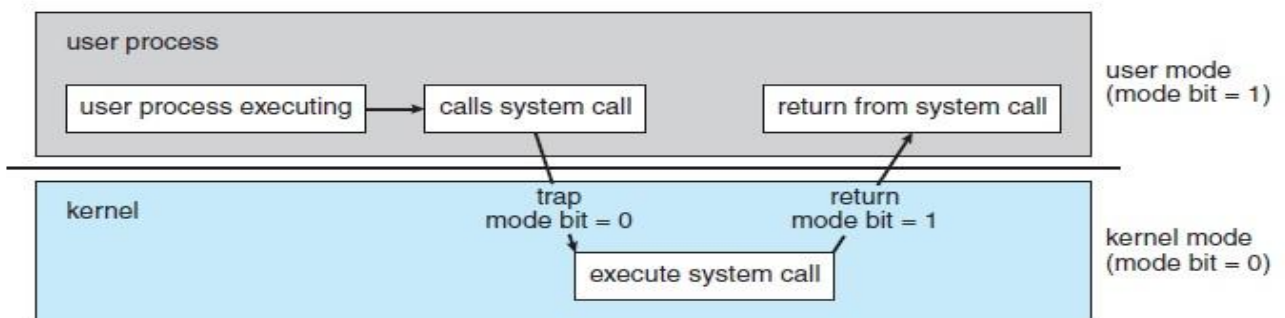
A **trap** (or an **exception**) is a software-generated interrupt caused either by an error (for example, division by zero or invalid memory access) or by a specific request from a user program that an operating-system service be performed. The interrupt-driven nature of an operating system defines that system's general structure.

For each type of interrupt, separate segments of code in the operating system determine what action should be taken. An interrupt service routine is provided to deal with the interrupt. Since the operating system and the users share the hardware and software resources of the computer system, we need to make sure that an error in a user program could cause problems only for the one program running. With sharing, many processes could be adversely affected by a bug in one program.

For example, if a process gets stuck in an infinite loop, this loop could prevent the correct operation of many other processes. More subtle errors can occur in a multiprogramming system, where one erroneous program might modify another program, the data of another program, or even the operating system itself. Without protection against these sorts of errors, either the computer must execute only one process at a time or all output must be suspect. A properly designed operating system must ensure that an incorrect (or malicious) program cannot cause other programs to execute incorrectly.

**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.



Transition from user to kernel mode.

At the very least, we need two separate *modes* of operation:

1. **User mode** and
2. **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 above Figure. 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.

**Timer:**

We must ensure that the operating system maintains control over the CPU. We cannot allow a user program to get stuck in an infinite loop or to fail to call system services and never return control to the operating system. To accomplish this goal, we can use a **timer**.

A timer can be set to interrupt the computer after a specified period. The period may be fixed (for example, 1/60 second) or variable (for example, from 1 millisecond to 1 second). A **variable timer** is generally implemented by a fixed-rate clock and a counter.

The operating system sets the counter. Every time the clock ticks, the counter is decremented. When the counter reaches 0, an interrupt occurs. For instance, a 10-bit counter with a 1-millisecond clock allows interrupts at intervals from 1 millisecond to 1,024 milliseconds, in steps of 1 millisecond. Before turning over control to the user, the operating system ensures that the timer is set to interrupt. If the timer interrupts, control transfers automatically to the operating system, which may treat the interrupt as a fatal error or may give the program more time. Clearly, instructions that modify the content of the timer are privileged.

# OPERATING SYSTEM MANAGEMENT TASKS

1. **Processor management** which involves putting the tasks into order and pairing them into manageable size before they go to the CPU.

2. **Memory management** which coordinates data to and from RAM (random-access memory) and determines the necessity for virtual memory.

3. **Device management** which provides interface between connected devices.

4. **Storage management** which directs permanent data storage.

5. **Production and Security.**

6. **Application** which allows standard communication between software and your computer.

7. **User interface** which allows you to communicate with your computer.

## Process Management:

A program does nothing unless its instructions are executed by a CPU.

A program in execution, as mentioned, is a process.

A time-shared user program such as a compiler is a process.

A word-processing program being run by an individual user on a PC is a process.

A system task, such as sending output to a printer, can also be a process .

A process needs certain resources-including CPU time, memory, files, and I/O devices-to accomplish its task

A process is the unit of work in a system. A system consists of a collection of processes, some of which are operating-system processes (those that execute system code) and the rest of which are user processes (those that execute user code). All these processes can potentially execute concurrently—by multiplexing on a single CPU, for example.

The operating system is responsible for the following activities in connection with process management:

• Scheduling processes and threads on the CPUs

• Creating and deleting both user and system processes

• Suspending and resuming processes

• Providing mechanisms for process synchronization

• Providing mechanisms for process communication

## Memory Management:

The main memory is generally the only large storage device that the CPU is able to address and access directly. For example, for the CPU to process data from disk, those data must first be transferred to main memory by CPU-generated I/O calls. In the same way, instructions must be in memory for the

CPU to execute them. For a program to be executed, it must be mapped to absolute addresses and loaded into memory. As the program executes, it accesses program instructions and data from memory by generating these absolute addresses. Eventually, the program terminates, its memory space is declared available, and the next program can be loaded and executed.

The operating system is responsible for the following activities in connection with memory management:

- Keeping track of which parts of memory are currently being used and who is using them
- Deciding which processes (or parts of processes) and data to move into and out of memory
- Allocating and deallocating memory space as needed

**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.

**File-System Management**:

File management is one of the most visible components of an operating system. Computers can store information on several different types of physical media. Magnetic disk, optical disk, and magnetic tape are the most common. Each of these media has its own characteristics and physical organization. Each medium is controlled by a device, such as a disk drive or tape drive, that also has its own unique characteristics. These properties include access speed, capacity, data-transfer rate, and access method (sequential or random).

A file is a collection of related information defined by its creator. Commonly, files represent programs (both source and object forms) and data. Data files may be numeric, alphabetic, alphanumeric, or binary. Files may be free-form (for example, text files), or they may be formatted rigidly (for example, fixed fields).

Clearly, the concept of a file is an extremely general one. The operating system implements the abstract concept of a file by managing mass-storage media, such as tapes and disks, and the devices that control them. In addition, files are normally organized into directories to make them easier to use. Finally, when multiple users have access to files, it may be desirable to control which user may access a file and how that user may access it (for example, read, write, append).

The operating system is responsible for the following activities in connection with file management:

- Creating and deleting files
- Creating and deleting directories to organize files
- Supporting primitives for manipulating files and directories
- Mapping files onto secondary storage
- Backing up files on stable (nonvolatile) storage media

**Mass-Storage Management:**

The 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.

**Caching:**

**Caching** is an important principle of computer systems. Here's how it works. Information is normally kept in some storage system (such as main memory). As it is used, it is copied into a faster storage system—the cache—on a temporary basis. When we need a particular piece of information, we first check whether it is in the cache. If it is, we use the information directly from the cache. If it is not, we use the information from the source, putting a copy in the cache under the assumption that we will need it again soon In a hierarchical storage structure, the same data may appear in different levels of the storage system. For example, suppose that an integer A that is to be incremented by 1 is located in file B, and file B resides on magnetic disk.

The increment operation proceeds by first issuing an I/O operation to copy the disk block on which A resides to main memory. This operation is followed by copying A to the cache and to an internal register. Thus, the copy of A appears in several places: on the magnetic disk, in main memory, in the cache, and in an internal register. Once the increment takes place in the internal

register, the value of A differs in the various storage systems. The value of A becomes the same only after the new value of A is written from the internal register back to the magnetic disk.

**I/O Systems:**

One of the purposes of an operating system is to hide the peculiarities of specific hardware devices from the user. For example, in UNIX, the peculiarities of I/O devices are hidden from the bulk of the operating system itself by the **I/O subsystem**.

The I/O subsystem consists of several components:

• A memory-management component that includes buffering, caching, and spooling

• A general device-driver interface

• Drivers for specific hardware devices

Only the device driver knows the peculiarities of the specific device to which it is assigned.

**Protection and Security:**

**Protection** – any mechanism for controlling access of processes or users to resources defined by the OS

**Security** – defense of the system against internal and external attacks

- Huge range, including denial-of-service, worms, viruses, identity theft, theft of service
- Systems generally first distinguish among users, to determine who can do what
- User identities (**user IDs**, security IDs) include name and associated number, one per user
- User ID then associated with all files, processes of that user to determine access control
- Group identifier (**group ID**) allows set of users to be defined and controls managed, then also associated with each process, file
- **Privilege escalation** allows user to change to effective ID with more rights

## COMPUTING ENVIRONMENTS

In computers, the term **environment** when unqualified usually refers to the combination of hardware and software in a **computer**. In this usage, the term platform is a synonym.

**Traditional Computing:**

Consider the "typical office environment." Just a few years ago, this environment consisted of PCs connected to a network, with servers providing file and print services. Remote access was awkward,

and portability was achieved by use of laptop computers. Terminals attached to mainframes were prevalent at many companies as well, with even fewer remote access and portability options.

The current trend is toward providing more ways to access these computing environments. Web technologies are stretching the boundaries of traditional computing. Companies establish portals, which provide web accessibility to their internal servers.

Network computers are essentially terminals that understand web-based computing. Handheld computers can synchronize with PCs to allow very portable use of company information. Handheld PDAs can also connect to wireless networks to use the company's web portal (as well as the myriad other web resources).

At home, most users had a single computer with a slow modem connection to the office, the Internet, or both. Today, network-connection speeds once available only at great cost are relatively inexpensive, giving home users more access to more data. These fast data connections are allowing home computers to serve up web pages and to run networks that include printers, client PCs, and servers. Some homes even have firewalls to protect their networks from security breaches. Those firewalls cost thousands of rupees a few years ago and did not even exist a decade ago. In the latter half of the previous century, computing resources were scarce.

**Mobile Computing:**
Mobile computing refers to computing on handheld smartphones and tablet computers. These devices share the distinguishing physical features of being portable and lightweight.

Historically, compared with desktop and laptop computers, mobile systems gave up screen size, memory capacity, and overall functionality in return for handheld mobile access to services such as e-mail and web browsing.

Today, mobile systems are used not only for e-mail and web browsing but also for playing music and video, reading digital books, taking photos, and recording high-definition video. Accordingly, tremendous growth continues in the wide range of applications that run on such devices.

Many developers are now designing applications that take advantage of the unique features of mobile devices, such as

global positioning system (GPS) chips,

accelerometers, and

gyroscopes.

The memory capacity and processing speed of mobile devices, however, are more limited than those of PCs.  Whereas a smartphone or tablet may have 64 GB in storage, it is not uncommon to find 1

TB in storage on a desktop computer. Similarly, becauseare smaller, are slower, and offer fewer processing cores than processors found on traditional desktop and laptop computers.

Two operating systems currently dominate mobile computing:

Apple iOS and Google Android. iOS was designed to run on Apple iPhone and iPad mobile devices. Android powers smartphones and tablet computers available from many manufacturers.

**Distributed Systems:**

A distributed system is a collection of physically separate, possibly heterogeneous, computer systems that are networked to provide users with access to the various resources that the system maintains. Access to a shared resource increases computation speed, functionality, data availability, and reliability.

Generally, systems contain a mix of the two modes—

for example FTP and NFS. The protocols that create a distributed system can greatly affect that system's utility and popularity.

A **network**, in the simplest terms, is a communication path between two or more systems. Distributed systems depend on networking for their functionality.
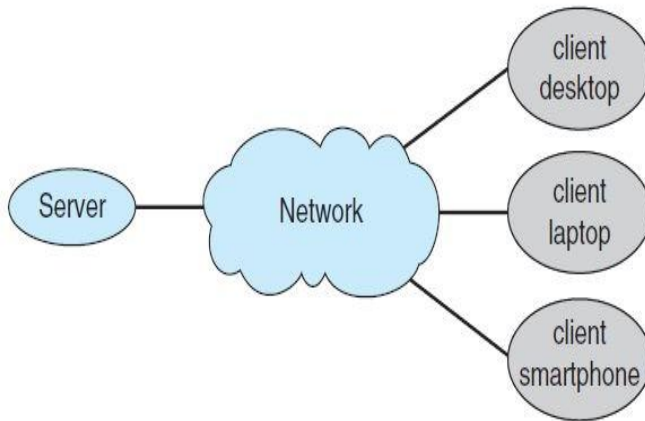
Networks vary by the protocols used, the distances between nodes, and the transport media. **TCP/IP** is the most common network protocol, and it provides the fundamental architecture of the Internet.

Networks are characterized based on the distances between their nodes.

- A **local-area network (LAN)** connects computers within a room, a building, or a campus.
- A **wide-area network (WAN)** usually links buildings, cities, or countries. A global company may have a WAN to connect its offices worldwide, for example. These networks may run one protocol or several protocols.
- The continuing advent of new technologies brings about new forms of networks. For example, a **metropolitan-area network (MAN)** could link buildings within a city.
- BlueTooth device use wireless technology to communicate over a distance of several feet, in essence creating a **personal-area network (PAN)** between a phone and a headset or a smartphone and a desktop computer.

**Client-Server Computing:**

As PCs have become faster, more powerful, and cheaper, designers have shifted away from centralized system architecture. Terminals connected to centralized systems are now being supplanted by PCs and mobile devices.
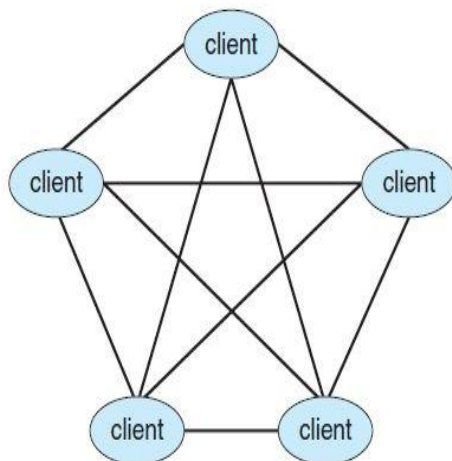
Correspondingly, user-interface functionality once handled directly by centralized systems is increasingly being handled by PCs, quite often through a web interface. As a result, many of today's systems act as **server systems** to satisfy requests generated by **client systems**.

This form of specialized distributed system, called a **client–server** system, has the general structure depicted in Figure

Server systems can be broadly categorized as compute servers and file servers:

General structure of a client–server system.

- The **compute-server system** provides an interface to which a client can send a request to perform an action (for example, read data). In response, the server executes the action and sends the results to the client. A server running a database that responds to client requests for data is an example of such a system.
- The **file-server system** provides a file-system interface where clients can create, update, read, and delete files. An example of such a system is a web server that delivers files to clients running web browsers.

**Peer-to-Peer Computing:**

Another structure for a distributed system is the peer-to-peer (P2P) system model. In this model, clients and servers are not distinguished from one another; instead, all nodes within the system are considered peers, and each may act as either a client or a server, depending on whether it is requesting or providing a service.

Peer-to-peer systems offer an advantage over traditional client-server systems. In a client-server system, the server is a bottleneck; but in a

Peer-to-peer system with no centralized service.

peer-to-peer system, services can be provided by several nodes distributed throughout the network. To participate in a peer-to-peer system, a node must first join the network of peers. Once a node has joined the network, it can begin providing services to—and requesting services from—other nodes in the network. Determining what services are available is accomplished in one of two general ways:

• When a node joins a network, it registers its service with a centralized lookup service on the network. Any node desiring a specific service first contacts this centralized lookup service to determine which node provides the service. The remainder of the communication takes place between the client and the service provider.

• A peer acting as a client must first discover what node provides a desired service by broadcasting a request for the service to all other nodes in the network. The node (or nodes) providing that service responds to the peer making the request. To support this approach, a discovery protocol must be provided that allows peers to discover services provided by other peers in the network.

Peer-to-peer networks gained widespread popularity in the late 1990s with several file-sharing services, such as Napster and Gnutella that enable peers to exchange files with one another.

**Web-Based Computing**

The Web has become ubiquitous, leading to more access by a wider variety of devices than was dreamt of a few years ago. PCs are still the most prevalent access devices, with workstations, handheld PDAs, and even cell phones also providing access.
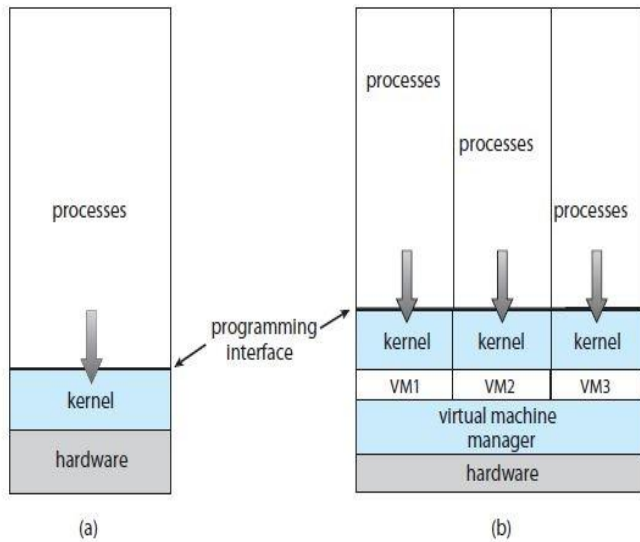
Web computing has increased the emphasis on networking. Devices that were not previously networked now include wired or wireless access. Devices that were networked now have faster network connectivity, provided by either improved networking technology, optimized network implementation code, or both.

The implementation of web-based computing has given rise to new categories of devices, such as load balancers, which distribute network connections among a pool of similar servers.

Operating systems like Windows 95, which acted as web clients, have evolved into Linux and Windows XP, which can act as web servers as well as clients. Generally, the Web has increased the complexity of devices, because their users require them to be web-enabled.


**Virtualization:**

Virtualization is a technology that allows operating systems to run as applications within other operating systems. But the virtualization industry is vast and growing, which is a testament to its utility and importance. Broadly speaking,
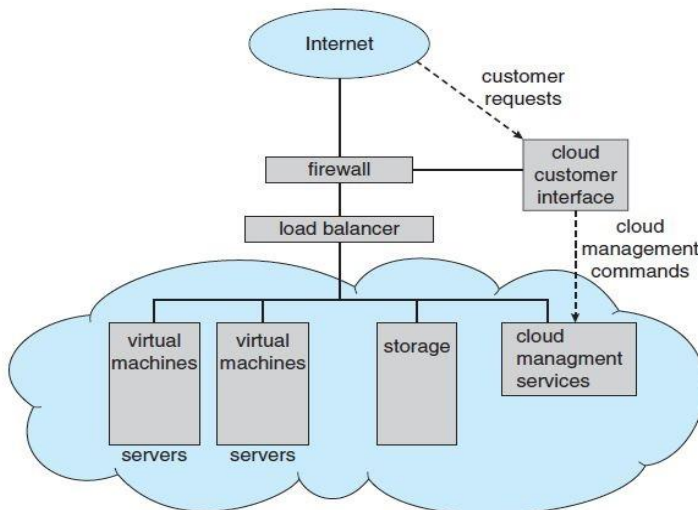
It is also called OS-level virtualization is a type of virtualization technology which work on OS layer. Here the kernel of an OS allows more than one isolated user-space instances to exist. Such instances are called containers/software containers or virtualization engines. In other words, OS kernel will run a single operating system & provide that operating system's functionality to replicate on each of the isolated partitions.

Virtualization is one member of a class of software that also includes emulation.

VMware.

**Emulation** is used when the source CPU type is different from the target CPU type.

For example, when Apple switched from the IBM Power CPU to the Intel x86 CPU for its desktop and laptop computers, it included an emulation facility called "Rosetta," which allowed applications compiled for the IBM CPU to run on the Intel CPU. That same concept can be extended to allow an entire operating system written for one platform to runon another. Emulation comes at a heavy price, however. Every machine-level instruction that runs natively on the source system must be translated to the equivalent function on the target system, frequently resulting in several target instructions. If the source and target CPUs have similar performance levels, the emulated code can run much slower than the native code.

## Cloud Computing:

**Cloud computing** is a type of computing that delivers computing, storage, and even applications as a service across a network. In some ways, it's a logical extension of virtualization, because it uses virtualization as a base for its functionality.

For example, the Amazon Elastic Compute Cloud **(EC2)** facility has thousands of servers, millions of virtual

Cloud computing.

machines, and peta bytes of storage available for use by anyone on the Internet. Users pay per month based on how much of those resources they use.

There are actually many types of cloud computing, including the following:

- **Public cloud**—a cloud available via the Internet to anyone willing to pay for the services
- **Private cloud**—a cloud run by a company for that company's own use
- **Hybrid cloud**—a cloud that includes both public and private cloud components
- Software as a service **(SaaS)**—one or more applications (such as word processors or spreadsheets) available via the Internet
- Platform as a service **(PaaS)**—a software stack ready for application use via the Internet (for example, a database server)
- Infrastructure as a service **(IaaS)**—servers or storage available over the Internet (for example, storage available for making backup copies of production data)

These cloud-computing types are not discrete, as a cloud computing environment may provide a combination of several types. For example, an organization may provide both SaaS and IaaS as a publicly available service.

**Real-Time Embedded Systems:**

Embedded computers are the most prevalent form of computers in existence.

These devices are found everywhere, from car engines and manufacturing robots to DVDs and microwave ovens. They tend to have very specific tasks.

The systems they run on are usually primitive, and so the operating systems provide limited features. Usually, they have little or no user interface, preferring to spend their time monitoring and managing hardware devices, such as automobile engines and robotic arms.

Embedded systems almost always run **real-time operating systems**. A real-time system is used when rigid time requirements have been placed on the operation of a processor or the flow of data; thus, it is often used as a control device in a dedicated application. Sensors bring data to the computer.

The computer must analyze the data and possibly adjust controls to modify the sensor inputs. Systems that control scientific experiments, medical imaging systems, industrial control systems, and certain display systems are real-time systems. Some automobile-engine fuel-injection systems, home-appliance controllers, and weapon systems are also real-time systems.

# HISTORY OF OPERATING SYSTEM

- Earliest computers had no Operating System
    - Applications loaded manually
    - Users were experts on the hardware
- First *System Software* was libraries of code to manage devices.
- This grew to batch processing systems, where some focused on application programming and some on systems programming.

**Batch Processing:**



A typical computer in the 1960s and 70s was a large machine. Its processing was managed by a human operator. The operator would organize various jobs from multiple users into batches. Human operators would organize jobs into batches

**Time Sharing:**

A timesharing system allows multiple users to interact with a computer at the same time.

Multiprogramming allowed multiple processes to be active at once, which gave rise to the ability for programmers to interact with the computer system directly, while still sharing its resources. In a timesharing system, each user has his or her own virtual machine, in which all system resources are (in effect) available for use.

Fun Pictures





first computer *bug*, a moth          The IBM 650 Magnetic Drum Data Processing System Machine

Cray I supercomputer, introduced in 1976

**Current Operating System Research Topics**

### Symmetric multiprocessing

Allows for several CPUs to process multiple jobs at the same time. CPUs are independent of one another, but each has access to the operating system.

### Asymmetric multiprocessing

Some operating systems functions are assigned to subordinate processors, which take their instructions from the main CPU.

### Distributed processing

Processors are placed at remote locations and are connected to each other via telecom devices. Different from symmetric multiprocessing systems as they do not share memory. Computations can be dispersed among several processors.

## EVOLUTION OF OPERATING SYSTEMS

The evolution of operating systems is directly dependent to the development of computer systems and how users use them. Here is a quick tour of computing systems through the past fifty years in the timeline.

### Early Evolution

- 1945: ENIAC, Moore School of Engineering, University of Pennsylvania.
- 1949: EDSAC and EDVAC
- 1949 BINAC - a successor to the ENIAC
- 1951: UNIVAC by Remington
- 1952: IBM 701
- 1956: The interrupt
- 1954-1957: FORTRAN was developed

### Operating Systems by the late 1950s

By the late 1950s Operating systems were well improved and started supporting following usages :

- It was able to Single stream batch processing
- It could use Common, standardized, input/output routines for device access
- Program transition capabilities to reduce the overhead of starting a new job was added
- Error recovery to clean up after a job terminated abnormally was added.
- Job control languages that allowed users to specify the job definition and resource requirements were made possible.

### Operating Systems In 1960s

- 1961: The dawn of minicomputers
- 1962 Compatible Time-Sharing System (CTSS) from MIT
- 1963 Burroughs Master Control Program (MCP) for the B5000 system
- 1964: IBM System/360
- 1960s: Disks become mainstream
- 1966: Minicomputers get cheaper, more powerful, and really useful
- 1967-1968: The mouse
- 1964 and onward: Multics
- 1969: The UNIX Time-Sharing System from Bell Telephone Laboratories

### Supported OS Features by 1970s

- Multi User and Multi tasking was introduced.
- Dynamic address translation hardware and Virtual machines came into picture.
- Modular architectures came into existence.
- Personal, interactive systems came into existence.

### Accomplishments after 1970

- 1971: Intel announces the microprocessor
- 1972: IBM comes out with VM: the Virtual Machine Operating System
- 1973: UNIX 4th Edition is published
- 1973: Ethernet

- 1974 The Personal Computer Age begins

- 1974: Gates and Allen wrote BASIC for the Altair

- 1976: Apple II

- August 12, 1981: IBM introduces the IBM PC

- 1983 Microsoft begins work on MS-Windows

- 1984 Apple Macintosh comes out

- 1990 Microsoft Windows 3.0 comes out

- 1991 GNU/Linux

- 1992 The first Windows virus comes out

- 1993 Windows NT

- 2007: iOS

- 2008: Android OS


And as the research and development work still goes on, with new operating systems being developed and existing ones getting improved and modified to enhance the overall user experience, making operating systems fast and efficient like never before.

# OPERATING SYSTEM SERVICES

An Operating System provides services to both the users and to the programs.

- It provides programs an environment to execute.
- It provides users the services to execute the programs in a convenient manner.

The specific services provided, of course, differ from one operating system to another, but we can identify common classes. These operating system services are provided for the convenience of the programmer, to make the programming task.



A view of operating system services.

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

- User Interface
- Program Execution
- I/O operations
- File System manipulation
- Communication
- Error Detection
- Resource Allocation
- Accounting
- Protection and Security

One set of operating system services provides functions that are helpful to the user.

**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 Operation:**

An I/O subsystem comprises of I/O devices and their corresponding driver software. Drivers hide the peculiarities of specific hardware devices from the users.

An Operating System manages the communication between user and device drivers.

- I/O operation means read or write operation with any file or any specific I/O device.
- Operating system provides the access to the required I/O device when required.

**File system manipulation:**

A file represents a collection of related information. Computers can store files on the disk (secondary storage), for long-term storage purpose. Examples of storage media include magnetic tape, magnetic disk and optical disk drives like CD, DVD. Each of these media has its own properties like speed, capacity, data transfer rate and data access methods.

A file system is normally organized into directories for easy navigation and usage. These directories may contain files and other directions.

Following are the major activities of an operating system with respect to file management −

- Program needs to read a file or write a file.
- The operating system gives the permission to the program for operation on file.
- Permission varies from read-only, read-write, denied and so on.
- Operating System provides an interface to the user to create/delete files.
- Operating System provides an interface to the user to create/delete directories.
- Operating System provides an interface to create the backup of file system.

**Communication:**

In case of distributed systems which are a collection of processors that do not share memory, peripheral devices, or a clock, the operating system manages communications between all the processes. Multiple processes communicate with one another through communication lines in the network.

The OS handles routing and connection strategies, and the problems of contention and security. Following are the major activities of an operating system with respect to communication −

- Two processes often require data to be transferred between them
- Both the processes can be on one computer or on different computers, but are connected through a computer network.
- Communication may be implemented by two methods, either by Shared Memory or by Message Passing.

**Error detection:**

Errors can occur anytime and anywhere. An error may occur in CPU, in I/O devices or in the memory hardware.

Following are the major activities of an operating system with respect to error handling −

- The OS constantly checks for possible errors.
- The OS takes an appropriate action to ensure correct and consistent computing.

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**

In case of multi-user or multi-tasking environment, resources such as main memory, CPU cycles and files storage are to be allocated to each user or job.

Following are the major activities of an operating system with respect to resource management −

- The OS manages all kinds of resources using schedulers.
- CPU scheduling algorithms are used for better utilization of CPU.

**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:**

Considering a computer system having multiple users and concurrent execution of multiple processes, the various processes must be protected from each other's activities.

Protection refers to a mechanism or a way to control the access of programs, processes, or users to the resources defined by a computer system.

Following are the major activities of an operating system with respect to protection −

- The OS ensures that all access to system resources is controlled.
- The OS ensures that external I/O devices are protected from invalid access attempts.
- The OS provides authentication features for each user by means of passwords.

# 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- a GUI -Gra**phical User Interfaces.**

## Command Line Interpreter:

Some operating systems include the command interpreter in the kernel. Some, such as the popular Windows and Linux operating systems, use the command interpreter as a special program that runs when a user logs on or a job is initiated.

- In Windows, this is the MS-DOS prompt.
- Linux has more options. The command interpreter in Linux is known as a shell. The most commonly used shell is the Bash shell, but others such as the Korn shell, C shell, and Bourne shell exist. Most shells provide similar functionality, personal preference usually dictates which shell is best.
- The main function of the command utility is to receive and execute the next user generated command.
- Many commands are intended to manipulate files.
- Operating systems such as UNIX implements commands through system programs. Often these programs are stored as text files, which allow programmers to add additional functionality to the utility.

Thus, the UNIX command to delete a file
<div align="center">

**Root  rm file.txt**

</div>
would search for a file called rm, load the file into memory, and execute it with the parameter file.txt.

## Graphical User Interfaces – GUI:

Here, rather than entering commands directly via a command-line interface, users employ a mouse-based window and- menu system characterized by a **desktop** metaphor.

The user moves the mouse to position its pointer on images, or **icons**, on the screen (the desktop) that represent programs, files, directories, and system functions. Depending on the mouse pointer's location, clicking a button on the mouse can invoke a program, select a file or directory—known as a **folder**—or pull down a menu that contains commands.

Because a mouse is impractical for most mobile systems, smartphones and handheld tablet computers typically use a touchscreen interface. Here, users interact by making **gestures** on the touchscreen—for example, pressing and swiping fingers across the screen

## Choice of Interface:

The choice of whether to use a command-line or GUI interface is mostly one of personal preference.

**System administrators** who manage computers and power users who have deep knowledge of a system frequently use the command-line interface. For them, it is more efficient, giving them faster access to the activities they need to perform.

Indeed, on some systems, only a subset of system functions is available via the GUI, leaving the less common tasks to those who are command-line knowledgeable. Further, command line interfaces usually make repetitive tasks easier, in part because they have their own programmability.

For example, if a frequent task requires a set of command-line steps, those steps can be recorded into a file, and that file can be run just like a program.

A GUI provides a mouse-based windows and menu system as an interface. Users of Windows are more likely to use the GUI rather than the command line interface of MS-DOS, while UNIX users generally prefer using the command line interface of the shell rather than the GUI.

## SYSTEM CALLS

**System calls** provide an interface to the services made available by an operating system.

Application developers often do not have direct access to the system calls, but can access them through an application programming interface (API). The functions that are included in the API invoke the actual system calls. By using the API, certain benefits can be gained:

- Portability: as long a system supports an API, any program using that API can compile and run.
- Ease of Use: using the API can be significantly easier than using the actual system call.

These calls are generally available as routines written in C and C++, although certain low-level tasks (for example, tasks where hardware must be accessed directly) may have to be written using assembly-language instructions.

Before we discuss how an operating system makes system calls available, let's first use an example to illustrate how system calls are used: writing a simple program to read data from one file and copy them to another file.

The first input that the program will need is the names of the two files: the input file and the output file. These names can be specified in many ways, depending on the operating-system design. One approach is for the program to ask the user for the names.

In an interactive system, this approach will require a sequence of system calls, first to write a prompting message on the screen and then to read from the keyboard the characters that define the two files.

On mouse-based and icon-based systems, a menu of file names is usually displayed in a window. The user can then use the mouse to select the source name, and a window can be opened for the destination name to be specified. This sequence requires many I/O system calls.

Once the two file names have been obtained, the program must open the input file and create the output file. Each of these operations requires another system call. Possible error conditions for each operation can require additional system calls. When the program tries to open the input file, for example, it may find that there is no file of that name or that the file is protected against access.



Example of how system calls are used.

In these cases, the program should print a message on the console (another sequence of system calls) and then terminate abnormally (another system call). If the input file exists, then we must create a new output file. We may find that there is already an output file with the same name. This situation may cause the program to abort (a system call), or we may delete the existing file (another system call) and create a new one (yet another system call). Another option, in an interactive system, is to ask the user (via a sequence of system calls to output the prompting message and to read the response from the terminal) whether to replace the existing file or to abort the program.

When both files are set up, we enter a loop that reads from the input file (a system call) and writes to the output file (another system call). Each read and write must return status information regarding various possible error conditions.

On input, the program may find that the end of the file has been reached or that there was a hardware failure in the read (such as a parity error). The write operation may encounter various errors, depending on the output device (for example, no more disk space).

Finally, after the entire file is copied, the program may close both files (another system call), write a message to the console or window (more system calls), and finally terminate normally (the final system call). This system-call sequence is shown in above Figure.

Frequently, systems execute thousands of system calls per second. Most programmers never see this level of detail, however. Typically, application developers design programs according to an **application programming interface (API)**.

- Programming interface to the services provided by the OS
- Typically written in a high-level language (C or C++)
- Mostly accessed by programs via a high-level Application Program Interface (API)
- Three most common APIs are
  1. Win32 API for Windows,
  2. POSIX API (all versions of UNIX, Linux, and Mac OS X), and
  3. Java API for the Java virtual machine (JVM)



The handling of a user application invoking the open () system call.

Thus, most of the details of the operating-system interface are hidden from the programmer by the API and are managed by the run-time support library. The relationship between an API, the system-call interface, and the operating system is shown in Figure, which illustrates how the operating system handles a user application invoking the open() system call.

**System calls Parameters:**

Three general methods exist for passing parameters to the OS:



1. Parameters can be passed in registers.

2. When there are more parameters than registers, parameters can be stored in a block and the block address can be passed as a parameter to a register.

3. Parameters can also be pushed on or popped off the stack by the operating system.

# 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.

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

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

- request device, release device
- read, write, reposition
- get device attributes,
- set device attributes
- logically attach or detach devices

Communication

- create, delete comm. connection
- send, receive messages
- transfer status information
- attach or detach remote devices

**Process Control:**

A running program needs to be able to stop execution either normally or abnormally. When execution is stopped abnormally, often a dump of memory is taken and can be examined with a debugger.

**File Management:**

Some common system calls are *create*, *delete*, *read*, *write*, *reposition*, or *close*. Also, there is a need to determine the file attributes – *get* and *set* file attribute. Many times the OS provides an API to make these system calls.

**Device Management:**

Process usually requires several resources to execute, if these resources are available, they will be granted and control returned to the user process. These resources are also thought of as devices. Some are physical, such as a video card, and others are abstract, such as a file.

User programs *request* the device, and when finished they *release* the device. Similar to files, we can *read*, *write*, and *reposition* the device.

**Information Management:**

Some system calls exist purely for transferring information between the user program and the operating system. An example of this is *time*, or *date*.

The OS also keeps information about all its processes and provides system calls to report this information.

**Communication:**

There are two models of interprocess communication, the message-passing model and the shared memory model.

- Message-passing uses a common mailbox to pass messages between processes.
- Shared memory use certain system calls to create and gain access to create and gain access to regions of memory owned by other processes. The two processes exchange information by reading and writing in the shared data.

**Protection**

Protection provides a mechanism for controlling access to the resources provided by a computer system. Historically, protection was a concern only on multiprogrammed computer systems with several users. However, with the advent of networking and the Internet, all computer systems, from servers to mobile handheld devices, must be concerned with protection.

Typically, system calls providing protection include set permission() and get permission(), which manipulate the permission settings of resources such as files and disks. The allow user()and deny user()system calls specify whether particular users can—or cannot—be allowed access to certain resources.

**Example of Standard C Library:**



- The standard C library provides a portion of the system call interface for many version of UNIX and Linux.
- As an example, let's assume a C program invokes printf() statement.
- The C library intercepts this call and invokes the necessary system call(s) in the OS.
- The C library takes the value returned by write() and passes it back to the user program

## OPERATING-SYSTEM DESIGN AND IMPLEMENTATION

In this section, we discuss problems we face in designing and implementing an operating system. There are, of course, no complete solutions to such problems, but there are approaches that have proved successful.

**Design Goals**

The first problem in designing a system is to define goals and specifications. At the highest level, the design of the system will be affected by the choice of hardware and the type of system: batch, time sharing, single user, multiuser, distributed, real time, or general purpose.

Beyond this highest design level, the requirements may be much harder to specify.

The requirements can, however, be divided into two basic groups:

**User goals** and

**System goals**.

Users want certain obvious properties in a system. The system should be convenient to use, easy to learn and to use, reliable, safe, and fast. Of course, these specifications are not particularly useful in the system design, since there is no general agreement on how to achieve them. A similar set of requirements can be defined by those people who must design, create, maintain, and operate the system. The system

should be easy to design, implement, and maintain; and it should be flexible, reliable, error free, and efficient. Again, these requirements are vague and may be interpreted in various ways.

**Mechanisms and Policies**

One important principle is the separation of **policy** from **mechanism**.

Mechanisms determine *how* to do something; policies determine *what* will be done.

For example, the timer construct is a mechanism for ensuring CPU protection, but deciding how long the timer is to be set for a particular user is a policy decision. The separation of policy and mechanism is important for flexibility.

Policies are likely to change across places or over time. In the worst case, each change in policy would require a change in the underlying mechanism. A general mechanism insensitive to changes in policy would be more desirable. A change in policy would then require redefinition of only certain parameters of the system. For instance, consider a mechanism for giving priority to certain types of programs over others. If the mechanism is properly separated from policy, it can be used either to support a policy decision that I/O-intensive programs should have priority over CPU-intensive ones or to support the opposite policy.

Policy decisions are important for all resource allocation. Whenever it is necessary to decide whether or not to allocate a resource, a policy decision must be made. Whenever the question is *how* rather than *what*, it is a mechanism that must be determined.

**Implementation**

Once an operating system is designed, it must be implemented. Because operating systems are collections of many programs, written by many people over a long period of time, it is difficult to make general statements about how they are implemented.

Early operating systems were written in assembly language. Now, although some operating systems are still written in assembly language, most are written in a higher-level language such as C or an even higher-level language such as C++. Actually, an operating system can be written in more than one language.

The lowest levels of the kernel might be assembly language. Higher-level routines might be in C, and system programs might be in C or C++, in interpreted scripting languages like PERL or Python, or in shell scripts. In fact, a given Linux distribution probably includes programs written in all of those languages.

# OPERATING-SYSTEM STRUCTURE

A system as large and complex as a modern operating system must be engineered carefully if it is to function properly and be modified easily. A common approach is to partition the task into small components, or modules, rather than have one **monolithic** system. Each of these modules should be a well-defined portion of the system, with carefully defined inputs, outputs, and functions.
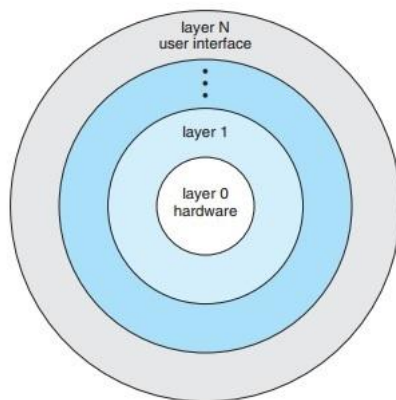
## Simple Structure

Many operating systems do not have well-defined structures. Frequently, such systems started as small, simple, and limited systems and then grew beyond their original scope. MS-DOS is an example of such a system. It was originally designed and implemented by a few people who had no idea that it would

become so popular. It was written to provide the most functionality in the least space, so it was not carefully divided into modules. Figure shows its structure.



MS-DOS layer structure.

In MS-DOS, the interfaces and levels of functionality are not well separated. For instance, application programs are able to access the basic I/O routines to write directly to the display and disk drives. Such freedom leaves MS-DOS vulnerable to errant (or malicious) programs, causing entire system crashes when user programs fail. Of course, MS-DOS was also limited by the hardware of its era. Because the Intel 8088 for which it was written provides no dual mode and no hardware protection, the designers of MS-DOS had no choice but to leave the base hardware accessible.

## Layered Approach



A layered operating system.

A system can be made modular in many ways. One method is the **layered approach**, in which the operating system is broken into a number of layers (levels). The bottom layer (layer 0) is the hardware; the highest (layer *N*) is the user interface. This layering structure is depicted in Figure.

An operating-system layer is an implementation of an abstract object made up of data and the operations that can manipulate those data. A typical operating-system layer—say, layer *M*—consists of data structures
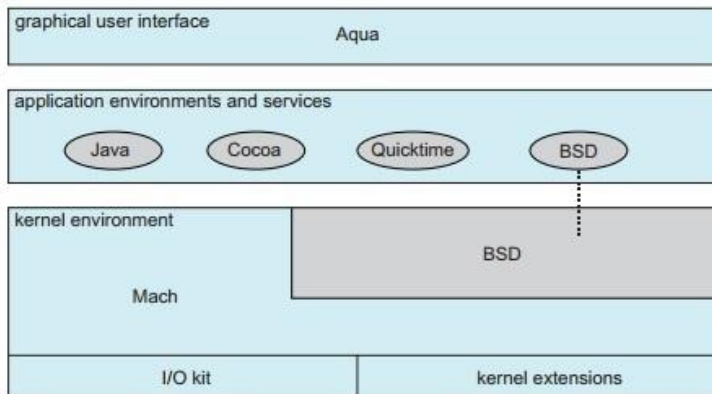
and a set of routines that can be invoked by higher-level layers. Layer *M,* in turn, can invoke operations on lower-level layers.

The main advantage of the layered approach is simplicity of construction and debugging.

The major difficulty with the layered approach involves appropriately defining the various layers.

**Microkernels**

In a **microkernel** (also known as µ-kernel) is the near-minimum amount of software that can provide the mechanisms needed to implement an**operating system** (**OS**). These mechanisms include low-level address space management, thread management, and inter-process communication (IPC).

We have already seen that as UNIX expanded, the kernel became large and difficult to manage. In the mid-1980s, researchers at Carnegie Mellon University developed an operating system called **Mach** that modularized the kernel using the **microkernel** approach.

This method structures the operating system by removing all nonessential components from the kernel and implementing them as system and user-level programs. The result is a smaller kernel. There is little consensus regarding which services should remain in the kernel and which should be implemented in user space.

Typically, however, microkernels provide minimal process and memory management, in addition to a communication facility. Figure illustrates the architecture of a typical microkernel.

The main function of the microkernel is to provide communication between the client program and the various services that are also running in user space. Communication is provided through **message passing,**

One benefit of the microkernel approach is that it makes extending the operating system easier.

**Example:** The Mac OS-X kernel (also known as **Darwin**) is also partly based on the Mach microkernel. Another example is QNX, a real-time operating system for embedded systems.



The Mac OS X structure.

**Modules**



Solaris loadable modules.

Perhaps the best current methodology for operating-system design involves using **loadable kernel modules**. Here, the kernel has a set of core components and links in additional services via modules, either at boot time or during run time. This type of design is common in modern implementations of UNIX, suchas Solaris, Linux, and Mac OS X, as well as Windows.

The idea of the design is for the kernel to provide core services while other services are implemented dynamically, as the kernel is running. Linking services dynamically is preferable to adding new features directly to the kernel, which would require recompiling the kernel every time a change was made.

Thus, for example, we might build CPU scheduling and memory management algorithms directly into the kernel and then add support for different file systems by way of loadable modules.

The Solaris operating system structure, shown in above Figure, is organized around a core kernel with seven types of loadable kernel modules:

1. Scheduling classes
2. File systems
3. Loadable system calls
4. Executable formats
5. STREAMS modules
6. Miscellaneous
7. Device and bus drivers

**Hybrid Systems**

A **hybrid** kernel is an **operating system** kernel architecture that attempts to combine aspects and benefits of microkernel and monolithic kernel architectures used in computer **operating systems**.

In practice, very few operating systems adopt a single, strictly defined structure. Instead, they combine different structures, resulting in hybrid systems that address performance, security, and usability issues. For example, both Linux and Solaris are monolithic, because having the operating system in a single address space provides very efficient performance.

43

However, they are also modular, so that new functionality can be dynamically added to the kernel. Windows is largely monolithic as well (again primarily for performance reasons), but it retains some behavior typical of microkernel systems, including providing support for separate subsystems (known as operating-system *personalities*) that run as user-mode processes.

## Mac OS X



The Mac OS X structure.

The Apple Mac OS X operating system uses a hybrid structure. As shown in Figure , it is a layered system. The top layers include the *Aqua* user interface and a set of application environments and services.

Notably, the **Cocoa** environment specifies an API for the Objective-C programming language, which is used for writing Mac OS X applications. Below these layers is the *kernel environment*, which consists primarily of the Mach microkernel and the BSD UNIX kernel. Mach provides memory management; support for remote procedure calls (RPCs) and interprocess communication (IPC) facilities, including message passing; and thread scheduling.

## iOS



Architecture of Apple's iOS.

iOS is a mobile operating system designed by Apple to run its smartphone, the *iPhone*, as well as its tablet computer, the *iPad*. iOS is structured on the Mac OS X operating system, with added functionality pertinent to mobile devices, but does not directly run Mac OS X applications. The structure of iOS appears in Figure.

**Cocoa Touch** is an API for Objective-C that provides several frameworks for developing applications that run on iOS devices. The fundamental difference between Cocoa, mentioned earlier, and Cocoa Touch is that the latter provides support for hardware features unique to mobile devices, such as touch screens.

The **media services** layer provides services for graphics, audio, and video.

**Android**



Architecture of Google's Android.

The Android operating system was designed by the Open Handset Alliance (led primarily by Google) and was developed for Android smartphones and tablet computers. Whereas iOS is designed to run on Apple mobile devices and is close-sourced, Android runs on a variety of mobile platforms and is open-sourced, partly explaining its rapid rise in popularity. The structure of Android appears in Figure 2.18.

Android is similar to iOS in that it is a layered stack of software that provides a rich set of frameworks for developing mobile applications. At the bottom of this software stack is the Linux kernel, although it has been modified by Google and is currently outside the normal distribution of Linux releases.

Software designers for Android devices develop applications in the Java language. However, rather than using the standard Java API, Google has designed a separate Android API for Java development. The Java class files are first compiled to Java bytecode and then translated into an executable file that runs on the Dalvik virtual machine. The Dalvik virtual machine was designed for Android and is optimized for mobile devices with limited memory and CPU processing capabilities.

## SYSTEM PROGRAMS

Another aspect of a modern system is its collection of system programs. In computer hierarchy, at the lowest level is hardware. Next is the operating system, then the system programs, and finally the application programs.

**System programs**, also known as **system utilities**, provide a convenient environment for program development and execution.

Some of them are simply user interfaces to system calls. Others are considerably more complex.

They can be divided into these categories:

• **File management**. These programs create, delete, copy, rename, print, dump, list, and generally manipulate files and directories.

• **Status information**. Some programs simply ask the system for the date, time, amount of available memory or disk space, number of users, or similar status information. Others are more complex, providing detailed performance, logging, and debugging information.

**File modification**. Several text editors may be available to create and modify the content of files stored on disk or other storage devices. There may also be special commands to search contents of files or perform transformations of the text.

• **Programming-language support**. Compilers, assemblers, debuggers, and interpreters for common programming languages (such as C, C++, Java, and PERL) are often provided with the operating system or available as a separate download.

• **Program loading and execution**. Once a program is assembled or compiled, it must be loaded into memory to be executed. The system may provide absolute loaders, relocatable loaders, linkage editors, and overlay loaders. Debugging systems for either higher-level languages or machine language are needed as well.

• **Communications**. These programs provide the mechanism for creating virtual connections among processes, users, and computer systems. They allow users to send messages to one another's screens, to browse Web pages, to send e-mail messages, to log in remotely, or to transfer files from one machine to another.

• **Background services**. All general-purpose systems have methods for launching certain system-program processes at boot time. Some of these processes terminate after completing their tasks, while others continue to run until the system is halted. Constantly running system-program processes are known as **services**, **subsystems**, or daemons.

One example is process schedulers that start processes according to a specified schedule, system error monitoring services, and print servers. Typical systems have dozens of daemons.

In addition, operating systems that run important activities in user context rather than in kernel context may use daemons to run these activities.

Along with system programs, most operating systems are supplied with programs that are useful in solving common problems or performing common operations. Such **application programs** include Web browsers, word processors and text formatters, spreadsheets, database systems, compilers, plotting and statistical-analysis packages, and games.

## OPEN-SOURCE OPERATING SYSTEMS

Operating systems has been made easier by the availability of a vast number of open-source format rather than as compiled binary code.

Linux is the most famous **open-source** operating system, while Microsoft Windows is a well-known example of the opposite **closed-source** approach.

Apple's Mac OS X and iOS operating systems comprise a hybrid approach. They contain an open-source kernel named Darwin yet include proprietary, closed-source components as well.

Starting with the source code allows the programmer to produce binary code that can be executed on a system. Doing the opposite—**reverse engineering** the source code from the binaries—is quite a lot of work, and useful items such as comments are never recovered. Learning operating systems by examining the source code has other benefits as well. With the source code in hand, a student can modify the operating system and then compile and run the code to try out those changes, which is an excellent learning tool.

There are many benefits to open-source operating systems, including a community of interested (and usually unpaid) programmers who contribute to the code by helping to debug it, analyze it, provide support, and suggest changes.

Arguably, open-source code is more secure than closed-source code because many more eyes are viewing the code.

## VIRTUAL MACHINE

A virtual machine is a program that acts as a virtual computer. It runs on your current operating system – the "host" operating system – and provides virtual hardware to "guest" operating systems.
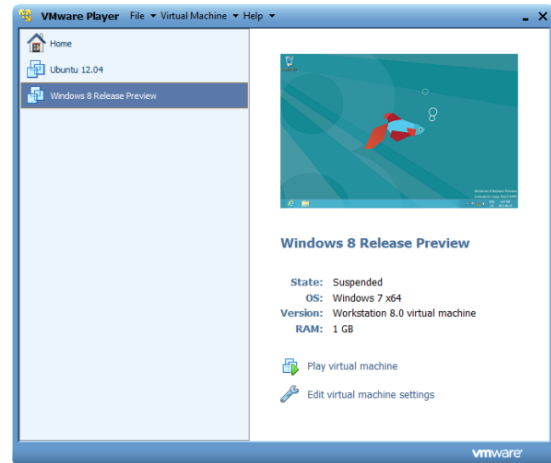
The guest operating systems run in windows on your host operating system, just like any other program on your computer.

The guest operating system runs normally, as if it were running on a physical computer – from the guest operating system's perspective, the virtual machine appears to be a real, physical computer.

Virtual machines provide their own virtual hardware, including a virtual CPU, memory, hard drive, network interface, and other devices. The virtual hardware devices provided by the virtual machine

are mapped to real hardware on your physical machine. For example, a virtual machine's virtual hard disk is stored in a file located on your hard drive.

You can have several virtual machines installed on your system; you're only limited by the amount of storage you have available for them. Once you've installed several operating systems, you can open your virtual machine program and choose which virtual machine you want to boot – the guest operating system starts up and runs in a window on your host operating system, although you can also run it in full-screen mode.



## Uses for Virtual Machines:

Virtual machines have a number of popular uses:

- **Test new versions of operating systems**: You can **run the development version of Windows 8 in a virtual machine** on your Windows 7 computer. This allows you to experiment with Windows 8 without **installing an unstable version of Windows on your computer**.
- **Experiment with other operating systems**: You can **install various distributions of Linux and other more obscure operating systems** in a virtual machine to experiment with them and learn how they work. If you're interested in **Ubuntu**, you can install it in a virtual machine and play with it at your own pace — in a window on your normal desktop.
- **Use software requiring an outdated operating system**: If you've got an important application that only runs on Windows XP, you can **install XP in a virtual machine** and run the application in the virtual machine. The virtual machine is actually running Windows XP, so compatibility shouldn't be a problem. This allows you to use an application that only works with Windows XP without actually installing Windows XP on your computer – especially important considering many new laptops and other hardware may not fully support Windows XP.

- **Run software designed for another operating system**s: Mac and Linux users can run Windows in a virtual machine to run Windows software on their computers without the compatibility headaches of **Wine** and Crossover. Unfortunately, games can be a problem – virtual machine programs introduce overhead and no virtual machine application will allow you to run the latest 3D games in a virtual machine. Some 3D effects are supported, but **3D graphics** are the least well supported thing you can do in a virtual machine.

- **Test software on multiple platforms**: If you need to test whether an application works on multiple operating systems – or just different versions of Windows – you can install each in a virtual machine instead of keeping separate computers around for each.

- **Consolidate servers**: For businesses running multiple **servers**, existing servers can be placed into virtual machines and run on a single computer. Each virtual machine is an isolated container, so this doesn't introduce the security headaches involved with running different servers on the same operating system. The virtual machines can also be moved between physical servers.

**Recommended Virtual Machine Software**

VirtualBox is a great, open-source application that runs on Windows, Mac OS X, and Linux. One of the best things about VirtualBox is that there's no commercial version – you get all the features for free, including advanced features like "snapshots," which allow you to take a snapshot of a virtual machine's state and revert to that state in the future – a great feature for testing.

VMware Player is another high-quality virtual machine program for Windows and Linux. VMware Player is the free counterpart to VMware Workstation, a commercial application, so you don't get all the advanced features you would with VirtualBox. However, both VirtualBox and VMware Player are solid programs that offer the basic features – creating and running virtual machines – for free. If one of them doesn't work quite right, try the other.

# OPERATING-SYSTEM GENERATION

It is possible to design, code, and implement an operating system specifically for one machine at one site. More commonly, however, operating systems are designed to run on any of a class of machines at a variety of sites with a variety of peripheral configurations.

The system must then be configured or generated for each specific computer site, a process sometimes known as **system generation SYSGEN**.

The operating system is normally distributed on disk, on CD-ROM or DVD-ROM, or as an "ISO" image, which is a file in the format of a CD-ROM or DVD-ROM. To generate a system, we use a special program. This SYSGEN program reads from a given file, or asks the operator of the system for information concerning the specific configuration of the hardware system, or probes the hardware directly to determine what components are there. The following kinds of information must be determined.

• **What CPU is to be used?**

What options (extended instruction sets, floating point arithmetic, and so on) are installed? For multiple CPU systems, each CPU may be described.

• **How will the boot disk be formatted?**

How many sections, or "partitions," will it be separated into, and what will go into each partition? How much memory is available? Some systems will determine this value
themselves by referencing memory location after memory location until an "illegal address" fault is generated. This procedure defines the final legal address and hence the amount of available memory.

• **What devices are available?**

The system will need to know how to address each device (the device number), the device interrupt number, the device's type and model, and any special device characteristics.

• **What operating-system options are desired, or what parameter values are to be used?**

These options or values might include how many buffers of which sizes should be used, what type of CPU-scheduling algorithm is desired, what the maximum number of processes to be supported is, and so on.

Once this information is determined, it can be used in several ways. At one extreme, a system administrator can use it to modify a copy of the source code of the operating system. The operating

system then is completely compiled. Data declarations, initializations, and constants, along with conditional compilation, produce an output-object version of the operating system that is tailored to the system described.

At a slightly less tailored level, the system description can lead to the creation of tables and the selection of modules from a precompiled library. These modules are linked together to form the generated operating system. Selection allows the library to contain the device drivers for all supported I/O devices, but only those needed are linked into the operating system. Because the system is not recompiled, system generation is faster, but the resulting system may be overly general. At the other extreme, it is possible to construct a system that is completely table driven. All the code is always part of the system, and selection occurs at execution time, rather than at compile or link time. System generation involves simply creating the appropriate tables to describe the system.

The major differences among these approaches are the size and generality of the generated system and the ease of modifying it as the hardware configuration changes. Consider the cost of modifying the system to support a newly acquired graphics terminal or another disk drive. Balanced against that cost, of course, is the frequency (or infrequency) of such changes.

## SYSTEM BOOT

[

- Booting the system is done by loading the kernel into main memory, and starting its execution.
- The CPU is given a reset event, and the instruction register is loaded with a predefined memory location, where execution starts.
  - o The initial bootstrap program is found in the BIOS read-only memory.
  - o This program can run diagnostics, initialize all components of the system, loads and starts the Operating System loader. (Called **boot strapping**)
  - o The loader program loads and starts the operating system.
  - o When the Operating system starts, it sets up needed data structures in memory, sets several registers in the CPU, and then creates and starts the first user level program. From this point, the operating system only runs in response to interrupts.

]

After an operating system is generated, it must be made available for use by the hardware. But how does the hardware know where the kernel is or how to load that kernel?

The procedure of starting a computer by loading the kernel is known as **booting** the system. On most computer systems, a small piece of code known as the **bootstrap program** or **bootstrap loader** locates the kernel, loads it into main memory, and starts its execution.

Some computer systems, such as PCs, use a two-step process in which a simple bootstrap loader fetches a more complex boot program from disk, which in turn loads the kernel. When a CPU receives a reset event—for instance, when it is powered up or rebooted—the instruction register is loaded with a predefined memory location, and execution starts there. At that location is the initial bootstrap program. This program is in the form of **read-only memory (ROM)**, because the RAM is in an unknown state at system startup. ROM is convenient because it needs no initialization and cannot easily be infected by a computer virus.

The bootstrap program can perform a variety of tasks.

Usually, one task is to run diagnostics to determine the state of the machine. If the diagnostics pass, the program can continue with the booting steps. It can also initialize all aspects of the system, from CPU registers to device controllers and the contents of main memory. Sooner or later, it starts the operating system.

Some systems—such as cellular phones, tablets, and game consoles—store the entire operating system in ROM. Storing the operating system in ROM is suitable for small operating systems, simple supporting hardware, and rugged operation. A problem with this approach is that changing the bootstrap code requires changing the ROM hardware chips. Some systems resolve this problem by using **erasable programmable read-only memory (EPROM)**, which is read only except when explicitly given a command to become writable. All forms of ROM are also known as **firmware**, since their characteristics fall somewhere between those of hardware and those of software. A problem with firmware in general is that executing code there is slower than executing code in RAM.

Some systems store the operating system in firmware and copy it to RAM for fast execution. A final issue with firmware is that it is relatively expensive, so usually only small amounts are available.

For large operating systems (including most general-purpose operating systems like Windows, Mac OS X, and UNIX) or for systems that change frequently, the bootstrap loader is stored in firmware, and the operating system is on disk. In this case, the bootstrap runs diagnostics and has a bit of code that can read a single block at a fixed location (**say block zero**) from disk into memory and execute the code from that **boot block**.

# UNIT – II         PROCESS MANAGEMENT

- A process can be thought of as a program in execution. A process will need certain resources - such as CPU time, memory, files, and I/O devices - to accomplish its task. These resources are allocated to the process either when it is created or while it is executing.

- A process is the unit of work in most systems. Systems consist of a collection of processes: operating-system processes execute system code, and user processes execute user code. All these processes may execute concurrently.

- Although traditionally a process contained only a single thread of control as it ran, most modern operating systems now support processes that have multiple threads.

- The operating system is responsible for several important aspects of process and thread management: the creation and deletion of both user and system processes; the scheduling of processes; and the provision of mechanisms for synchronization, communication, and deadlock handling for processes.

## Process Concept:

A question that arises in discussing operating systems involves what to call all the CPU activities. A batch system executes **jobs**, whereas a time-shared system has **user programs**, or **tasks**.

Even on a single-user system, a user may be able to run several programs at one time: a word processor, a Web browser, and an e-mail package.

And even if a user can execute only one program at a time, such as on an embedded device that does not support multitasking, the operating system may need to support its own internal programmed activities, such as memory management.

In many respects, all these activities are similar, so we call all of them **processes**.

The terms *job* and *process* are used almost interchangeably in this text.

## The Process:



- A process is more than the program code, which is sometimes known as the **text section**.

- It also includes the current activity, as represented by the value of the **program counter** and the contents of the processor's registers.

- A process generally also includes the process **stack**, which contains temporary data (such as function parameters, return addresses, and local variables), and

- A **data section**, which contains global variables.

- A process may also include a **heap**, which is memory that is dynamically allocated during process run time. The structure of a process in memory is shown in Figure. We emphasize that a program by itself is not a process.
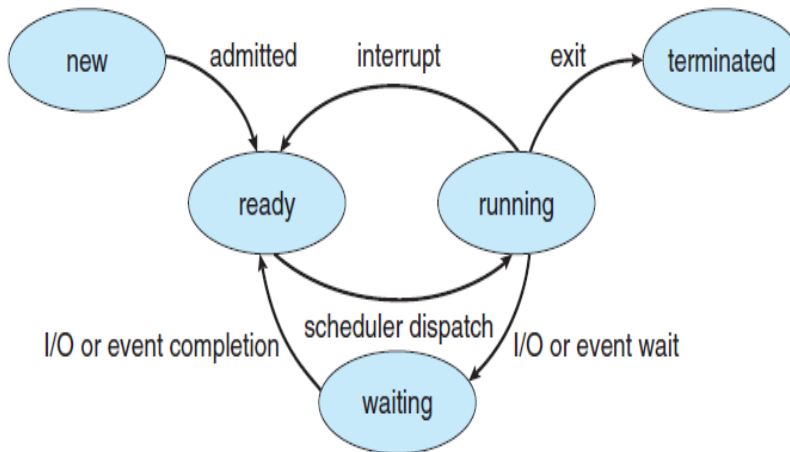
- A program is a *passive* entity, such as a file containing a list of instructions stored on disk (often called an **executable file**).
- In contrast, a process is an *active* entity, with a program counter specifying the next instruction to execute and a set of associated resources. A program becomes a process when an executable file is loaded into memory.

Two common techniques for loading executable files are

- Double-clicking an icon representing the executable file and
- Entering the name of the executable file on the command line (as in prog.exe or a.out).

## 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:  The state diagram corresponding to these states is presented in above Figure.



**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)

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

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

It is important to realize that only one process can be *running* on any processor at any instant. Many processes may be *ready* and *waiting.*

## Process Control Block-[PCB]:



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.

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.
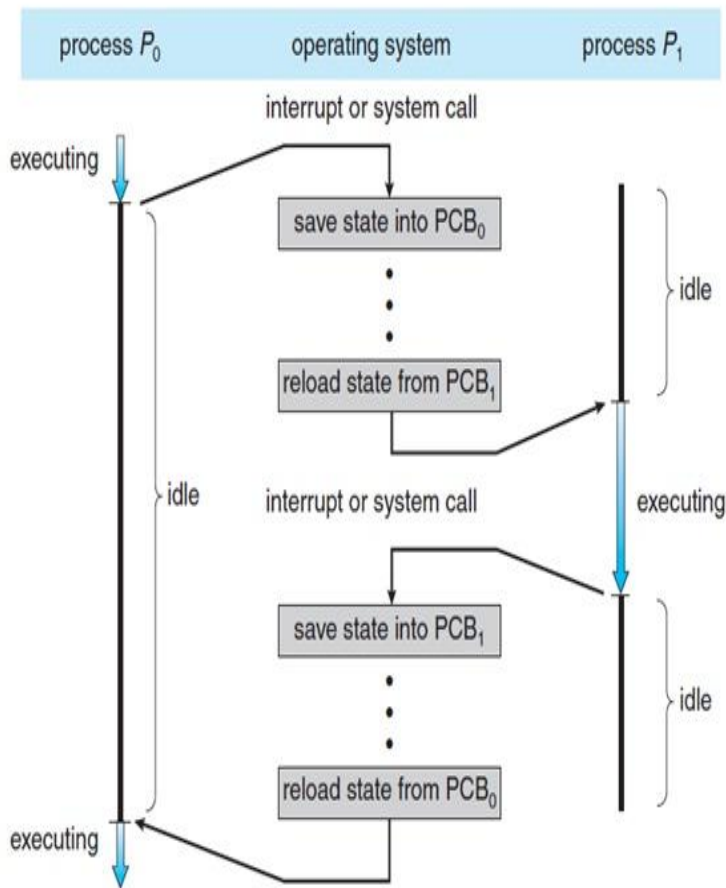
- **CPU registers**. The registers vary in number and type, depending on the computer architecture. They include accumulators, index registers, stack pointers, and general-purpose registers, plus any condition-code information.

- **CPU-scheduling information**. This information includes a process priority, pointers to scheduling queues, and any other scheduling parameters.

- **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.

- **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.

Diagram showing CPU switch from process to process.

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

## 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.
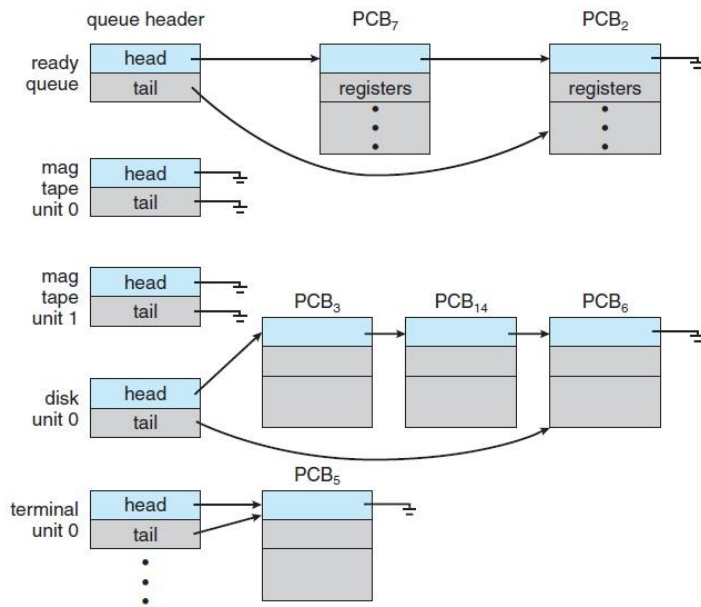
To meet these objectives, the **process scheduler** selects an available process (possibly from a set of several available processes) for program execution on the CPU.

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.

**Scheduling Queues:**

As processes enter the system, they are put into a **job queue**, which consists of all processes in the system.

55

The processes that are residing in main memory and are ready and waiting to execute are kept on a list called the **ready queue**.
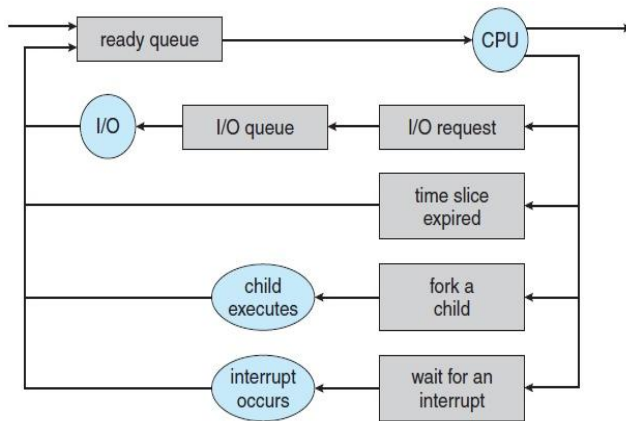


The ready queue and various I/O device queues.

This queue is generally stored as a linked list. A ready-queue header contains pointers to the first and final PCBs in the list. Each PCB includes a pointer field that points to the next PCB in the ready queue.

The system also includes other queues. When a process is allocated the CPU, it executes for a while and eventually quits, is interrupted, or waits for the occurrence of a particular event, such as the completion of an I/O request. Suppose the process makes an I/O request to a shared device, such as a disk.

Since there are many processes in the system, the disk may be busy with the I/O request of some other process. The process therefore may have to wait for the disk. The list of processes waiting for a particular I/O device is called a **device queue**. Each device has its own device queue (following Figure).



Queueing-diagram representation of process scheduling.

A common representation of process scheduling is a **queuing diagram**, such as that in Figure. Each rectangular box represents a queue. Two types of queues are present: the ready queue and a set of device queues. The circles represent the resources that serve the queues, and the arrows indicate the flow of processes in the system.

A new process is initially put in the ready queue. It waits there until it is selected for execution, or **dispatched**. Once the process is allocated the CPU and is executing, one of several events could occur:

• The process could issue an I/O request and then be placed in an I/O queue.

• The process could create a new child process and wait for the child's termination.

• The process could be removed forcibly from the CPU, as a result of an interrupt, and be put back in the ready queue.

In the first two cases, the process eventually switches from the waiting state to the ready state and is then put back in the ready queue. A process continues this cycle until it terminates, at which time it is removed from all queues and has its PCB and resources deallocated.

## 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.

The long-term scheduler executes much less frequently; minutes may separate the creation of one new process and the next. The long-term scheduler controls the **degree of multiprogramming** (the number of processes in memory).

Some operating systems, such as time-sharing systems, may introduce an additional, intermediate level of scheduling. This **medium-term scheduler** is diagrammed in Figure. The key idea behind a medium-term scheduler is that sometimes it can be advantageous to remove a process from memory (and from active contention for the CPU) and thus reduce the degree of multiprogramming. Later, the process can be reintroduced into memory, and its execution can be continued where it left off. This scheme is called **swapping**. The process is swapped out, and is later swapped in, by the medium-term scheduler. Swapping may be necessary to improve the process mix or because a change in memory requirements has overcommitted available memory, requiring memory to be freed up.

## Context Switch:

Switching the CPU to another process requires performing a state save of the current process and a state restore of a different process. This task is known as a **context switch**.

When a context switch occurs, the kernel saves the context of the old process in its PCB and loads the saved context of the new process scheduled to run.
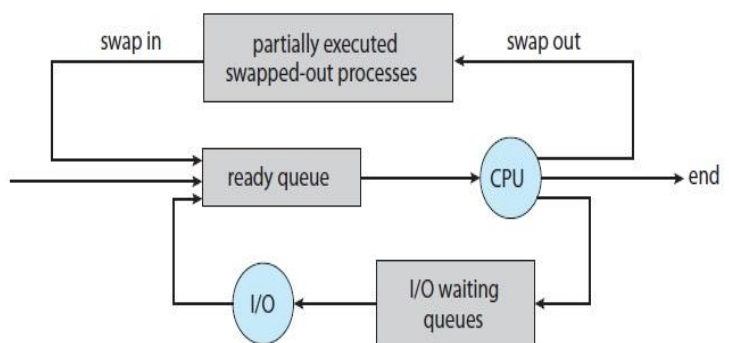


Figure 3.7  Addition of medium-term scheduling to the queueing diagram.

Context-switch time is pure overhead, because the system does no useful work while switching. Switching speed varies from machine to machine, depending on the memory speed, the number of registers that must be copied, and the existence of special instructions (such as a single instruction to load or store all registers). A typical speed is a few milliseconds.

## OPERATIONS ON PROCESSES

The processes in most systems can execute concurrently, and they may be created and deleted dynamically. Thus, these systems must provide a mechanism for process creation and termination.

In this section, we explore the mechanisms involved in creating processes and illustrate process creation on UNIX and Windows systems.

### Process Creation

During the course of execution, a process may create several new processes.

As mentioned earlier,

the creating process is called a **parent process**, and the new processes are called the **children of that process.**

Each of these new processes may in turn create other processes, forming a **tree** of processes.



A tree of processes on a typical Linux system.

Most operating systems (including UNIX, Linux, and Windows) identify processes according to a unique **process identifier** (or **pid**), which is typically an integer number.

The pid provides a unique value for each process in the system, and it can be used as an index to access various attributes of a process With in the kernel.

In general, when a process creates a child process, that child process will need certain resources (CPU time, memory, files, I/O devices) to accomplish its task.

A child process may be able to obtain its resources directly from the operating system, or it may be constrained to a subset of the resources of the parent process. The parent may have to partition its resources among its children, or it may be able to share some resources (such as memory or files) among several of its children.

Restricting a child process to a subset of the parent's resources prevents any process from overloading the system by creating too many child processes.

When a process creates a new process, two possibilities for execution exist:

**1.** The parent continues to execute concurrently with its children.

**2.** The parent waits until some or all of its children have terminated.

There are also two address-space possibilities for the new process:

**1.** The child process is a duplicate of the parent process (same program and data as the parent).

**2.** The child process has a new program loaded into it.

```
#include <sys/types.h>
#include <stdio.h>
#include <unistd.h>

int main()
{
pid_t pid;

    /* fork a child process */
    pid = fork();

    if (pid < 0) { /* error occurred */
       fprintf(stderr, "Fork Failed");
       return 1;
    }
    else if (pid == 0) { /* child process */
       execlp("/bin/ls","ls",NULL);
    }
    else { /* parent process */
       /* parent will wait for the child to complete */
       wait(NULL);
       printf("Child Complete");
    }

    return 0;
}
```

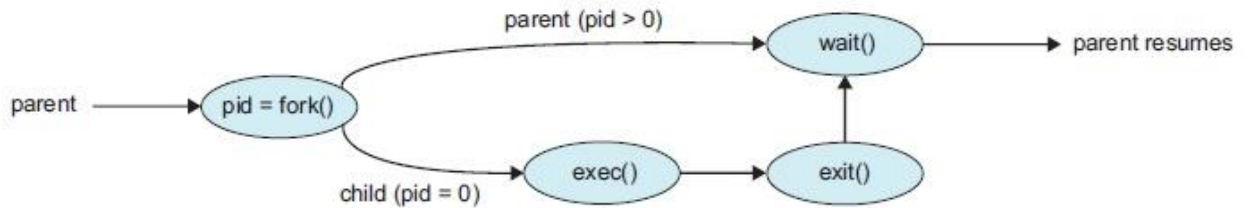**Figure -** Creating a separate process using the UNIX fork() system call.

The C program shown in above Figure illustrates the UNIX system calls previously described. We now have two different processes running copies of the same program. The only difference is that the value of pid (the process identifier) for the child process is zero, while that for the parent is an integer value greater than zero (in fact, it is the actual pid of the child process).

The child process inherits privileges and scheduling attributes from the parent, as well certain resources, such as open files. The child process then overlays its address space with the UNIX command /bin/ls (used to get a directory listing) using the execlp() system call (execlp() is a version of the exec() system call).

The parent waits for the child process to complete with the wait() system call. When the child process completes (by either implicitly or explicitly invoking exit()), the parent process resumes from the call to wait(), where it completes using the exit() system call. This is also illustrated in the following Figure.



Process creation using the fork() system call.

Of course, there is nothing to prevent the child from *not* invoking exec()and continuing to execute as a copy of the parent process. In this scenario, the parent and child are concurrent processes running the same code instructions. Because the child is a copy of the parent, each process has its own copy of any data.

**Process Termination**

A process terminates when it finishes executing its final statement and asks the operating system to delete it by using the exit() system call. At that point, the process may return a status value (typically an integer) to its parent process (via the wait() system call).

All the resources of the process—including physical and virtual memory, open files, and I/O buffers—are deallocated by the operating system.

Termination can occur in other circumstances as well. A process can cause the termination of another process via an appropriate system call (for example, TerminateProcess() in Windows).

Usually, such a system call can be invoked only by the parent of the process that is to be terminated. Otherwise, users could arbitrarily kill each other's jobs. Note that a parent needs to know the identities of its children if it is to terminate them.

Thus, when one process creates a new process, the identity of the newly created process is passed to the parent.

A parent may terminate the execution of one of its children for a variety of reasons, such as these:

- The child has exceeded its usage of some of the resources that it has been allocated. (To determine whether this has occurred, the parent must have a mechanism to inspect the state of its children.)
- The task assigned to the child is no longer required.
- The parent is exiting, and the operating system does not allow a child to continue if its parent terminates.

Some systems do not allow a child to exist if its parent has terminated. In such systems, if a process terminates (either normally or abnormally), then all its children must also be terminated. This phenomenon, referred to as **cascading termination**, is normally initiated by the operating system.

To illustrate process execution and termination, consider that, in Linux and UNIX systems, we can terminate a process by using the exit() system call, providing an exit status as a parameter:

**/* exit with status 1 */**

**exit(1);**

In fact, under normal termination, exit() may be called either directly (as shown above) or indirectly (by a return statement in main()).

## 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. 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.

• **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.



Communications models. (a) Message passing. (b) Shared memory.

Cooperating processes require an **interprocess communication (IPC)** mechanism that will allow them to exchange data and information. There are two fundamental models of interprocess communication:

- **Shared memory**
- **Message passing**.

In the shared-memory model, a region of memory that is shared by cooperating processes is established. Processes can then exchange information by reading and writing data to the shared region.

In the message-passing model, communication takes place by means of messages exchanged between the cooperating processes.

The two communications models are contrasted in above Figure.

Both of the models just mentioned are common in operating systems, and many systems implement both.

Message passing is useful for exchanging smaller amounts of data, because no conflicts need be avoided. Message passing is also easier to implement in a distributed system than shared memory.

Shared memory can be faster than message passing, since message-passing systems are typically implemented using system calls and thus require the more time-consuming task of kernel intervention. In shared-memory systems, system calls are required only to establish shared memory regions. Once shared memory is established, all accesses are treated as routine memory accesses, and no assistance from the kernel is required.

## Shared-Memory Systems

Interprocess communication using shared memory requires communicating processes to establish a region of shared memory.

They can then exchange information by reading and writing data in the shared areas.

The form of the data and the location are determined by these processes and are not under the operating system's control.

The processes are also responsible for ensuring that they are not writing to the same location simultaneously.

To illustrate the concept of cooperating processes, let's consider the producer – consumer problem,

which is a common paradigm for cooperating processes. A **producer** process produces information that is consumed by a **consumer** process.

For example, a compiler may produce assembly code that is consumed by an assembler. The assembler, in turn, may produce object modules that are consumed by the loader.

The producer – consumer problem also provides a useful metaphor for the client – server paradigm. We generally think of a server as a producer and a client as a consumer. For example, a web server produces (that is, provides) HTML files and images, which are consumed (that is, read) by the client web browser requesting the resource.

One solution to the producer – consumer problem uses shared memory.

To allow producer and consumer processes to run concurrently, we must have available a buffer of items that can be filled by the producer and emptied by the consumer. This buffer will reside in a region of memory that is shared by the producer and consumer processes. A producer can produce one item while the consumer is consuming another item. The producer and consumer must be synchronized, so that the consumer does not try to consume an item that has not yet been produced.

```
while (true)
{
        /* produce an item in next produced */
        while (((in + 1) % BUFFER SIZE) == out)
            ;    /* do nothing */
        buffer[in] = next produced;
        in = (in + 1) % BUFFER SIZE;
}
```
        The producer process using shared memory.

Two types of buffers can be used.

        The **unbounded buffer** places no practical limit on the size of the buffer. The consumer may have to wait for new items, but the producer can always produce new items.

        The **bounded buffer** assumes a fixed buffer size. In this case, the consumer must wait if the buffer is empty, and the producer must wait if the buffer is full.

## Message-Passing Systems

        The shared-memory environment requires that these processes share a region of memory and that the code for accessing and manipulating the shared memory be written explicitly by the application programmer.

        Another way to achieve the same effect is for the operating system to provide the means for cooperating processes to communicate with each other via a message-passing facility.

        Message passing provides a mechanism to allow processes to communicate and to synchronize their actions without sharing the same address space. It is particularly useful in a distributed environment, where the communicating processes may reside on different computers connected by a network.

        For example, an Internet chat program could be designed so that chat participants communicate with one another by exchanging messages.

        A message-passing facility provides at least two operations:

$$\text{send(message)} \qquad \text{receive(message)}$$

Messages sent by a process can be either fixed or variable in size.

If only fixed-sized messages can be sent, the system-level implementation is straight-forward. This restriction, however, makes the task of programming more difficult.

Conversely, variable-sized messages require a more complex system- level implementation, but the programming task becomes simpler.

If processes *P* and *Q* want to communicate, they must send messages to and receive messages from each other: a ***communication link*** must exist between them.

This link can be implemented in a variety of ways. We are concerned here not with the link's physical implementation (such as shared memory, hardware bus, or network) but rather with its logical implementation.

Here are several methods for logically implementing a link and the send()/receive() operations:

- Direct or indirect communication
- Synchronous or asynchronous communication
- Automatic or explicit buffering


## Naming

Processes that want to communicate must have a way to refer to each other.

They can use either direct or indirect communication.

Under **direct communication**, each process that wants to communicate must explicitly name the recipient or sender of the communication. In this scheme, the send() and receive() primitives are defined as:

send(P, message)— Send a message to process P.

receive(Q, message)— Receive a message from process Q.

A communication link in this scheme has the following properties:

- A link is established automatically between every pair of processes that want to communicate. The processes need to know only each other 's identity to communicate.
- A link is associated with exactly two processes.
- Between each pair of processes, there exists exactly one link.


## Synchronization

Communication between processes takes place through calls to send() and receive() primitives. There are different design options for implementing each primitive. Message passing may be either **blocking** or **nonblocking**— also known as **synchronous** and **asynchronous**. (Throughout this text, you will encounter the concepts of synchronous and asynchronous behavior in relation to various operating-system algorithms.)

**Blocking send**. The sending process is blocked until the message is received by the receiving process or by the mailbox.

**Nonblocking send**. The sending process sends the message and resumes operation.

**Blocking receive**. The receiver blocks until a message is available.

**Nonblocking receive**. The receiver retrieves either a valid message or a null.

## Buffering

Whether communication is direct or indirect, messages exchanged by communicating processes reside in a temporary queue. Basically, such queues can be implemented in three ways:

**Zero capacity**. The queue has a maximum length of zero; thus, the link cannot have any messages waiting in it. In this case, the sender must block until the recipient receives the message.

**Bounded capacity**. The queue has finite length $n;$ thus, at most $n$ messages can reside in it. If the queue is not full when a new message is sent, the message is placed in the queue (either the message is copied or a pointer to the message is kept), and the sender can continue execution without waiting. The link's capacity is finite, however. If the link is full, the sender must block until space is available in the queue.

**Unbounded capacity**. The queue's length is potentially infinite; thus, any number of messages can wait in it. The sender never blocks.

## COMMUNICATION IN CLIENT – SERVER SYSTEMS

how processes can communicate using shared memory and message passing?.

These techniques can be used for communication in client – server systems. Communication in client – server systems may use

**(1) sockets,**

**(2) remote procedure calls (RPCs),**

**(3) pipes.**

Pipes provide a relatively simple ways for processes to communicate with one another. Ordinary pipes allow communication between parent and child processes, while named pipes permit unrelated processes to communicate.

## Sockets:

A **socket** is defined as an endpoint for communication.

A socket is defined as an endpoint for communication. A connection between a pair of applications consists of a pair of sockets, one at each end of the communication channel.
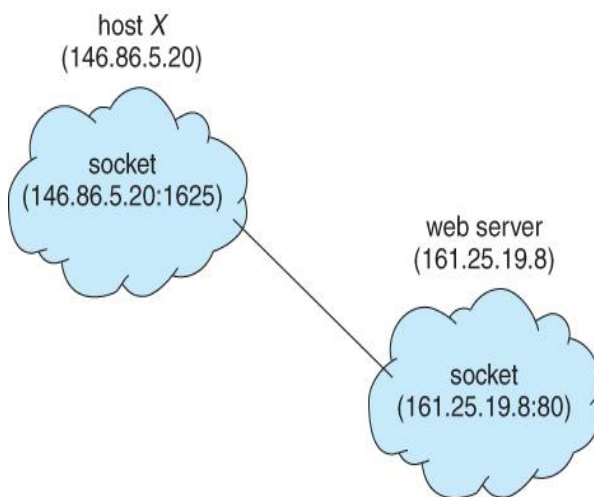
A pair of processes communicating over a network employs a pair of sockets — one for each process. A socket is identified by an IP address concatenated with a port number.

In general, sockets use a client – server architecture. The server waits for incoming client requests by listening to a specified port.

Once a request is received, the server accepts a connection from the client socket to complete the connection.

Servers implementing specific services (such as telnet, FTP, and HTTP) listen to well-known ports (a telnet server listens to port 23; an FTP server listens to port 21; and a web, or HTTP, server listens to port 80).

All ports below 1024 are considered *well known;* we can use them to implement standard services.

host X
(146.86.5.20)

socket
(146.86.5.20:1625)

web server
(161.25.19.8)

socket
(161.25.19.8:80)

When a client process initiates a request for a connection, it is assigned a port by its host computer. This port has some arbitrary number greater than 1024. For example, if a client on host X with IP address 146.86.5.20 wishes to establish a connection with a web server (which is listening on port 80) at address 161.25.19.8, host X may be assigned port 1625. The connection will consist of a pair of sockets: (146.86.5.20:1625) on host X and (161.25.19.8:80) on the web server. This situation is illustrated in Figure 3.20.

The packets traveling between the hosts are delivered to the appropriate process based on the destination port number.

- Communication channels via sockets may be of one of two major forms:
    - **Connection-oriented (TCP, Transmission Control Protocol)** connections emulate a telephone connection. All packets sent down the connection are guaranteed to arrive in good condition at the other end, and to be delivered to the receiving process in the order in which they were sent. The TCP layer of the network protocol takes steps to verify all packets sent, re-send packets if necessary, and arrange the received packets in the proper order before delivering them to the receiving process. There is a certain amount of overhead involved in this procedure, and if one packet is missing or delayed, then any packets which follow will have to wait until the errant packet is delivered before they can continue their journey.
    - **Connectionless (UDP, User Datagram Protocol)** emulate individual telegrams. There is no guarantee that any particular packet will get through undamaged (or at all), and no

guarantee that the packets will get delivered in any particular order. There may even be duplicate packets delivered, depending on how the intermediary connections are configured. UDP transmissions are much faster than TCP, but applications must implement their own error checking and recovery procedures.

- Sockets are considered a low-level communications channel, and processes may often choose to use something at a higher level, such as those covered in the next two sections.

## Remote Procedure Calls

In contrast IPC messages, the messages exchanged in RPC communication are well structured and are thus no longer just packets of data.

RPCs are another form of distributed communication. An RPC occurs when a process (or thread) calls a procedure on a remote application.

Each message is addressed to an RPC daemon listening to a port on the remote system, and each contains an identifier specifying the function to execute and the parameters to pass to that function.

The function is then executed as requested, and any output is sent back to the requester in a separate message.

A **port** is simply a number included at the start of a message packet. Whereas a system normally has one network address, it can have many ports within that address to differentiate the many network services it supports.

If a remote process needs a service, it addresses a message to the proper port. For instance, if a system wished to allow other systems to be able to list its current users, it would have a daemon supporting such an RPC attached to a port — say, port 3027. Any remote system could obtain the needed information (that is, the list of current users) by sending an RPC message to port 3027 on the server. The data would be received in a reply message.

The semantics of RPCs allows a client to invoke a procedure on a remote host as it would invoke a procedure locally.

The RPC system hides the details that allow communication to take place by providing a **stub** on the client side.

Typically, a separate stub exists for each separate remote procedure. When the client invokes a remote procedure, the RPC system calls the appropriate stub, passing it the parameters provided to the remote procedure.

This stub locates the port on the server and **marshals** the parameters. Parameter marshalling involves packaging the parameters into a form that can be transmitted over a network. The stub then transmits a message to the server using message passing.

A similar stub on the server side receives this message and invokes the procedure on the server. If necessary, return values are passed back to the client using the same technique.

On Windows systems, stub code is compiled from a specification written in the **Microsoft Interface Definition Language (MIDL)**, which is used for defining the interfaces between client and server programs.

Two approaches are common. First, the binding information may be predetermined, in the form of fixed port addresses. At compile time, an RPC call has a fixed port number associated with it. Once a program is compiled, the server cannot change the port number of the requested service. Second, binding can be done dynamically by a rendezvous mechanism. Typically, an operating system provides a rendezvous (also called a **matchmaker**) daemon on a fixed RPC port. A client then sends a message containing the name of the RPC to the rendezvous daemon requesting the port address of the RPC it needs to execute. The port number is returned, and the RPC calls can be sent to that port until the process terminates (or the server crashes). This method requires the extra overhead of the initial request but is more flexible than the first approach. Figure 3.23 shows a sample interaction.



**Figure 3.23** Execution of a remote procedure call (RPC).

**Pipes**

A **pipe** acts as a conduit allowing two processes to communicate. Pipes were one of the first IPC mechanisms in early UNIX systems. They typically provide one of the simpler ways for processes to communicate with one another, although they also have some limitations. In implementing a pipe, four issues must be considered:
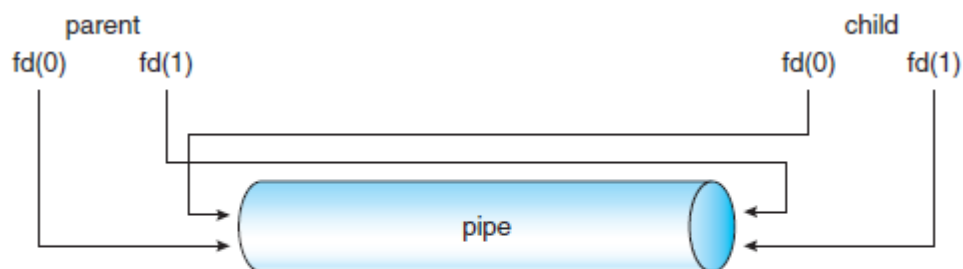


**Figure 3.24** File descriptors for an ordinary pipe.

- Does the pipe allow bidirectional communication, or is communication unidirectional?
- If two-way communication is allowed, is it half duplex (data can travel only one way at a time) or full duplex (data can travel in both directions at the same time)?
- Must a relationship (such as *parent–child*) exist between the communicating processes?
- Can the pipes communicate over a network, or must the communicating processes reside on the same machine?
- In the following sections, we explore two common types of pipes used on both UNIX and Windows systems: ordinary pipes and named pipes.

## Ordinary Pipes

Ordinary pipes allow two processes to communicate in standard producer – consumer fashion:

the producer writes to one end of the pipe (the **write-end**) and the consumer reads from the other end (the **read-end**). As a result, ordinary pipes are unidirectional, allowing only one-way communication. If two-way communication is required, two pipes must be used, with each pipe sending data in a different direction. We next illustrate constructing ordinary pipes on both UNIX and Windows systems. In both program examples, one process writes the message Greetings to the pipe, while the other process reads this message from the pipe.

On UNIX systems, ordinary pipes are constructed using the function

**pipe(int fd[])**

This function creates a pipe that is accessed through the int fd[] file descriptors: fd[0] is the read-end of the pipe, and fd[1] is the write-end. UNIX treats a pipe as a special type of file. Thus, pipes can be accessed using ordinary read() and write() system calls.

## Named Pipes

Ordinary pipes provide a simple mechanism for allowing a pair of processes to communicate. However, ordinary pipes exist only while the processes are communicating with one another.

On both UNIX and Windows systems, once the processes have finished communicating and have terminated, the ordinary pipe ceases to exist.

Named pipes provide a much more powerful communication tool. Communication can be bidirectional, and no parent – child relationship is required. Once a named pipe is established, several processes can use it for communication.

In fact, in a typical scenario, a named pipe has several writers. Additionally, named pipes continue to exist after communicating processes have finished. Both UNIX and Windows systems support named pipes, although the details of implementation differ greatly. Next, we explore named pipes in each of these systems.

# MULTI-THREADED  PROGRAMMING

## What is Thread?

A thread is a flow of execution through the process code, with its own program counter that keeps track of which instruction to execute next, system registers which hold its current working variables, and a stack which contains the execution history.

A thread shares with its peer threads few information like code segment, data segment and open files. When one thread alters a code segment memory item, all other threads see that.

A thread is also called a **lightweight process**. Threads provide a way to improve application performance through parallelism. Threads represent a software approach to improving performance of operating system by reducing the overhead thread is equivalent to a classical process.

Each thread belongs to exactly one process and no thread can exist outside a process. Each thread represents a separate flow of control. Threads have been successfully used in implementing network servers and web server. They also provide a suitable foundation for parallel execution of applications on shared memory multiprocessors. The following figure shows the working of a single-threaded and a multithreaded process.

## Difference between Process and Thread

| S.N. | Process | Thread |
|------|---------|--------|
| 1 | Process is heavy weight or resource intensive. | Thread is light weight, taking lesser resources than a process. |
| 2 | Process switching needs interaction with operating system. | Thread switching does not need to interact with operating system. |
| 3 | In multiple processing environments, each process executes the same code but has its own memory and file resources. | All threads can share same set of open files, child processes. |

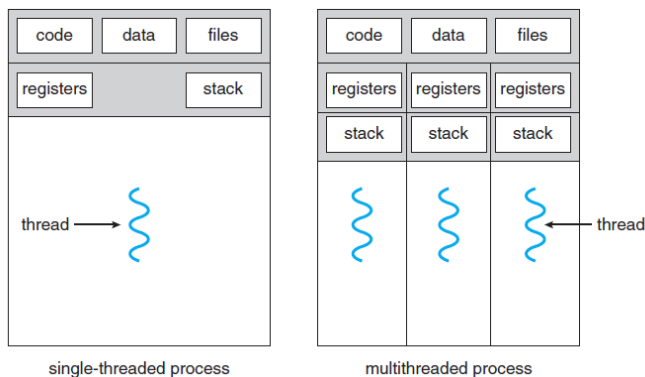| 4 | If one process is blocked, then no other process can execute until the first process is unblocked. | While one thread is blocked and waiting, a second thread in the same task can run. |
|---|---|---|
| 5 | Multiple processes without using threads use more resources. | Multiple threaded processes use fewer resources. |
| 6 | In multiple processes each process operates independently of the others. | One thread can read, write or change another thread's data. |

## Advantages of Thread

- Threads minimize the context switching time.
- Use of threads provides concurrency within a process.
- Efficient communication.
- It is more economical to create and context switch threads.
- Threads allow utilization of multiprocessor architectures to a greater scale and efficiency.

## Types of Thread

Threads are implemented in following two ways −

- **User Level Threads** − User managed threads.
- **Kernel Level Threads** − Operating System managed threads acting on kernel, an operating system core.

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. The following Figure illustrates the difference between a traditional **single-threaded** process and a **multithreaded** process.

single-threaded process        multithreaded process

## Motivation

Most software applications that run on modern computers are multithreaded. An application typically is implemented as a separate process with several threads of control.

A web browser might have one thread display images or text while another thread retrieves data from the network.

For example. A word processor may have a thread for displaying graphics, another thread for responding to keystrokes from the user, and a third thread for performing spelling and grammar checking in the background. Applications can also be designed to leverage processing capabilities on multicore systems. Such applications can perform several CPU-intensive tasks in parallel across the multiple computing cores.



**Figure 4.2**   Multithreaded server architecture.

One solution is to have the server run as a single process that accepts requests. When the server receives a request, it creates a separate process to service that request. In fact, this process-creation method was in common use before threads became popular. Process creation is time consuming and resource intensive, however. If the new process will perform the same tasks as the existing process, why incur all that overhead? It is generally more efficient to use one process that contains multiple threads. If the web-server process is multithreaded, the server will create a separate thread that listens for client requests. When a request is made, rather than creating another process, the server creates a new thread to service the request and resume listening for additional requests. This is illustrated in Figure 4.2.

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.

### Benefits

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

- **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.

- **Resource sharing**. Processes can only share resources through techniques such as shared memory and message passing.

- **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.

- **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. We explore this issue further in the following section.

### Multithreading Models

Some operating system provide a combined user level thread and Kernel level thread facility.

Solaris is a good example of this combined approach. In a combined system, multiple threads within the same application can run in parallel on multiple processors and a blocking system call need not block the entire process. Multithreading models are three types

- Many to many relationship.
- Many to one relationship.
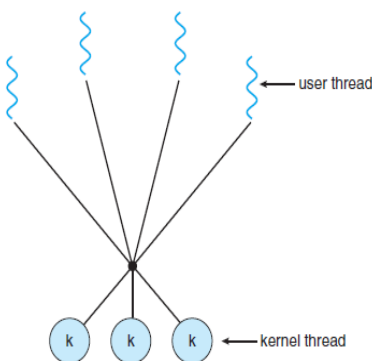- One to one relationship.

### Many to Many Model



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

The many-to-many model multiplexes any number of user threads onto an equal or smaller number of kernel threads.

The following diagram shows the many-to-many threading model where 6 user level threads are multiplexing with 6 kernel level threads. In this model, developers can create as many user threads as necessary and the corresponding Kernel threads can run in parallel on a multiprocessor machine. This model provides the best accuracy on concurrency and when a thread performs a

blocking                                system                        call,                                    the
kernel can schedule another thread for execution.
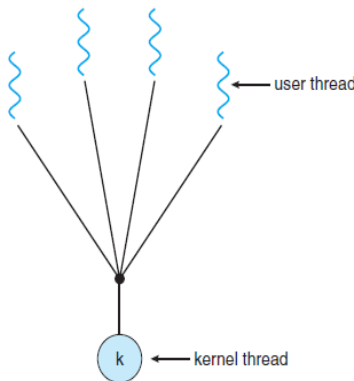
## Many to One Model



Figure 4.5  Many-to-one model.

Many-to-one model maps many user level threads to one Kernel-level thread. Thread management is done in user space by the thread library. When thread makes a blocking system call, the entire process will be blocked. Only one thread can access the Kernel at a time, so multiple threads are unable to run in parallel on multiprocessors.

If the user-level thread libraries are implemented in the operating system in such a way that the system does not support them, then the Kernel threads use the many-to-one relationship modes.

## One to One Model



Figure 4.6  One-to-one model.

There is one-to-one relationship of user-level thread to the kernel-level thread. This model provides more concurrency than the many-to-one model. It also allows another thread to run when a thread makes a blocking system call. It supports multiple threads to execute in parallel on microprocessors.

Disadvantage of this model is that creating user thread requires the corresponding Kernel thread. OS/2, windows NT and windows 2000 use one to one relationship model.

## Difference between User-Level & Kernel-Level Thread

| S.N. | User-Level Threads | Kernel-Level Thread |
|------|--------------------|---------------------|
| 1 | User-level threads are faster to create and manage. | Kernel-level threads are slower to create and manage. |
| 2 | Implementation is by a thread library at the user level. | Operating system supports creation of Kernel threads. |

| 3 | User-level thread is generic and can run on any operating system. | Kernel-level thread is specific to the operating system. |
|---|---|---|
| 4 | Multi-threaded applications cannot take advantage of multiprocessing. | Kernel routines themselves can be multithreaded. |

## Threading Issues

### The fork( ) and exec( ) System Calls

Recall that when fork() is called, a separate, duplicate process is created

• How should fork() behave in a multithreaded program? - Should all threads be duplicated?

    - Should only the thread that made the call to fork() be duplicated?

• In some systems, diff erent versions of fork() exist depending on the desired behavior

    - Some UNIX systems have fork1() and forkall() • fork1() only duplicates the calling thread

• forkall() duplicates all of the threads in a process

    - In a POSIX-compliant system, fork() behaves the same as fork1()

• The exec() system call continues to behave as expected - Replaces the entire process that called it, including all threads

• If planning to call exec() after fork(), then there is no need to duplicate all of the threads in the calling process - All threads in the child process will be terminated when exec() is called

    - Use fork1(), rather than forkall() if using in conjunction with exec()

### Signal Handling

Signals are used in UNIX systems to notify a process that a particular event has occurred

    - CTRL-C is an example of an asynchronous signal that might be sent to a process

• An asynchronous signal is one that is generated from outside the process that receives it

    - Divide by 0 is an example of a synchronous signal that might be sent to a process

• A synchronous signal is delivered to the same process that caused the signal to occur

• All signals follow the same basic pattern: - A signal is generated by particular event

    - The signal is delivered to a process

    - The signal is handled by a signal handler (all signals are handled exactly once)

Signal handling is straightforward in a single-threaded process - The one (and only) thread in the process receives and handles the signal

In a multithreaded program, where should signals be delivered? - Options:

    (1) Deliver the signal to the thread to which the signal applies

(2) Deliver the signal to every thread in the process

(3) Deliver the signal only to certain threads in the process

(4) Assign a specific thread to receive all signals for the process

• Option 1 - Deliver the signal to the thread to which the signal applies - Most likely option when handling synchronous signals (e.g. only the thread that attempts to divide by zero needs to know of the error)

• Option 2 - Deliver the signal to every thread in the process - Likely to be used in the event that the process is being terminated (e.g. a CTRLC is sent to terminate the process, all threads need to receive this signal and terminate)

## Thread Cancellation

• Thread cancellation is the act of terminating a thread before it has completed - Example - clicking the stop button on your web browser will stop the thread that is rendering the web page

• The thread to be cancelled is called the target thread • Threads can be cancelled in a couple of ways - Asynchronous cancellation terminates the target thread immediately

• Thread may be in the middle of writing data ... not so good

- Deferred cancellation allows the target thread to periodically    check if it should be cancelled

• Allows thread to terminate itself in an orderly fashion

- Threads that are no longer needed may be cancelled by another thread in one of two ways:
    1. **Asynchronous Cancellation** cancels the thread immediately.
    2. **Deferred Cancellation** sets a flag indicating the thread should cancel itself when it is convenient. It is then up to the cancelled thread to check this flag periodically and exit nicely when it sees the flag set.
- ( Shared ) resource allocation and inter-thread data transfers can be problematic with asynchronous cancellation.

## Thread-Local Storage

- Most data is shared among threads, and this is one of the major benefits of using threads in the first place.
- However sometimes threads need thread-specific data also.
- Most major thread libraries ( pThreads, Win32, Java ) provide support for thread-specific data, known as **thread-local storage** or **TLS.** Note that this is more like static data than local variables, because it does not cease to exist when the function ends.

**Scheduler Activations**



- Many implementations of threads provide a virtual processor as an interface between the user thread and the kernel thread, particularly for the many-to-many or two-tier models.
  - This virtual processor is known as a "Lightweight Process", LWP.
    - There is a one-to-one correspondence between LWPs and kernel threads.
    - The number of kernel threads available, ( and hence the number of LWPs ) may change dynamically.
    - The application ( user level thread library ) maps user threads onto available LWPs.
    - kernel threads are scheduled onto the real processor(s) by the OS.
  - The kernel communicates to the user-level thread library when certain events occur ( such as a thread about to block ) via an **upcall**, which is handled in the thread library by an **upcall handler**. The upcall also provides a new LWP for the upcall handler to run on, which it can then use to reschedule the user thread that is about to become blocked. The OS will also issue upcalls when a thread becomes unblocked, so the thread library can make appropriate adjustments.
- If the kernel thread blocks, then the LWP blocks, which blocks the user thread.
- Ideally there should be at least as many LWPs available as there could be concurrently blocked kernel threads. Otherwise if all LWPs are blocked, then user threads will have to wait for one to become available.

# PROCESS SCHEDULING

## 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.
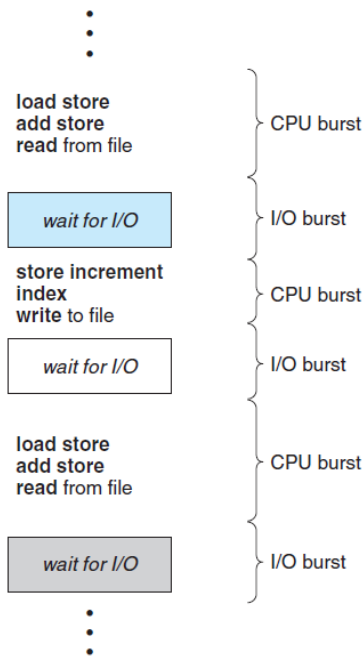
## Basic Concepts



**Figure 6.1** Alternating sequence of CPU and I/O bursts.

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.

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.

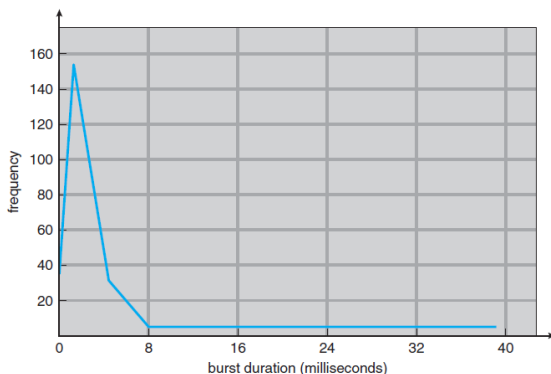## CPU – I/O Burst Cycle



**Figure 6.2** Histogram of CPU-burst durations.

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

78

burst ends with a system request to terminate execution (Figure 6.1).

The durations of CPU bursts have been measured extensively. Although they vary greatly from process to process and from computer to computer, they tend to have a frequency curve similar to that shown in Figure 6.2. The curve is generally characterized as exponential or hyper exponential, with a large number of short CPU bursts and a small number of long CPU bursts.

An I/O-bound program typically has many short CPU bursts. A CPU-bound program might have a few long CPU bursts. This distribution can be important in the selection of an appropriate CPU-scheduling algorithm.

**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.

**Preemptive Scheduling**

CPU-scheduling decisions may take place under the following four circum-stances:

When a process switches from the running state to the waiting state (for example, as the result of an I/O request or an invocation of wait() for the termination of a child process)

When a process switches from the running state to the ready state (for example, when an interrupt occurs)

When a process switches from the waiting state to the ready state (for example, at completion of I/O)

When a process terminates

For situations 1 and 4, there is no choice in terms of scheduling. A new process (if one exists in the ready queue) must be selected for execution. There is a choice, however, for situations 2 and 3.

When scheduling takes place only under circumstances 1 and 4, we say that the scheduling scheme is **nonpreemptive** or **cooperative**. Otherwise, it is **preemptive**. Under nonpreemptive scheduling, once the

CPU has been allocated to a process, the process keeps the CPU until it releases the CPU either by terminating or by switching to the waiting state. This scheduling method was used by Microsoft Windows 3.x. Windows 95 introduced preemptive scheduling, and all subsequent versions of Windows operating systems have used preemptive scheduling. The Mac OS X operating system for the Macintosh also uses preemptive scheduling; previous versions of the Macintosh operating system relied on cooperative scheduling. Cooperative scheduling is the only method that can be used on certain hardware platforms, because it does not require the special hardware (for example, a timer) needed for preemptive scheduling.

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**.

**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 algorithms. 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. Conceptually, 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.

## CPU-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.
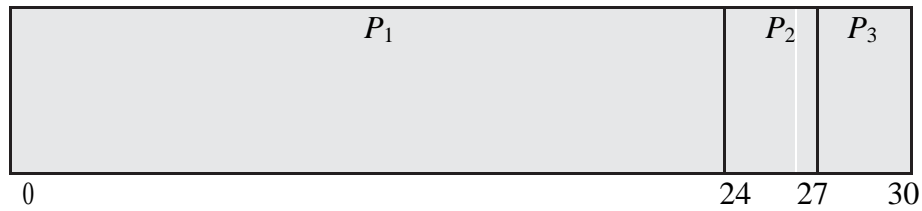
## First-Come, First-Served Scheduling

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:
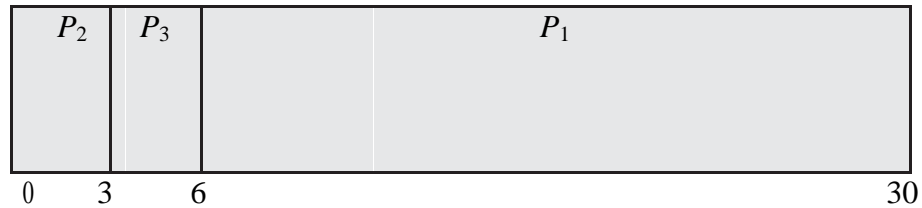
| Process | Burst Time |
|---------|-----------|
| $P_1$   | 24        |

| | |
|---|---|
| $P_2$ | 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.

**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.

## 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 $P_1$, $P_2$, $\cdots$, $P_5$, with the length of the CPU burst given in milliseconds:

| Process | Burst Time | Priority |
|:---:|:---:|:---:|
| $P1$ | 10 | 3 |
| $P_2$ | 1 | 1 |
| $P_3$ | 2 | 4 |
| $P_4$ | 1 | 5 |

83

$P5 \qquad 5 \qquad 2$

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.

## 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.
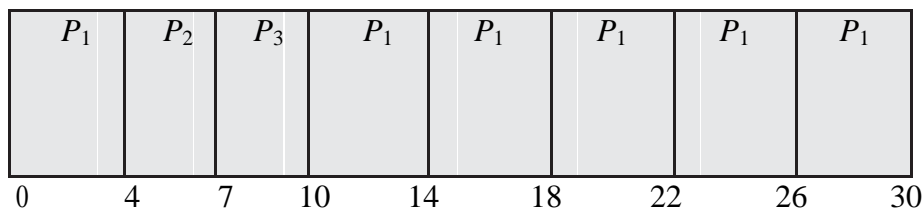
One of two things will then happen. The process may have a CPU burst of less than 1 time quantum. In this case, the process itself will release the CPU voluntarily. The scheduler will then proceed to the next process in the ready queue. If the CPU burst of the currently running process is longer than 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 $(n - 1) \times q$ 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.
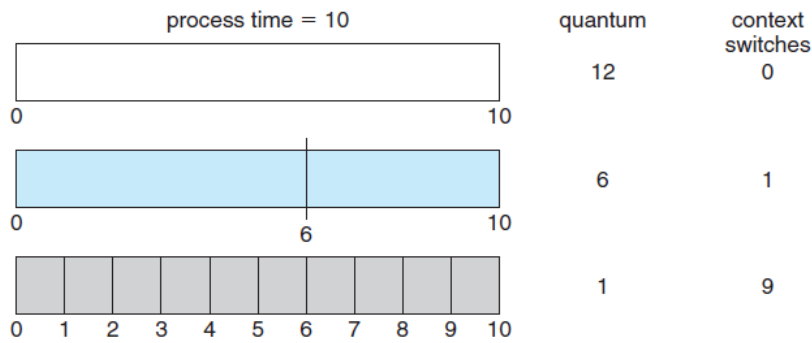
Figure 6.4 How a smaller time quantum increases context switches.

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 6.4).

## Multiple-Level Queues Scheduling

Multiple-level queues are not an independent scheduling algorithm. They make use of other existing algorithms to group and schedule jobs with common characteristics.

- Multiple queues are maintained for processes with common characteristics.
- Each queue can have its own scheduling algorithms.
- Priorities are assigned to each queue.

For example, CPU-bound jobs can be scheduled in one queue and all I/O-bound jobs in another queue. The Process Scheduler then alternately selects jobs from each queue and assigns them to the CPU based on the algorithm assigned to the queue.

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

System processes
Interactive processes
Interactive editing processes
Batch processes
Student processes



Figure 6.6 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.

## 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



**Figure 6.7** Multilevel feedback queues.

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 6.7). 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.

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.

# PROCESS SYNCHRONIZATION

A **cooperating process** is one that can affect or be affected by other processes executing in the system. Cooperating processes can either directly share a logical address space (that is, both code and data) or be allowed to share data only through files or messages. The former case is achieved through the use of threads,

## The Critical-Section Problem

We begin our consideration of process synchronization by discussing the so-called critical-section problem. Consider a system consisting of $n$ processes $\{P_0, P_1, ..., P_{n-1}\}$. Each process has a segment of code, called a **critical section**, in which the process may be changing common variables, updating a table, writing a file, and soon.

The important feature of the system is that, when one process is executing in its critical section, no other process is allowed to execute in its critical section. That is, no two processes are executing in their critical sections at the same time.

```
do {

    entry section

    critical section

    exit section

    remainder section

} while (true);
```

**Figure 5.1** General structure of a typical process $P_i$.

The **critical-section problem** is to design a protocol that the processes can use to cooperate. Each process must request permission to enter its critical section.

The section of code implementing this request is the **entry section**.

The critical section may be followed by an **exit section**.

The remaining code is the **remainder section**.

The general structure of a typical process $P_i$ is shown in Figure 5.1. The entry section and exit section are enclosed in boxes to highlight these important segments of code.

A solution to the critical-section problem must satisfy the following three requirements:

**Mutual exclusion**. If process $P_i$ is executing in its critical section, then no other processes can be executing in their critical sections.

**Progress**. If no process is executing in its critical section and some processes wish to enter their critical sections, then only those processes that are not executing in their remainder

sections can participate in deciding which will enter its critical section next, and this selection cannot be postponed indefinitely.

**Bounded waiting**. There exists a bound, or limit, on the number of times that other processes are allowed to enter their critical sections after a process has made a request to enter its critical section and before that request is granted.

We assume that each process is executing at a nonzero speed. However, we can make no assumption concerning the relative speed of the $n$ processes.

At a given point in time, many kernel-mode processes may be active in the operating system. As a result, the code implementing an operating system (***kernel code***) is subject to several possible race conditions. Consider as an example a kernel data structure that maintains a list of all open files in the system. This list must be modified when a new file is opened or closed (adding the file to the list or removing it from the list) . If two processes were to open files simultaneously, the separate updates to this list could result in a race condition. Other kernel data structures that are prone to possible race conditions include structures for maintaining memory allocation, for maintaining process lists, and for interrupt handling. It is up to kernel developers to ensure that the operating system is free from such race conditions.

Two general approaches are used to handle critical sections in operating systems:

**preemptive kernels**. A preemptive kernel allows a process to be preempted while it is running in kernel mode.

A **nonpreemptive kernel** does not allow a process running in kernel mode to be preempted; a kernel- mode process will run until it exits kernel mode, blocks, or voluntarily yields control of the CPU.

Obviously, a nonpreemptive kernel is essentially free from race conditions on kernel data structures, as only one process is active in the kernel at a time. We cannot say the same about preemptive kernels, so they must be carefully designed to ensure that shared kernel data are free from race conditions. Preemptive kernels are especially difficult to design for SMP architectures, since in these environments it is possible for two kernel-mode processes to run simultaneously on different processors.

## Peterson's Solution

The classic software-based solution to the critical-section problem known as **Peterson's solution**. Because of the way modern computer architectures perform basic machine-language instructions, such as load and store, there are no guarantees that Peterson's solution will work correctly on such architectures. However, we present the solution because it provides a good algorithmic description of solving the critical-section problem and illustrates some of the complexities involved in designing software that addresses the requirements of mutual exclusion, progress, and bounded waiting

```
do {

    flag[i] = true;
    turn = j;
    while (flag[j] && turn == j);

        critical section

    flag[i] = false;

        remainder section

} while (true);
```

**Figure 5.2** The structure of process $P_i$ in Peterson's solution.

Peterson's solution is restricted to two processes that alternate execution between their critical sections and remainder sections. The processes are numbered $P_0$ and $P_1$. For convenience, when presenting $P_i$, we use $P_j$ to denote the other process; that is, j equals $1 - i$.

Peterson's solution requires the two processes to share two data items:

**int turn;**

**boolean flag[2];**

The variable turn indicates whose turn it is to enter its critical section. That is, if turn == i, then process $P_i$ is allowed to execute in its critical section. The flag array is used to indicate if a process is ready to enter its critical section. For example, if flag[i] is true, this value indicates that $P_i$ is ready to enter its critical section. With an explanation of these data structures complete, we are now ready to describe the algorithm shown in Figure 5.2.

To enter the critical section, process $P_i$ first sets flag[i] to be true and then sets turn to the value j , thereby asserting that if the other process wishes to enter the critical section, it can do so. If both processes try to enter at the same time, turn will be set to both i and j at roughly the same time. Only one of these assignments will last; the other will occur but will be overwritten immediately. The eventual value of turn determines which of the two processes is allowed to enter its critical section first.

We now prove that this solution is correct. We need to show that:

Mutual exclusion is preserved.

The progress requirement is satisfied.

The bounded-waiting requirement is met.

91

To prove property 1, we note that each $P_i$ enters its critical section only if either flag[j] == false or turn == i. Also note that, if both processes can be executing in their critical sections at the same time, then flag[0] == flag[1] == true. These two observations imply that $P_0$ and $P_1$ could not have successfully executed their while statements at about the same time, since the value of turn can be either 0 or 1 but cannot be both. Hence, one of the processes — say, $P_j$ — must have successfully executed the while statement, whereas $P_i$ had to execute at least one additional statement ("turn == j"). However, at that time, flag[j] == true and turn == j, and this condition will persist as long as $P_j$ is in its critical section; as a result, mutual exclusion is preserved.

## *Synchronization Hardware*

- To generalize the solution(s) expressed above, each process when entering their critical section must set some sort of ***lock***, to prevent other processes from entering their critical sections simultaneously, and must release the lock when exiting their critical section, to allow other processes to proceed. Obviously it must be possible to attain the lock only when no other process has already set a lock. Specific implementations of this general procedure can get quite complicated, and may include hardware solutions as outlined in this section.

```
boolean TestAndSet(boolean *target) {
    boolean rv = *target;
    *target = TRUE;
    return rv;
}
```

**Figure 5.3** The definition of the `TestAndSet()` instruction.

```
do {
    while (TestAndSetLock(&lock))
        ; // do nothing

        // critical section

    lock = FALSE;

        // remainder section
}while (TRUE);
```

**Figure 5.4** Mutual-exclusion implementation with `TestAndSet()`.

- One simple solution to the critical section problem is to simply prevent a process from being interrupted while in their critical section, which is the approach taken by non preemptive kernels. Unfortunately this does not work well in multiprocessor environments, due to the difficulties in disabling and the re-enabling interrupts on all processors. There is also a question as to how this approach affects timing if the clock interrupt is disabled.

```
int compare_and_swap(int *value, int expected, int new_value) {
    int temp = *value;

    if (*value == expected)
        *value = new_value;

    return temp;
}
```

**Figure 5.5** The definition of the `compare_and_swap()` instruction.

```
do {
    while (compare_and_swap(&lock, 0, 1) != 0)
        ; /* do nothing */

    /* critical section */

    lock = 0;

    /* remainder section */
} while (true);
```

**Figure 5.6** Mutual-exclusion implementation with the `compare_and_swap()` instruction.

- Another approach is for hardware to provide certain ***atomic*** operations. These operations are guaranteed to operate as a single instruction, without interruption. One such operation is the "Test and Set", which simultaneously sets a boolean lock variable and returns its previous

value, as shown in Figures 5.3 and 5.4:

Another variation on the test-and-set is an atomic swap of two booleans, as shown in Figures 5.5 and 5.6:

- The above examples satisfy the mutual exclusion requirement, but unfortunately do not guarantee bounded waiting. If there are multiple processes trying to get into their critical sections, there is no guarantee of what order they will enter, and any one process could have the bad luck to wait forever until they got their turn in the critical section. ( Since there is no guarantee as to the relative *rates* of the processes, a very fast process could theoretically release the lock, whip through their remainder section, and re-lock the lock before a slower process got a chance. As more

```
do {
   waiting[i] = TRUE;
   key = TRUE;
   while (waiting[i] && key)
      key = TestAndSet(&lock);
   waiting[i] = FALSE;

      // critical section

   j = (i + 1) % n;
   while ((j != i) && !waiting[j])
      j = (j + 1) % n;

   if (j == i)
      lock = FALSE;
   else
      waiting[j] = FALSE;

      // remainder section
}while (TRUE);
```

  and more processes are involved vying for the same resource, the odds of a slow process getting locked out completely increase. )
- Figure 5.7 illustrates a solution using test-and-set that does satisfy this requirement, using two shared data structures, `boolean lock` and `boolean waiting[ N ]`, where N is the number of processes in contention for critical sections:


The key feature of the above algorithm is that a process blocks on the AND of the critical section being locked and that this process is in the waiting state. When exiting a critical section, the exiting process does not just unlock the critical section and let the other processes have a free-for-all trying to get in. Rather it first looks in an orderly progression ( starting with the next process on the list ) for a process that has been waiting, and if it finds one, then it releases that particular process from its waiting state, without unlocking the critical section, thereby allowing a specific process into the critical section while continuing to block all the others. Only if there are no other processes currently waiting is the general lock removed, allowing the next process to come along access to the critical section.

- Unfortunately, hardware level locks are especially difficult to implement in multi-processor architectures. Discussion of such issues is left to books on advanced computer architecture.

## Mutex Locks :

- The hardware solutions presented above are often difficult for ordinary programmers to access, particularly on multi-processor machines, and particularly because they are often platform-dependent.

- Therefore most systems offer a software API equivalent called *mutex locks* or simply *mutexes.* ( For mutual exclusion )

- The terminology when using mutexes is to *acquire* a lock prior to entering a critical section, and to *release* **it when exiting, as shown in Figure**

```
do {

    acquire lock

        critical section

    release lock

        remainder section

} while (TRUE);
```

- Just as with hardware locks, the acquire step will block the process if the lock is in use by another process, and both the acquire and release operations are atomic.

- Acquire and release can be implemented as shown here, based on a boolean variable "available":

**Acquire:**
```
acquire() {
    while (!available)
        ; /* busy wait */
    available = false;;
}
```

**Release:**
```
release() {
    available = true;
}
```

- One problem with the implementation shown here, ( and in the hardware solutions presented earlier ), is the busy loop used to block processes in the acquire phase. These types of locks are referred to as *spinlocks*, because the CPU just sits and spins while blocking the process.

- Spinlocks are wasteful of cpu cycles, and are a really bad idea on single-cpu single-threaded machines, because the spinlock blocks the entire computer, and doesn't allow any other process to release the lock. ( Until the scheduler kicks the spinning process off of the cpu. )

- On the other hand, spinlocks do not incur the overhead of a context switch, so they are effectively used on multi-threaded machines when it is expected that the lock will be released after a short time.

### Semaphores

A **semaphore** S is an integer variable that, apart from initialization, is accessed only through two standard atomic operations: wait() and signal(). The wait() operation was originally termed P (from the Dutch *proberen,* "to test"); signal() was originally called V (from *verhogen,* "to increment"). The definition of wait() is as follows:

```
wait(S) {
        while (S <= 0)
        // busy wait
        S--;
        }
```

The definition of signal() is as follows:

```
signal(S) {
        S++;
        }
```

All modifications to the integer value of the semaphore in the wait() and signal() operations must be executed indivisibly. That is, when one process modifies the semaphore value, no other process can simultaneously modify that same semaphore value. In addition, in the case of wait(S), the testing of the integer value of S (S $\leq$ 0), as well as its possible modification (S-- ), must be executed without interruption. We shall see how these operations can be implemented in Section 5.6.2. First, let's see how semaphores can be used.

### Semaphore Usage

Operating systems often distinguish between counting and binary semaphores. The value of a **counting semaphore** can range over an unrestricted domain. The value of a **binary semaphore** can range only between 0 and 1. Thus, binary semaphores behave similarly to mutex locks. In fact, on systems that do not provide mutex locks, binary semaphores can be used instead for providing mutual exclusion.

Counting semaphores can be used to control access to a given resource consisting of a finite number of instances. The semaphore is initialized to the number of resources available. Each process that wishes to use a resource performs a wait() operation on the semaphore (thereby decrementing the count). When a process releases a resource, it performs a signal() operation (incrementing the count). When the count for the semaphore goes to 0, all resources are being used. After that, processes that wish to use a resource will block until the count becomes greater than 0.

We can also use semaphores to solve various synchronization problems. For example, consider two concurrently running processes: $P_1$ with a statement $S_1$ and $P_2$ with a statement $S_2$ . Suppose we

require that $S_2$ be executed only after $S_1$ has completed. We can implement this scheme readily by letting $P_1$ and $P_2$ share a common semaphore synch, initialized to 0. In process $P_1$, we insert the statements

$S_1$ ;
signal(synch);

In process $P_2$ , we insert the statements

wait(synch);
$S_2$ ;

Because synch is initialized to 0, $P_2$ will execute $S_2$ only after $P_1$ has invoked signal(synch), which is after statement $S_1$ has been executed.


**Semaphore Implementation**

The definitions of the wait() and signal() semaphore operations just described present the same problem. To overcome the need for busy waiting, we can modify the definition of the wait() and signal() operations as follows: When a process executes the wait() operation and finds that the semaphore value is not positive, it must wait. However, rather than engaging in busy waiting, the process can block itself. The block operation places a process into a waiting queue associated with the semaphore, and the state of the process is switched to the waiting state. Then control is transferred to the CPU scheduler, which selects another process to execute.

A process that is blocked, waiting on a semaphore S, should be restarted when some other process executes a signal() operation. The process is restarted by a wakeup() operation, which changes the process from the waiting state to the ready state. The process is then placed in the ready queue. (The CPU may or may not be switched from the running process to the newly ready process, depending on the CPU-scheduling algorithm.)

To implement semaphores under this definition, we define a semaphore as follows:

```
typedef struct {
    int value;
    struct process *list;
} semaphore;
```

Each semaphore has an integer value and a list of processes list. When a process must wait on a semaphore, it is added to the list of processes. A signal() operation removes one process from the list of waiting processes and awakens that process.

Now, the wait() semaphore operation can be defined as

```
wait(semaphore *S) {
        S->value--;
        if (S->value < 0) {
                add this process to S->list;
                block();
                }
        }
```

and the signal() semaphore operation can be defined as

```
signal(semaphore *S) {
        S->value++;
                if (S->value <= 0) {
                remove a process P from S->list;
                wakeup(P);
                        }
                }
```

The block() operation suspends the process that invokes it. The wakeup(P) operation resumes the execution of a blocked process P. These two operations are provided by the operating system as basic system calls.

Note that in this implementation, semaphore values may be negative, whereas semaphore values are never negative under the classical definition of semaphores with busy waiting. If a semaphore value is negative, its magnitude is the number of processes waiting on that semaphore. This fact results from switching the order of the decrement and the test in the implementation of the wait() operation.

**Deadlocks and Starvation**

The implementation of a semaphore with a waiting queue may result in a situation where two or more processes are waiting indefinitely for an event that can be caused only by one of the waiting processes. The event in question is the execution of a signal() operation. When such a state is reached, these processes are said to be **deadlocked**.

To illustrate this, consider a system consisting of two processes, $P_0$ and $P_1$, each accessing two semaphores, S and Q, set to the value 1:

```
        P0                  P1
    wait(S);            wait(Q);
    wait(Q);            wait(S);
        .                  .
        .                  .
        .                  .
    signal(S);          signal(Q);
    signal(Q);          signal(S);
```

Suppose that $P_0$ executes wait(S) and then $P_1$ executes wait(Q). When $P_0$ executes wait(Q), it must wait until $P_1$ executes signal(Q). Similarly, when $P_1$ executes wait(S) , it must wait until $P_0$ executes signal(S). Since these signal() operations cannot be executed, $P_0$ and $P_1$ are deadlocked.

We say that a set of processes is in a deadlocked state when every process in the set is waiting for an event that can be caused only by another process in the set. The events with which we are mainly concerned here are resource acquisition and release. Other types of events may result in deadlocks, as we show in Chapter 7. In that chapter, we describe various mechanisms for dealing with the deadlock problem.

Another problem related to deadlocks is **indefinite blocking** or **starvation**, a situation in which processes wait indefinitely within the semaphore. Indefinite blocking may occur if we remove processes from the list associated with a semaphore in LIFO (last-in, first-out) order.

**Priority Inversion**

A scheduling challenge arises when a higher-priority process needs to read or modify kernel data that are currently being accessed by a lower-priority process — or a chain of lower-priority processes. Since kernel data are typically protected with a lock, the higher-priority process will have to wait for a lower-priority one to finish with the resource. The situation becomes more complicated if the lower-priority process is preempted in favor of another process with a higher priority.

As an example, assume we have three processes — $L$ , $M$, and $H$ — whose priorities follow the order $L < M < H$. Assume that process $H$ requires resource $R$, which is currently being accessed by process $L$. Ordinarily, process $H$ would wait for $L$ to finish using resource $R$. However, now suppose that process $M$ becomes runnable, thereby preempting process $L$. Indirectly, a process with a lower priority — process $M$— has affected how long process $H$ must wait for $L$ to relinquish resource $R$.

This problem is known as **priority inversion**. It occurs only in systems with more than two priorities, so one solution is to have only two priorities. That is insufficient for most general-purpose operating systems, however. Typically these systems solve the problem by implementing a **priority-inheritance protocol**.

According to this protocol, all processes that are accessing resources needed by a higher-priority process inherit the higher priority until they are finished with the resources in question. When they are finished, their priorities revert to their original values. In the example above, a priority-inheritance protocol would allow process $L$ to temporarily inherit the priority of process $H$, thereby preventing process $M$ from preempting its execution. When process had finished using resource $R$, it would

relinquish its inherited priority from $H$ and assume its original priority. Because resource $R$ would now be available, process $H$ — not $M$— would run next.

# Monitors

- Semaphores can be very useful for solving concurrency problems, *but only if programmers use them properly.* If even one process fails to abide by the proper use of semaphores, either accidentally or deliberately, then the whole system breaks down. ( And since concurrency problems are by definition rare events, the problem code may easily go unnoticed and/or be heinous to debug. )
- For this reason a higher-level language construct has been developed, called *monitors*.

```
monitor  monitor name
{
    // shared variable declarations
    procedure P1 ( . . . ) {
        . . .
    }
    procedure P2 ( . . . ) {
        . . .
    }
            .
            .
            .
    procedure Pn ( . . . ) {
        . . .
    }
    initialization code ( . . . ) {
        . . .
    }
}
```

**Figure - Syntax of a monitor.**

## Monitor Usage

- A monitor is essentially a class, in which all data is private, and with the special restriction that only one method within any given monitor object may be active at the same time. An additional restriction is that monitor methods may only access the shared data within the monitor and any data passed to them as parameters. I.e. they cannot access any data external to the monitor.



**Figure - Schematic view of a monitor**



- In order to fully realize the potential of monitors, we need to introduce one additional new data type, known as a *condition*.
  - A variable of type condition has only two legal operations, *wait* and *signal*. I.e. if X was defined as type condition, then legal operations would be X.wait( ) and X.signal( )
  - The wait operation blocks a process

101

until some other process calls signal, and adds the blocked process onto a list associated with that condition.

- o The signal process does nothing if there are no processes waiting on that condition. Otherwise it wakes up exactly one process from the condition's list of waiting processes. ( Contrast this with counting semaphores, which always affect the semaphore on a signal call. )

But now there is a potential problem - If process P within the monitor issues a signal that would wake up process Q also within the monitor, then there would be two processes running simultaneously within the monitor, violating the exclusion requirement. Accordingly there are two possible solutions to this dilemma:

**Signal and wait** - When process P issues the signal to wake up process Q, P then waits, either for Q to leave the monitor or on some other condition.

**Signal and continue** - When P issues the signal, Q waits, either for P to exit the monitor or for some other condition.

There are arguments for and against either choice. Concurrent Pascal offers a third alternative - The signal call causes the signaling process to immediately exit the monitor, so that the waiting process can then wake up and proceed.

- Java and C# ( C sharp ) offer monitors bulit-in to the language. Erlang offers similar but different constructs.

## CLASSIC PROBLEMS OF SYNCHRONIZATION

In this section, we present a number of synchronization problems as examples of a large class of concurrency-control problems. These problems are used for testing nearly every newly proposed synchronization scheme. In our solutions to the problems, we use semaphores for synchronization, since that is the traditional way to present such solutions. However, actual implementations of these solutions could use mutex locks in place of binary semaphores

**The Bounded-Buffer Problem**

The ***bounded-buffer problem*** is commonly used to illustrate the power of synchronization primitives. Here, we present a general structure of this scheme without committing ourselves to any particular implementation. We provide a related programming project in the exercises at the end of the chapter.

In our problem, the producer and consumer processes share the following data structures:

**int n;**
**semaphore mutex = 1;**
**semaphore empty = n;**
**semaphore full = 0**

We assume that the pool consists of n buffers, each capable of holding one item. The mutex semaphore provides mutual exclusion for accesses to the buffer pool and is initialized to the value 1. The empty and full semaphores count the number of empty and full buffers. The semaphore empty is initialized to the value n; the semaphore full is initialized to the value 0.

The code for the producer process is shown in Figure 5.9, and the code for the consumer process is shown in Figure 5.10. Note the symmetry between the producer and the consumer. We can interpret this code as the producer producing full buffers for the consumer or as the consumer producing empty buffers for the producer.

```
do {
    wait(full);
    wait(mutex);
      . . .
                                  /* remove an item from buffer to next consumed */

      . . .
    signal(mutex);
    signal(empty);
      . . .
    /* consume the item in next consumed */
      . . .
} while (true);
```

The structure of the consumer process.

**The Readers – Writers Problem**

Suppose that a database is to be shared among several concurrent processes. Some of these processes may want only to read the database, whereas others may want to update (that is, to read and write) the database. We distinguish between these two types of processes by referring to the former as ***readers*** and to the latter as ***writers***. Obviously, if two readers access the shared data simultaneously, no

adverse effects will result. However, if a writer and some other process (either a reader or a writer) access the database simultaneously, chaos may ensue.

To ensure that these difficulties do not arise, we require that the writers have exclusive access to the shared database while writing to the database. This synchronization problem is referred to as the **readers–writers problem**. Since it was originally stated, it has been used to test nearly every new synchronization primitive. The readers – writers problem has several variations, all involving priorities. The simplest one, referred to as the *first* readers – writers problem, requires that no reader be kept waiting unless a writer has already obtained permission to use the shared object. In other words, no reader should wait for other readers to finish simply because a writer is waiting. The *second* readers–writers problem requires that, once a writer is ready, that writer perform its write as soon as possible. In other words, if a writer is waiting to access the object, no new readers may start reading.

A solution to either problem may result in starvation. In the first case, writers may starve; in the second case, readers may starve. For this reason, other variants of the problem have been proposed. Next, we present a solution to the first readers – writers problem. See the bibliographical notes at the end of the chapter for references describing starvation-free solutions to the second readers – writers problem.

In the solution to the first readers – writers problem, the reader processes share the following data structures:

```
semaphore rw mutex = 1;
semaphore mutex = 1;
int read count = 0;
```

The semaphores mutex and rw mutex are initialized to 1; read count is initialized to 0. The semaphore rw mutex is common to both reader and writer

```
do {
    wait(rw mutex);

        . . .
    /* writing is performed */
        . . .
    signal(rw mutex);

} while (true);
```

**Figure 5.11**   The structure of a writer process.

processes. The mutex semaphore is used to ensure mutual exclusion when the variable read count is updated. The read count variable keeps track of how many processes are currently reading the object. The semaphore rw_mutex functions as a mutual exclusion semaphore for the writers. It is also used by the

first or last reader that enters or exits the critical section. It is not used by readers who enter or exit while other readers are in their critical sections.

The code for a writer process is shown in Figure 5.11; the code for a reader process is shown in Figure 5.12. Note that, if a writer is in the critical section and $n$ readers are waiting, then one reader is queued on rw mutex, and $n - 1$ readers are queued on mutex. Also observe that, when a writer executes signal(rw mutex), we may resume the execution of either the waiting readers or a single waiting writer. The selection is made by the scheduler.

The readers – writers problem and its solutions have been generalized to provide **reader–writer** locks on some systems. Acquiring a reader – writer lock requires specifying the mode of the lock: either *read* or *write* access. When a process wishes only to read shared data, it requests the reader – writer lock in read mode. A process wishing to modify the shared data must request the lock in write mode. Multiple processes are permitted to concurrently acquire a reader – writer lock in read mode, but only one process may acquire the lock for writing, as exclusive access is required for writers.

**The Dining-Philosophers Problem**

Consider five philosophers who spend their lives thinking and eating. The philosophers share a circular table surrounded by five chairs, each belonging to one philosopher. In the center of the table is a bowl of rice, and the table is laid with five single chopsticks (Figure 5.13).

When a philosopher thinks, she does not interact with her colleagues. From time to time, a philosopher gets hungry and tries to pick up the two chopsticks that are closest to her (the chopsticks that are between her and her left and right neighbors).

A philosopher may pick up only one chopstick at a time. Obviously, she cannot pick up a chopstick that is already in the hand of a neighbor. When a hungry philosopher has both her chopsticks at the same time, she eats without releasing the chopsticks. When she is finished eating, she puts down both chopsticks and starts thinking again.

The **dining-philosophers problem** is considered a classic synchronization problem neither because of its practical importance nor because computer scientists dislike philosophers but because it is an example of a large class of concurrency-control problems. It is a simple representation of the need to allocate several resources among several processes in a deadlock-free and starvation-free manner.

**Figure 5.13** The situation of the dining philosophers.

One simple solution is to represent each chopstick with a semaphore. A philosopher tries to grab a chopstick by executing a wait() operation on that semaphore. She releases her chopsticks by executing the signal() operation on the appropriate semaphores. Thus, the shared data are

**semaphore chopstick[5];**

where all the elements of chopstick are initialized to 1.

The Structure of Philosopher i as follows:

```
do {
    wait(chopstick[i]);
    wait(chopstick[(i+1) % 5]);

      . . .
    /* eat for awhile */

      . . .
    signal(chopstick[i]);
    signal(chopstick[(i+1) % 5]);

      . . .
    /* think for awhile */

      . . .

} while (true);
```

Although this solution guarantees that no two neighbors are eating simultaneously, it nevertheless must be rejected because it could create a deadlock. Suppose that all five philosophers become hungry at the same time and each grabs her left chopstick. All the elements of chopstick will now be equal to 0. When each philosopher tries to grab her right chopstick, she will be delayed forever.

Several possible remedies to the deadlock problem are replaced by:

- Allow at most four philosophers to be sitting simultaneously at the table.
- Allow a philosopher to pick up her chopsticks only if both chopsticks are available (to do this, she must pick them up in a critical section).
- Use an asymmetric solution — that is, an odd-numbered philosopher picks up first her left chopstick and then her right chopstick, whereas an even-numbered philosopher picks up her right chopstick and then her left chopstick.

we present a solution to the dining-philosophers problem that ensures freedom from deadlocks. Note, however, that any satisfactory solution to the dining-philosophers problem must guard against the possibility that one of the philosophers will starve to death. A deadlock-free solution does not necessarily eliminate the possibility of starvation.

# 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."**

Although some applications can identify programs that may deadlock, operating systems typically do not provide deadlock-prevention facilities, and it remains the responsibility of programmers to ensure that they design deadlock-free programs.

Deadlock problems can only become more common, given current trends, including larger numbers of processes, multithreaded programs, many more resources within a system, and an emphasis on long-lived file and database servers rather than batch systems.

## System Model

A system consists of a finite number of resources to be distributed among a number of competing processes. The resources may be partitioned into several types (or classes), each consisting of some number of identical instances. CPU cycles, files, and I/O devices (such as printers and DVD drives) are examples of resource types. If a system has two CPUs, then the resource type *CPU* has two instances. Similarly, the resource type *printer* may have five instances.

If a process requests an instance of a resource type, the allocation of *any* instance of the type should satisfy the request. If it does not, then the instances are not identical, and the resource type classes have not been defined properly.

For example, a system may have two printers. These two printers may be defined to be in the same resource class if no one cares which printer prints which output. However, if one printer is on the ninth floor and the other is in the basement, then people on the ninth floor may not see both printers as equivalent, and separate resource classes may need to be defined for each printer.

A process must request a resource before using it and must release the resource after using it. A process may request as many resources as it requires to carry out its designated task. Obviously, the number of resources requested may not exceed the total number of resources available in the system. In other words, a process cannot request three printers if the system has only two.

Under the normal mode of operation, a process may utilize a resource in only the following sequence:

- **Request**. The process requests the resource. If the request cannot be granted immediately (for example, if the resource is being used by another process), then the requesting process must wait until it can acquire the resource.
- **Use**. The process can operate on the resource (for example, if the resource is a printer, the process can print on the printer).
- **Release**. The process releases the resource.

The request and release of resources may be system calls,. **Examples are** the request() and release() device, open() and close() file, and allocate() and free() memory system calls.

A set of processes is in a deadlocked state when every process in the set is waiting for an event that can be caused only by another process in the set. The events with which we are mainly concerned here are resource acquisition and release. The resources may be either physical resources (for example, printers, tape drives, memory space, and CPU cycles) or logical resources (for example, semaphores, mutex locks, and files). However, other types of events may result in deadlocks (for example, the IPC facilities).

To illustrate a deadlocked state, consider a system with three CD RW drives. Suppose each of three processes holds one of these CD RW drives. If each process now requests another drive, the three processes will be in a deadlocked state. Each is waiting for the event "CD RW is released," which can be caused only by one of the other waiting processes. This example illustrates a deadlock involving the same resource type.

Deadlocks may also involve different resource types. For example, consider a system with one printer and one DVD drive. Suppose that process $P_i$ is holding the DVD and process $P_j$ is holding the printer. If $P_i$ requests the printer and $P_j$ requests the DVD drive, a deadlock occurs.

Developers of multithreaded applications must remain aware of the possibility of deadlocks. The locking tools presented in Chapter 5 are designed to avoid race conditions. However, in using these tools, developers must pay careful attention to how locks are acquired and released. Otherwise, deadlock can occur.

**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.

**Necessary Conditions**

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

- **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.
- **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.
- **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.
- **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.

**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 7.1 depicts the following situation.

The sets *P, R,* and *E*:

$P = \{P_1, P_2, P_3\}$
$R = \{R_1, R_2, R_3, R_4\}$
$\circ\ E = \{P_1 \rightarrow R_1, P_2 \rightarrow R_3, R_1 \rightarrow P_2, R_2 \rightarrow P_2, R_2 \rightarrow P_1, R_3 \rightarrow P_3\}$

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 $R_4$



**Figure 7.1** Resource-allocation graph.

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 $R_3$.

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 7.1. Suppose that process $P_3$ requests an instance of resource



**Figure 7.2** Resource-allocation graph with a deadlock.



**Figure 7.3** Resource-allocation graph with a cycle but no deadlock.

type $R_2$. Since no resource instance is currently available, we add a request edge $P_3 \rightarrow R_2$ to the graph (Figure 7.2). At this point, two minimal cycles exist in the system:

$$P_1 \rightarrow R_1 \rightarrow P_2 \rightarrow R_3 \rightarrow P_3 \rightarrow R_2 \rightarrow P_1 \; P_2 \rightarrow R_3 \rightarrow P_3 \rightarrow R_2 \rightarrow P_2$$

Processes $P_1$ , $P_2$, and $P_3$ are deadlocked. Process $P_2$ is waiting for the resource $R_3$ , which is held by process $P_3$. Process $P_3$ is waiting for either process $P_1$ or process $P_2$ to release resource $R_2$ . In addition, process $P_1$ is waiting for process $P_2$ to release resource $R_1$ .

Now consider the resource-allocation graph in Figure 7.3. In this example, we also have a cycle:

$$P_1 \rightarrow R_1 \rightarrow P_3 \rightarrow R_2 \rightarrow P_1$$

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.

**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.

**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.

**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.

**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 utilization 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.

**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.

**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.

To illustrate, we let $R = \{ R_1, R_2, ..., R_m\}$ be the set of resource types. We assign to each resource type a unique integer number, which allows us to compare two resources and to determine whether one precedes another in our ordering. Formally, we define a one-to-one function $F: R \rightarrow N,$ where $N$ is the set of natural numbers. For example, if the set of resource types $R$ includes tape drives, disk drives, and printers, then the function $F$ might be defined as follows:

$$F \text{ (tape drive)} = 1$$
$$F \text{ (disk drive)} = 5$$
$$F \text{ (printer)} = 12$$

We can now consider the following protocol to prevent deadlocks: Each process can request resources only in an increasing order of enumeration. That is, a process can initially request any number of instances of a resource type — say, $R_i$ . After that, the process can request instances of resource type $R_j$ if and only if $F( R_j ) > F( R_i )$. For example, using the function defined previously, a process that wants to use the tape drive and printer at the same time must first request the tape drive and then request the printer. Alternatively, we can require that a process requesting an instance of resource type $R_j$ must have released any resources $R_i$ such that $F( R_i ) \geq F( R_j )$. Note also that if several instances of the same resource type are needed, a *single* request for all of them must be issued.

**Deadlock Avoidance**

An alternative method for avoiding deadlocks is to require additional information about how resources are to be requested.

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. In the following sections, we explore two deadlock-avoidance algorithms.

**Safe State**



Figure 7.6 Safe, unsafe, and deadlocked state spaces.

A state is *safe* if the system can allocate resources to each process (up to its maximum) in some order and still avoid a deadlock. More formally, a system is in a safe state only if there exists a **safe sequence**.

A sequence of processes $<P_1, P_2, ..., P_n>$ is a safe sequence for the current allocation state if, for each $P_i$ , the resource requests that $P_i$ can still make can be satisfied by the currently available resources plus the resources held by all $P_j$ , with $j < i$.

In this situation, if the resources that $P_i$ needs are not immediately available, then $P_i$ can wait until all $P_j$ have finished. When they have finished, $P_i$ can obtain all of its needed resources, complete its designated task, return its allocated resources, and terminate. When $P_i$ terminates, $P_{i+1}$ can obtain its needed resources, and so on. If no such sequence exists, then the system state is said to be *unsafe.*

A safe state is not a deadlocked state. Conversely, a deadlocked state is an unsafe state. Not all unsafe states are deadlocks, however (Figure 7.6). An unsafe state *may* lead to a deadlock. As long as the state is safe, the operating system can avoid unsafe (and deadlocked) states. In an unsafe state, the operating system cannot prevent processes from requesting resources in such a way that a deadlock occurs. The behavior of the processes controls unsafe states.

To illustrate, we consider a system with twelve magnetic tape drives and three processes: $P_0$, $P_1$, and $P_2$ . Process $P_0$ requires ten tape drives, process $P_1$ may need as many as four tape drives, and process $P_2$ may need up to nine tape drives. Suppose that, at time $t_0$ , process $P_0$ is holding five tape drives, process $P_1$ is holding two tape drives, and process $P_2$ is holding two tape drives. (Thus, there are three free tape drives.)

|  | Maximum Needs | Current Needs |
|---|---|---|
| $P_0$ | 10 | 5 |
| $P_1$ | 4 | 2 |
| P2 | 9 | 2 |

At time $t_0$, the system is in a safe state. The sequence $<P_1, P_0, P_2>$ satisfies the safety condition. Process $P_1$ can immediately be allocated all its tape drives and then return them (the system will then have five available tape drives); then process $P_0$ can get all its tape drives and return them (the system will then have ten available tape drives); and finally process $P_2$ can get all its tape drives and return them (the system will then have all twelve tape drives available).

**Resource-Allocation-Graph Algorithm**



**Figure 7.7** Resource-allocation graph for deadlock avoidance.

If we have a resource-allocation system with only one instance of each resource type, we can use a variant of the resource-allocation graph for deadlock avoidance.

In addition to the request and assignment edges already described, we introduce a new type of edge, called a **claim edge**.

A claim edge $P_i \rightarrow R_j$ indicates that process $P_i$ may request resource $R_j$ at some time in the future. This edge resembles a request edge in direction but is represented in the graph by a dashed line. When process $P_i$ requests resource $R_j$, the claim edge $P_i \rightarrow R_j$ is converted to a request edge. Similarly, when a resource $R_j$ is released by $P_i$, the assignment edge $R_j \rightarrow P_i$ is reconverted to a claim edge $P_i \rightarrow R_j$.

Note that the resources must be claimed a priori in the system. That is, before process $P_i$ starts executing, all its claim edges must already appear in the resource-allocation graph. We can relax this condition by allowing a claim edge $P_i \rightarrow R_j$ to be added to the graph only if all the edges associated with process $P_i$ are claim edges.

Now suppose that process $P_i$ requests resource $R_j$. The request can be granted only if converting the request edge $P_i \rightarrow R_j$ to an assignment edge $R_j \rightarrow P_i$ does not result in the formation of a cycle in the resource-allocation graph. We check for safety by using a cycle-detection algorithm. An algorithm for detecting a cycle in this graph requires an order of $n^2$ operations, where $n$ is the number of processes in the system.



**Figure 7.8** An unsafe state in a resource-allocation graph.

If no cycle exists, then the allocation of the resource will leave the system in a safe state. If a cycle is found, then the allocation will put the system in an unsafe state. In that case, process $P_i$ will have to wait for its requests to be

117

satisfied.

To illustrate this algorithm, we consider the resource- allocation graph of Figure 7.7. Suppose that $P_2$ requests $R_2$. Although $R_2$ is currently free, we cannot allocate it to $P_2$, since this action will create a cycle in the graph (Figure 7.8). A cycle, as mentioned, indicates that the system is in an unsafe state. If $P_1$ requests $R_2$, and $P_2$ requests $R_1$, then a deadlock will occur.

## Banker's Algorithm

The resource-allocation-graph algorithm is not applicable to a resource-allocation system with multiple instances of each resource type. The deadlock-avoidance algorithm that we describe next is applicable to such a system but is less efficient than the resource-allocation graph scheme. This algorithm is commonly known as the **banker's algorithm.** The name was chosen because the algorithm could be used in a banking system to ensure that the bank never allocated its available cash in such a way that it could no longer satisfy the needs of all its customers.

When a new process enters the system, it must declare the maximum number of instances of each resource type that it may need. This number may not exceed the total number of resources in the system. When a user requests a set of resources, the system must determine whether the allocation of these resources will leave the system in a safe state. If it will, the resources are allocated; otherwise, the process must wait until some other process releases enough resources.

Several data structures must be maintained to implement the banker 's algorithm. These data structures encode the state of the resource-allocation system. We need the following data structures, where $n$ is the number of processes in the system and $m$ is the number of resource types:

**Available**. A vector of length $m$ indicates the number of available resources of each type. If *Available*[$j$] equals $k$, then $k$ instances of resource type $R_j$ are available.

**Max**. An $n \times m$ matrix defines the maximum demand of each process. If *Max*[$i$][$j$] equals $k$, then process $P_i$ may request at most $k$ instances of resource type $R_j$ .

**Allocation**. An $n \times m$ matrix defines the number of resources of each type currently allocated to each process. If *Allocation*[$i$][$j$] equals $k$, then process $P_i$ is currently allocated $k$ instances of resource type $R_j$ .

**Need**. An $n \times m$ matrix indicates the remaining resource need of each process. If **Need**[i][j] equals *k,* then process $P_i$ may need $k$ more instances of resource type $R_j$ to complete its task. Note that **Need**[i][j] equals **Max**[i][j] − **Allocation**[i][j].

These data structures vary over time in both size and value.

To simplify the presentation of the banker 's algorithm, we next establish some notation. Let $X$ and $Y$ be vectors of length $n$. We say that $X \leq Y$ if and only if $X[i] \leq Y[i]$ for all $i = 1, 2, ..., n$. For example, if $X = (1,7,3,2)$ and $Y = (0,3,2,1)$, then $Y \leq X$. In addition, $Y < X$ if $Y \leq X$ and $Y \neq X$.

We can treat each row in the matrices **Allocation** and **Need** as vectors and refer to them as **Allocation**$_i$ and **Need**$_i$ . The vector **Allocation**$_i$ specifies the resources currently allocated to process $P_i$ ; the vector **Need**$_i$ specifies the additional resources that process $P_i$ may still request to complete its task.

## Safety Algorithm

We can now present the algorithm for finding out whether or not a system is in a safe state. This algorithm can be described as follows:

Let **Work** and **Finish** be vectors of length $m$ and $n,$ respectively. Initialize **Work** = **Available** and **Finish**[i] = *false* for $i = 0, 1, ..., n − 1$.

Find an index $i$ such that both

**Finish**[i] == *false*
**Need**$_i$ ≤ **Work**

If no such $i$ exists, go to step 4.
**Work** = **Work** + **Allocation**$_i$ **Finish**[i] = *true*
Go to step 2.

If **Finish**[i] == *true* for all $i,$ then the system is in a safe state.

This algorithm may require an order of $m \times n^2$ operations to determine whether a state is safe.

## Resource-Request Algorithm

Next, we describe the algorithm for determining whether requests can be safely granted.

Let **Request**$_i$ be the request vector for process $P_i$ . If **Request**$_i$ [ j ] == $k$, then process $P_i$ wants $k$ instances of resource type $R_j$ . When a request for resources is made by process $P_i$ , the following actions are taken:

If **Request**$_i$ ≤ **Need**$_i$ , go to step 2. Otherwise, raise an error condition, since the process has exceeded its maximum claim.

If **Request**$_i$ ≤ **Available,** go to step 3. Otherwise, $P_i$ must wait, since the resources are not available.

Have the system pretend to have allocated the requested resources to process $P_i$ by modifying the state as follows:

$$Available = Available - Request_i;$$
$$Allocation_i = Allocation_i + Request_i;$$
$$Need_i = Need_i - Request_i;$$

If the resulting resource-allocation state is safe, the transaction is com-pleted, and process $P_i$ is allocated its resources. However, if the new state is unsafe, then $P_i$ must wait for **Request$_i$**, and the old resource-allocation state is restored.

## An Illustrative Example

- Consider the following situation:

| | Allocation | Max | Available | Need |
|---|---|---|---|---|
| | A B C | A B C | A B C | A B C |
| $P_0$ | 0 1 0 | 7 5 3 | 3 3 2 | 7 4 3 |
| $P_1$ | 2 0 0 | 3 2 2 | | 1 2 2 |
| $P_2$ | 3 0 2 | 9 0 2 | | 6 0 0 |
| $P_3$ | 2 1 1 | 2 2 2 | | 0 1 1 |
| $P_4$ | 0 0 2 | 4 3 3 | | 4 3 1 |

The system is in a safe state since the sequence < P1, P3, P4, P2, P0> satisfies safety criteria

- And now consider what happens if process P1 requests 1 instance of A and 2 instances of C. ( Request[ 1 ] = ( 1, 0, 2 ) )

| | Allocation | Need | Available |
|---|---|---|---|
| | A B C | A B C | A B C |
| $P_0$ | 0 1 0 | 7 4 3 | 2 3 0 |
| $P_1$ | 3 0 2 | 0 2 0 | |
| $P_2$ | 3 0 2 | 6 0 0 | |
| $P_3$ | 2 1 1 | 0 1 1 | |
| $P_4$ | 0 0 2 | 4 3 1 | |

- What about requests of ( 3, 3,0 ) by P4? or ( 0, 2, 0 ) by P0? Can these be safely granted? Why or why not?

## Deadlock Detection

If a system does not employ either a deadlock-prevention or a deadlock-avoidance algorithm, then a deadlock situation may occur. In this environment, the system may provide:

- An algorithm that examines the state of the system to determine whether a deadlock has occurred

- An algorithm to recover from the deadlock

In the following discussion, we elaborate on these two requirements as they pertain to systems with only a single instance of each resource type, as well as to systems with several instances of each resource type. At this point, however, we note that a detection-and-recovery scheme requires overhead that includes not only the run-time costs of maintaining the necessary information and executing the detection algorithm but also the potential losses inherent in recovering from a deadlock.

**Single Instance of Each Resource Type**

If all resources have only a single instance, then we can define a deadlock-detection algorithm that uses a variant of the resource-allocation graph, called a **wait-for** graph. We obtain this graph from the resource-allocation graph by removing the resource nodes and collapsing the appropriate edges.



**Figure 7.9** (a) Resource-allocation graph. (b) Corresponding wait-for graph.

More precisely, an edge from $P_i$ to $P_j$ in a wait-for graph implies that process $P_i$ is waiting for process $P_j$ to release a resource that $P_i$ needs. An edge $P_i \rightarrow P_j$ exists in a wait-for graph if and only if the corresponding resource-allocation graph contains two edges $P_i \rightarrow R_q$ and $R_q \rightarrow P_j$ for some resource $R_q$. In Figure 7.9, we present a resource-allocation graph and the corresponding wait-for graph.

As before, a deadlock exists in the system if and only if the wait-for graph contains a cycle. To detect deadlocks, the system needs to *maintain* the wait-for graph and periodically *invoke an algorithm* that searches for a cycle in the graph. An algorithm to detect a cycle in a graph requires an order of $n^2$ operations, where $n$ is the number of vertices in the graph.

**Several Instances of a Resource Type**

The wait-for graph scheme is not applicable to a resource-allocation system with multiple instances of each resource type. We turn now to a deadlock-detection algorithm that is applicable to such a system. The algorithm employs several time-varying data structures that are similar to those used in the banker's algorithm:

- **Available**. A vector of length $m$ indicates the number of available resources of each type.

- **Allocation**. An $n \times m$ matrix defines the number of resources of each type currently allocated to each process.

- **Request**. An $n \times m$ matrix indicates the current request of each process. If **Request**$[i][j]$ equals $k$, then process $P_i$ is requesting $k$ more instances of resource type $R_j$.

The $\leq$ relation between two vectors is defined as in Section 7.5.3. To simplify notation, we again treat the rows in the matrices *Allocation* and *Request* as vectors; we refer to them as *Allocation$_i$* and *Request$_i$*. The detection algorithm described here simply investigates every possible allocation sequence for the processes that remain to be completed. Compare this algorithm with the banker 's algorithm.

Let *Work* and *Finish* be vectors of length $m$ and $n$, respectively. Initialize *Work* = *Available.* For $i = 0, 1, ..., n - 1,$ if **Allocation**$_i \neq 0$, then **Finish**$[i]$ = *false.* Otherwise, **Finish**$[i]$ = *true.*
Find an index $i$ such that both

  **Finish**$[i]$ == *false*

  **Request**$_i \leq$ **Work**

If no such $i$ exists, go to step 4.

 **Work** = **Work** + **Allocation**$_i$ **Finish**$[i]$ = *true*
 Go to step 2.

 If **Finish**$[i]$ == *false* for some $i$, $0 \leq i < n$, then the system is in a deadlocked state. Moreover, if **Finish**$[i]$ == *false*, then process $P_i$ is deadlocked.
This algorithm requires an order of $m \times n^2$ operations to detect whether the system is in a deadlocked state.

You may wonder why we reclaim the resources of process $P_i$ (in step 3) as soon as we determine that **Request**$_i \leq$ **Work** (in step 2b). We know that $P_i$ is currently ***not*** involved in a deadlock (since **Request**$_i \leq$ **Work**). Thus, we take an optimistic attitude and assume that $P_i$ will require no more resources to complete its task; it will thus soon return all currently allocated resources to the system. If our assumption is incorrect, a deadlock may occur later. That deadlock will be detected the next time the deadlock-detection algorithm is invoked.

To illustrate this algorithm, we consider a system with five processes $P_0$ through $P_4$ and three resource types $A$, $B$, and $C$. Resource type $A$ has seven instances, resource type $B$ has two instances, and resource type $C$ has six instances. Suppose that, at time $T_0$, we have the following resource-allocation state:

| | Allocation | Request | Available |
|---|---|---|---|
| | A B C | A B C | A B C |
| $P_0$ | 0 1 0 | 0 0 0 | 0 0 0 |
| $P_1$ | 2 0 0 | 2 0 2 | |
| $P_2$ | 3 0 3 | 0 0 0 | |
| $P_3$ | 2 1 1 | 1 0 0 | |
| $P_4$ | 0 0 2 | 0 0 2 | |

We claim that the system is not in a deadlocked state. Indeed, if we execute our algorithm, we will find that the sequence <$P_0$, $P_2$, $P_3$, $P_1$, $P_4$> results in **Finish**[$i$] == *true* for all $i$.

Suppose now that process $P_2$ makes one additional request for an instance of type $C$. The **Request** matrix is modified as follows:

| | Request |
|---|---|
| | A B C |
| $P_0$ | 0 0 0 |
| $P_1$ | 2 0 2 |
| $P_2$ | 0 0 1 |
| $P_3$ | 1 0 0 |
| $P_4$ | 0 0 2 |

We claim that the system is now deadlocked. Although we can reclaim the resources held by process $P_0$, the number of available resources is not sufficient to fulfill the requests of the other processes. Thus, a deadlock exists, consisting of processes $P_1$, $P_2$, $P_3$, and $P_4$.

**Detection-Algorithm Usage**

When should we invoke the detection algorithm? The answer depends on two factors:

How **often** is a deadlock likely to occur?

How **many** processes will be affected by deadlock when it happens?

If deadlocks occur frequently, then the detection algorithm should be invoked frequently. Resources allocated to deadlocked processes will be idle until the deadlock can be broken. In addition, the number of processes involved in the deadlock cycle may grow.

Deadlocks occur only when some process makes a request that cannot be granted immediately. This request may be the final request that completes a chain of waiting processes. In the extreme, then, we can invoke the deadlock-detection algorithm every time a request for allocation cannot be granted immediately. In this case, we can identify not only the deadlocked set of processes but also the specific process that "caused" the deadlock. (In reality, each of the deadlocked processes is a link in the cycle in the resource graph, so all of them, jointly, caused the deadlock.) If there are many different resource

types, one request may create many cycles in the resource graph, each cycle completed by the most recent request and "caused" by the one identifiable process.

## Recovery from Deadlock

When a detection algorithm determines that a deadlock exists, several alter-natives are available. One possibility is to inform the operator that a deadlock has occurred and to let the operator deal with the deadlock manually. Another possibility is to let the system **recover** from the deadlock automatically.

There are two options for breaking a deadlock.

One is simply to abort one or more processes to break the circular wait.

The other is to preempt some resources from one or more of the deadlocked processes.

## Process Termination

To eliminate deadlocks by aborting a process, we use one of two methods. In both methods, the system reclaims all resources allocated to the terminated processes.

- **Abort all deadlocked processes**.
  This method clearly will break the deadlock cycle, but at great expense. The deadlocked processes may have computed for a long time, and the results of these partial computations must be discarded and probably will have to be recomputed later.
- **Abort one process at a time until the deadlock cycle is eliminated**.
  This method incurs considerable overhead, since after each process is aborted, a deadlock-detection algorithm must be invoked to determine whether any processes are still deadlocked.

Aborting a process may not be easy. If the process was in the midst of updating a file, terminating it will leave that file in an incorrect state. Similarly, if the process was in the midst of printing data on a printer, the system must reset the printer to a correct state before printing the next job.

If the partial termination method is used, then we must determine which deadlocked process (or processes) should be terminated. This determination is a policy decision, similar to CPU-scheduling decisions. The question is basically an economic one; we should abort those processes whose termination will incur the minimum cost.

Unfortunately, the term *minimum cost* is not a precise one. Many factors may affect which process is chosen, including:

What the priority of the process is
How long the process has computed and how much longer the process will compute before completing its designated task

How many and what types of resources the process has used (for example, whether the resources are simple to preempt)

How many more resources the process needs in order to complete

How many processes will need to be terminated

Whether the process is interactive or batch

## Resource Preemption

To eliminate deadlocks using resource preemption, we successively preempt some resources from processes and give these resources to other processes until the deadlock cycle is broken.

If preemption is required to deal with deadlocks, then three issues need to be addressed:

- **Selecting a victim**. Which resources and which processes are to be preempted? As in process termination, we must determine the order of preemption to minimize cost. Cost factors may include such parameters as the number of resources a deadlocked process is holding and the amount of time the process has thus far consumed.

- **Rollback**. If we preempt a resource from a process, what should be done with that process? Clearly, it cannot continue with its normal execution; it is missing some needed resource. We must roll back the process to some safe state and restart it from that state.

- **Starvation**. How do we ensure that starvation will not occur? That is, how can we guarantee that resources will not always be preempted from the same process?

In a system where victim selection is based primarily on cost factors, it may happen that the same process is always picked as a victim. As a result, this process never completes its designated task, a starvation situation any practical system must address. Clearly, we must ensure that a process can be picked as a victim only a (small) finite number of times. The most common solution is to include the number of rollbacks in the cost factor.

# Memory Management

**Background**

- Obviously memory accesses and memory management are a very important part of modern computer operation. Every instruction has to be fetched from memory before it can be executed, and most instructions involve retrieving data from memory or storing data in memory or both.

- The advent of multi-tasking Operating Systems compounds the complexity of memory management, because as processes are swapped in and out of the CPU, so must their code and data be swapped in and out of memory, all at high speeds and without interfering with any other processes.

- Shared memory, virtual memory, the classification of memory as read-only versus read-write, and concepts like copy-on-write forking all further complicate the issue.

**Basic Hardware**

- It should be noted that from the memory chips point of view, all memory accesses are equivalent. The memory hardware doesn't know what a particular part of memory is being used for, nor does it care. This is almost true of the OS as well, although not entirely.

- The CPU can only access its registers and main memory. It cannot, for example, make direct access to the hard drive, so any data stored there must first be transferred into the main memory chips before the CPU can work with it. (Device drivers communicate with their hardware via interrupts and "memory" accesses, sending short instructions for example to transfer data from the hard drive to a specified location in main memory. The disk controller monitors the bus for such instructions, transfers the data, and then notifies the CPU that the data is there with another interrupt, but the CPU never gets direct access to the disk.)

- Memory accesses to registers are very fast, generally one clock tick, and a CPU may be able to execute more than one machine instruction per clock tick.

- Memory accesses to main memory are comparatively slow, and may take a number of clock ticks to complete. This would require intolerable waiting by the CPU if it were not for an intermediary fast memory *cache* built into most modern CPUs. The basic idea of the cache is to transfer chunks of memory at a time from the main memory to the cache, and then to access individual memory locations one at a time from the cache.

- User processes must be restricted so that they only access memory locations that "belong" to that particular process. This is usually implemented using a base register and a limit register for each process, as shown in Figures 8.1 and 8.2 below.

- *Every* memory access made by a user process is checked against these two registers, and if a memory access is attempted outside the valid range, then a fatal error is generated. The OS obviously has access to all existing memory locations, as this is necessary to swap users' code and

data in and out of memory. It should also be obvious that changing the contents of the base and limit registers is a privileged activity, allowed only to the OS kernel.



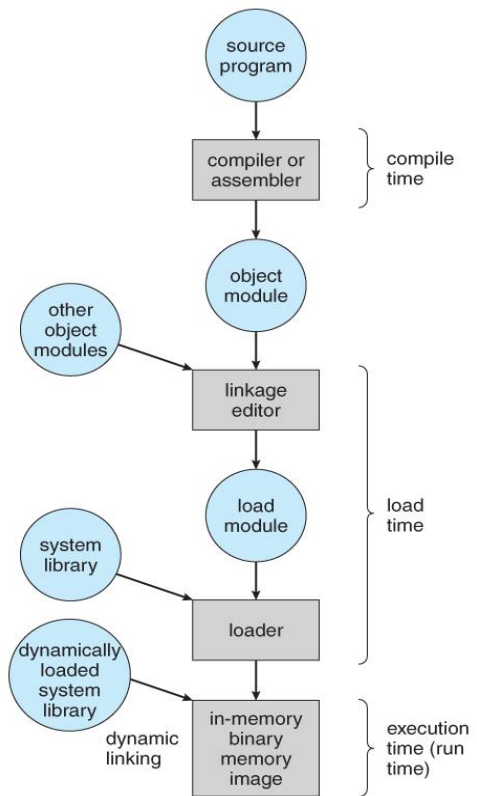**Figure 8.1 - A base and a limit register define a logical addresss space**



**Figure 8.2 - Hardware address protection with base and limit registers**

**Address Binding:**



- User programs typically refer to memory addresses with symbolic names such as "i", "count", and "average Temperature". These symbolic names must be mapped or **bound** to physical memory addresses, which typically occurs in several stages:

  - **Compile Time** - If it is known at compile time where a program will reside in physical memory, then **absolute code** can be generated by the compiler, containing actual physical addresses. However if the load address changes at some later time, then the program will have to be recompiled. DOS .COM programs use compile time binding.

  - **Load Time** - If the location at which a program will be loaded is not known at compile time, then the compiler must generate **relocatable code**, which references addresses relative to the start of the program. If that starting address changes, then the program must be reloaded but not recompiled.

  - **Execution Time** - If a program can be moved around in memory during the course of its execution, then binding must be delayed until execution time. This requires special hardware, and is the method implemented by most modern Operating Systems.

- Figure 8.3 shows the various stages of the binding processes and the units involved in each stage:

**Figure 8.3 - Multistep processing of a user program**

127

## Logical Versus Physical Address Space:

- The address generated by the CPU is a *logical address*, whereas the address actually seen by the memory hardware is a *physical address*.
- Addresses bound at compile time or load time have identical logical and physical addresses.
- Addresses created at execution time, however, have different logical and physical addresses.
  - In this case the logical address is also known as a *virtual address*, and the two terms are used interchangeably by our text.
  - The set of all logical addresses used by a program composes the *logical address space*, and the set of all corresponding physical addresses composes the *physical address space.*
- The run time mapping of logical to physical addresses is handled by the *memory-management unit, MMU*.
  - The MMU can take on many forms. One of the simplest is a modification of the base-register scheme described earlier.
  - The base register is now termed a *relocation register*, whose value is added to every memory request at the hardware level.
- Note that user programs never see physical addresses. User programs work entirely in logical address space, and any memory references or manipulations are done using purely logical addresses. Only when the address gets sent to the physical memory chips is the physical memory address generated.



**Figure 8.4 - Dynamic relocation using a relocation register**

## Dynamic Loading:

- Rather than loading an entire program into memory at once, dynamic loading loads up each routine as it is called. The advantage is that unused routines need never be loaded; reducing total memory usage and generating faster program startup times. The downside is the added complexity and overhead of checking to see if a routine is loaded every time it is called and then loading it up if it is not already loaded.

## Dynamic Linking and Shared Libraries

- With *static linking* library modules get fully included in executable modules, wasting both disk space and main memory usage, because every program that included a certain routine from the library would have to have their own copy of that routine linked into their executable code.

- With *dynamic linking*, however, only a stub is linked into the executable module, containing references to the actual library module linked in at run time.

    o This method saves disk space, because the library routines do not need to be fully included in the executable modules, only the stubs.

    o We will also learn that if the code section of the library routines is *reentrant*, ( meaning it does not modify the code while it runs, making it safe to re-enter it ), then main memory can be saved by loading only one copy of dynamically linked routines into memory and sharing the code amongst all processes that are concurrently using it. ( Each process would have their own copy of the *data* section of the routines, but that may be small relative to the code segments. ) Obviously the OS must manage shared routines in memory.

    o An added benefit of *dynamically linked libraries* (*DLLs* also known as *shared libraries* or *shared objects* on UNIX systems) involves easy upgrades and updates. When a program uses a routine from a standard library and the routine changes, then the program must be re-built ( re-linked ) in order to incorporate the changes. However if DLLs are used, then as long as the stub doesn't change, the program can be updated merely by loading new versions of the DLLs onto the system. Version information is maintained in both the program and the DLLs, so that a program can specify a particular version of the DLL if necessary.

    o In practice, the first time a program calls a DLL routine, the stub will recognize the fact and will replace itself with the actual routine from the DLL library. Further calls to the same routine will access the routine directly and not incur the overhead of the stub access. (Following the UML *Proxy Pattern*.)

## Swapping:

- A process must be loaded into memory in order to execute.
- If there is not enough memory available to keep all running processes in memory at the same time, then some processes who are not currently using the CPU may have their memory swapped out to a fast local disk called the *backing store.*

## Standard Swapping

- If compile-time or load-time address binding is used, then processes must be swapped back into the same memory location from which they were swapped out. If execution time binding is used, then the processes can be swapped back into any available location.
- Swapping is a very slow process compared to other operations.

    ▪ For example, if a user process occupied 10 MB and the transfer rate for the backing store were 40 MB per second, then it would take 1/4 second ( 250 milliseconds ) just to do the data transfer. Adding in a latency lag of 8

milliseconds and ignoring head seek time for the moment, and further recognizing that swapping involves moving old data out as well as new data in, the overall transfer time required for this swap is 512 milliseconds, or over half a second. For efficient processor scheduling the CPU time slice should be significantly longer than this lost transfer time.

- To reduce swapping transfer overhead, it is desired to transfer as little information as possible, which requires that the system know how much memory a process *is* using, as opposed to how much it *might* use. Programmers can help with this by freeing up dynamic memory that they are no longer using.

- It is important to swap processes out of memory only when they are idle, or more to the point, only when there are no pending I/O operations. (Otherwise the pending I/O operation could write into the wrong process's memory space.) The solution is to either swap only totally idle processes, or do I/O operations only into and out of OS buffers, which are then transferred to or from process's main memory as a second step.

- Most modern Operating Systems no longer use swapping, because it is too slow and there are faster alternatives available. (e.g. Paging. ) However some UNIX systems will still invoke swapping if the system gets extremely full, and then discontinue swapping when the load reduces again. Windows 3.1 would use a modified version of swapping that was somewhat controlled by the user, swapping process's out if necessary and then only swapping them back in when the user focused on that particular window.



**Figure 8.5 - Swapping of two processes using a disk as a backing store**

**Swapping on Mobile Systems (New Section in 9th Edition)**

- Swapping is typically not supported on mobile platforms, for several reasons:
  - Mobile devices typically use flash memory in place of more spacious hard drives for persistent storage, so there is not as much space available.
  - Flash memory can only be written to a limited number of times before it becomes unreliable.
  - The bandwidth to flash memory is also lower.
- Apple's IOS asks applications to voluntarily free up memory
  - Read-only data, e.g. code, is simply removed, and reloaded later if needed.
  - Modified data, e.g. the stack, is never removed, but . . .

- Apps that fail to free up sufficient memory can be removed by the OS
- Android follows a similar strategy.
  - Prior to terminating a process, Android writes its *application state* to flash memory for quick restarting.

## Contiguous Memory Allocation

- One approach to memory management is to load each process into a contiguous space. The operating system is allocated space first, usually at either low or high memory locations, and then the remaining available memory is allocated to processes as needed. ( The OS is usually loaded low, because that is where the interrupt vectors are located, but on older systems part of the OS was loaded high to make more room in low memory ( within the 640K barrier ) for user processes. )

### Memory Protection (was Memory Mapping and Protection)



- The system shown in Figure 8.6 below allows protection against user programs accessing areas that they should not, allows programs to be relocated to different memory starting addresses as needed, and allows the memory space devoted to the OS to grow or shrink dynamically as needs change.

**Figure 8.6 - Hardware support for relocation and limit registers**

### Memory Allocation

- One method of allocating contiguous memory is to divide all available memory into equal sized partitions, and to assign each process to their own partition. This restricts both the number of simultaneous processes and the maximum size of each process, and is no longer used.
- An alternate approach is to keep a list of unused (free) memory blocks ( holes ), and to find a hole of a suitable size whenever a process needs to be loaded into memory. There are many different strategies for finding the "best" allocation of memory to processes, including the three most commonly discussed:
  1. **First fit** - Search the list of holes until one is found that is big enough to satisfy the request, and assign a portion of that hole to that process. Whatever fraction of the hole not needed by the request is left on the free list as a smaller hole. Subsequent requests may start looking either from the beginning of the list or from the point at which this search ended.

2. **Best fit** - Allocate the *smallest* hole that is big enough to satisfy the request. This saves large holes for other process requests that may need them later, but the resulting unused portions of holes may be too small to be of any use, and will therefore be wasted. Keeping the free list sorted can speed up the process of finding the right hole.

3. **Worst fit** - Allocate the largest hole available, thereby increasing the likelihood that the remaining portion will be usable for satisfying future requests.

- Simulations show that either first or best fit are better than worst fit in terms of both time and storage utilization. First and best fits are about equal in terms of storage utilization, but first fit is faster.

**Fragmentation:**

- All the memory allocation strategies suffer from ***external fragmentation***, though first and best fits experience the problems more so than worst fit. External fragmentation means that the available memory is broken up into lots of little pieces, none of which is big enough to satisfy the next memory requirement, although the sum total could.

- The amount of memory lost to fragmentation may vary with algorithm, usage patterns, and some design decisions such as which end of a hole to allocate and which end to save on the free list.

- Statistical analysis of first fit, for example, shows that for N blocks of allocated memory, another 0.5 N will be lost to fragmentation.

- ***Internal fragmentation*** also occurs, with all memory allocation strategies. This is caused by the fact that memory is allocated in blocks of a fixed size, whereas the actual memory needed will rarely be that exact size. For a random distribution of memory requests, on the average 1/2 block will be wasted per memory request, because on the average the last allocated block will be only half full.

  o Note that the same effect happens with hard drives, and that modern hardware gives us increasingly larger drives and memory at the expense of ever larger block sizes, which translates to more memory lost to internal fragmentation.

  o Some systems use variable size blocks to minimize losses due to internal fragmentation.

- If the programs in memory are relocatable, ( using execution-time address binding ), then the external fragmentation problem can be reduced via ***compaction***, i.e. moving all processes down to one end of physical memory. This only involves updating the relocation register for each process, as all internal work is done using logical addresses.

  - Another solution as we will see in upcoming sections is to allow processes to use non-contiguous blocks of physical memory, with a separate relocation register for each block.

**Segmentation**



subroutine

stack

symbol table

Sqrt

main program

logical address

132

**Basic Method:**

- Most users (programmers ) do not think of their programs as existing in one continuous linear address space.

- Rather they tend to think of their memory in multiple *segments*, each dedicated to a particular use, such as code, data, the stack, the heap, etc.

- Memory *segmentation* supports this view by providing addresses with a segment number (mapped to a segment base address) and an offset from the beginning of that segment.

- For example, a C compiler might generate 5 segments for the user code, library code, global ( static ) variables, the stack, and the heap, as shown in Figure 8.7: **Programmer's view of a program.**

**Segmentation** is a memory-management scheme that supports this programmer view of memory. A logical address space is a collection of segments.

Each segment has a name and a length. The addresses specify both the segment name and the offset within the segment. The programmer therefore specifies each address by two quantities: a segment name and an offset.

For simplicity of implementation, segments are numbered and are referred to by a segment number, rather than by a segment name. Thus, a logical address consists of a *two tuple:*

<segment-number, offset>.

Normally, when a program is compiled, the compiler automatically constructs segments reflecting the input program.

A C compiler might create separate segments for the following:

The code
Global variables
The heap, from which memory is allocated
The stacks used by each thread
The standard C library

Libraries that are linked in during compile time might be assigned separate segments. The loader would take all these segments and assign them segment numbers.

**Segmentation Hardware**

Although the programmer can now refer to objects in the program by a two-dimensional address, the actual physical memory is still, of course, a one-dimensional sequence of bytes. Thus, we must define an implementation to map two-dimensional user-defined addresses into one-dimensional physical addresses. This mapping is effected by a **segment table**. Each entry in the segment table has a **segment base** and a **segment limit**. The segment base contains the starting

physical address where the segment resides in memory, and the segment limit specifies the length of the segment.

The use of a segment table is illustrated in Figure 8.8. A logical address consists of two parts: a segment number, *s,* and an offset into that segment, *d.* The segment number is used as an index to the segment table. The offset *d* of the logical address must be between 0 and the segment limit. If it is not, we trap to the operating system (logical addressing attempt beyond end of segment). When an offset is legal, it is added to the segment base to produce the address in physical memory of the desired byte. The segment table is thus essentially an array of base − limit register pairs.



**Figure 8.8 - Segmentation hardware**



**Figure 8.9 - Example of segmentation**

As an example, consider the situation shown in Figure 8.9. We have five segments numbered from 0 through 4. The segments are stored in physical memory as shown. The segment table has a separate entry for each segment, giving the beginning address of the segment in physical memory (or base) and the length of that segment (or limit). For example, segment 2 is 400 bytes long and begins at location 4300. Thus, a reference to byte 53 of segment 2 is mapped onto location $4300 + 53 = 4353$. A reference to segment 3, byte 852, is mapped to 3200 (the base of segment 3) $+ 852 = 4052$. A reference to byte 1222 of segment 0 would result in a trap to the operating system, as this segment is only 1,000 bytes long.

**Paging:**

- Paging is a memory management scheme that allows processes physical memory to be discontinuous, and which eliminates problems with fragmentation by allocating memory in equal sized blocks known as *pages*.
- Paging eliminates most of the problems of the other methods discussed previously, and is the predominant memory management technique used today.

**Basic Method**

- The basic idea behind paging is to divide physical memory into a number of equal sized blocks called *frames*, and to divide programs logical memory space into blocks of the same size called *pages.*
- Any page (from any process) can be placed into any available frame.
- The *page table* is used to look up what frame a particular page is stored in at the moment. In the following example, for instance, page 2 of the program's logical memory is currently stored in frame 3 of physical memory:

**Figure 8.10 - Paging hardware**



**Figure 8.11 - Paging model of logical and physical memory**



- A logical address consists of two parts: A page number in which the address resides, and an offset from the beginning of that page. (The number of bits in the page number limits how many pages a single process can address. The number of bits in the offset determines the maximum size of each page, and should correspond to the system frame size. )
- The page table maps the page number to a frame number, to yield a physical address which also has two parts: The frame number and the offset within that frame. The number of bits in the frame number determines how many frames the system can address, and the number of bits in the offset determines the size of each frame.

- Page numbers, frame numbers, and frame sizes are determined by the architecture, but are typically powers of two, allowing addresses to be split at a certain number of bits. For example, if the logical address size is $2^m$ and the page size is $2^n$, then the high-order m-n bits of a logical address designate the page number and the remaining n bits represent the offset.

- Note also that the number of bits in the page number and the number of bits in the frame number do not have to be identical. The former determines the address range of the logical address space, and the latter relates to the physical address space.



- (DOS used to use an addressing scheme with 16 bit frame numbers and 16-bit offsets, on hardware that only supported 24-bit hardware addresses. The result was a resolution of starting frame addresses finer than the size of a single frame, and multiple frame-offset combinations that mapped to the same physical hardware address.)

- Consider the following micro example, in which a process has 16 bytes of logical memory, mapped in 4 byte pages into 32 bytes of physical memory. (Presumably some other processes would be consuming the remaining 16 bytes of physical memory. )



**Figure 8.12 - Paging example for a 32-byte memory with 4-byte pages**

- Note that paging is like having a table of relocation registers, one for each page of the logical memory.

- There is no external fragmentation with paging. All blocks of physical memory are used, and there are no gaps in between and no problems with finding the right sized hole for a particular chunk of memory.

- There is, however, internal fragmentation. Memory is allocated in chunks the size of a page, and on the average, the last page will only be half full, wasting on the average half a page of memory per process. ( Possibly more, if processes keep their code and data in separate pages. )

- Larger page sizes waste more memory, but are more efficient in terms of overhead. Modern trends have been to increase page sizes, and some systems even have multiple size pages to try and make the best of both worlds.

- Page table entries ( frame numbers ) are typically 32 bit numbers, allowing access to 2^32 physical page frames. If those frames are 4 KB in size each, that translates to 16 TB of addressable physical memory. ( 32 + 12 = 44 bits of physical address space. )
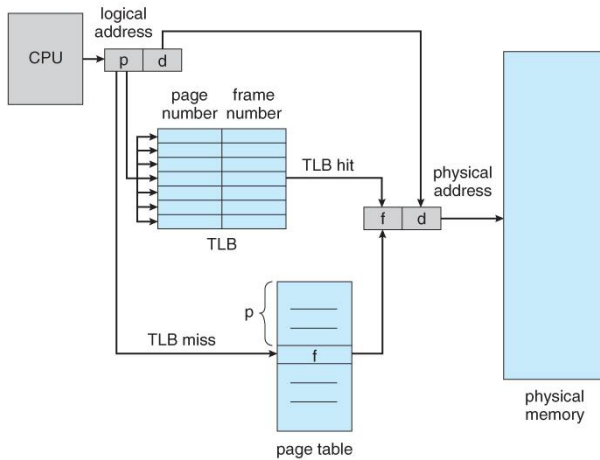


**Figure 8.13 - Free frames (a) before allocation and (b) after allocation**

- When a process requests memory ( e.g. when its code is loaded in from disk ), free frames are allocated from a free-frame list, and inserted into that process's page table.

- Processes are blocked from accessing anyone else's memory because all of their memory requests are mapped through their page table. There is no way for them to generate an address that maps into any other process's memory space.

- The operating system must keep track of each individual process's page table, updating it whenever the process's pages get moved in and out of memory, and applying the correct page table when processing system calls for a particular process. This all increases the overhead involved when swapping processes in and out of the CPU. ( The currently active page table must be updated to reflect the process that is currently running. )

**Hardware Support**

- Page lookups must be done for every memory reference, and whenever a process gets swapped in or out of the CPU, its page table must be swapped in and out too, along with the instruction registers, etc. It is therefore appropriate to provide hardware support for this operation, in order to make it as fast as possible and to make process switches as fast as possible also.

- One option is to use a set of registers for the page table. For example, the DEC PDP-11 uses 16-bit addressing and 8 KB pages, resulting in only 8 pages per process. ( It takes 13 bits to address 8 KB of offset, leaving only 3 bits to define a page number. )

- An alternate option is to store the page table in main memory, and to use a single register ( called the *page-table base register, PTBR* ) to record where in memory the page table is located.

  o Process switching is fast, because only the single register needs to be changed.

  o However memory access just got half as fast, because every memory access now requires *two* memory accesses - One to fetch the frame number from memory and then another one to access the desired memory location.

  o The solution to this problem is to use a very special high-speed memory device called the *translation look-aside buffer, TLB.*

**The benefit of the TLB is that it can search an entire table for a key value in parallel, and if it is found anywhere in the table, then the corresponding lookup value is returned.**

**Figure 8.14 - Paging hardware with TLB**

- The TLB is very expensive, however, and therefore very small. ( Not large enough to hold the entire page table. ) It is therefore used as a cache device.
- Addresses are first checked against the TLB, and if the info is not there (a TLB miss ), then the frame is looked up from main memory and the TLB is updated.
- If the TLB is full, then replacement strategies range from *least-recently used, LRU* to random.
- Some TLBs allow some entries to be *wired down*, which means that they cannot be removed from the TLB. Typically these would be kernel frames.
- Some TLBs store *address-space identifiers, ASIDs*, to keep track of which process "owns" a particular entry in the TLB. This allows entries from multiple processes to be stored simultaneously in the TLB without granting one process access to some other process's memory location. Without this feature the TLB has to be flushed clean with every process switch.
- The percentage of time that the desired information is found in the TLB is termed the *hit ratio*.

**Protection**

- The page table can also help to protect processes from accessing memory that they shouldn't, or their own memory in ways that they shouldn't.
- A bit or bits can be added to the page table to classify a page as read-write, read-only, read-write-execute, or some combination of these sorts of things. Then each memory reference can be checked to ensure it is accessing the memory in the appropriate mode.
- Valid / invalid bits can be added to "mask off" entries in the page table that are not in use by the current process, as shown by example in Figure 8.12 below.
- Note that the valid / invalid bits described above cannot block all illegal memory accesses, due to the internal fragmentation. ( Areas of memory in the last page that are not entirely filled by the process, and may contain data left over by whoever used that frame last. )
- Many processes do not use all of the page table available to them, particularly in modern systems with very large potential page tables. Rather than waste memory by creating a full-size page table for every process, some systems use a *page-table length register, PTLR*, to specify the length of the page table.

<u>**Structure of the Page Table**</u>

**1. Hierarchical Paging**

- Most modern computer systems support logical address spaces of 2^32 to 2^64.
- With a 2^32 address space and 4K ( 2^12 ) page sizes, this leave 2^20 entries in the page table. At 4 bytes per entry, this amount to a 4 MB page table, which is too large to reasonably keep in contiguous memory. (And to swap in and out of memory with each process switch. ) Note that with 4K pages, this would take 1024 pages just to hold the page table!
- One option is to use a two-tier paging system, i.e. to page the page table.
- For example, the 20 bits described above could be broken down into two 10-bit page numbers. The first identifies an entry in the outer page table, which identifies where in memory to find one page of an inner page table. The second 10 bits finds a specific entry in that inner page table, which in turn identifies a particular frame in physical memory. (The remaining 12 bits of the 32 bit logical address are the offset within the 4K frame.)

**Figure 8.17 A two-level page-table**  **scheme**



**Figure 8.18 - Address translation for a two-level 32-bit paging architecture**

- VAX Architecture divides 32-bit addresses into 4 equal sized sections, and each page is 512 bytes, yielding an address form of:

With a 64-bit logical address space and 4K pages, there are 52 bits worth of page numbers, which is still too many even for two-level paging. One could increase the paging level, but with 10-bit page tables it would take 7 levels of indirection, which would be prohibitively slow memory access. So some other approach must be used.



| section | page | offset |
|---------|------|--------|
| s | p | d |
| 2 | 21 | 9 |

139

| outer page | inner page | offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 42 | 10 | 12 |

64-bits Two-tiered leaves 42 bits in outer table

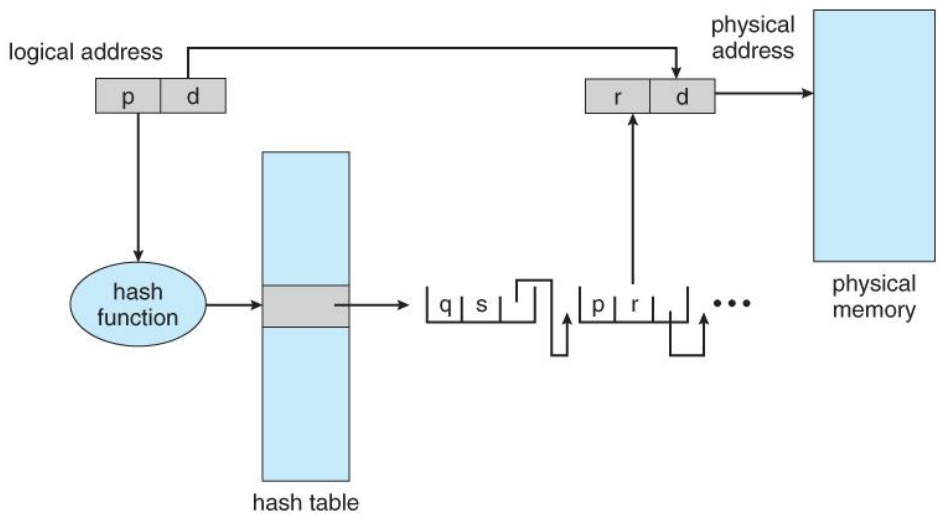| 2nd outer page | outer page | inner page | offset |
|:---:|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $p_3$ | $d$ |
| 32 | 10 | 10 | 12 |

Going to a fourth level still leaves 32 bits in the outer table.

## 2. Hashed Page Tables

- One common data structure for accessing data that is sparsely distributed over a broad range of possible values is with *hash tables*. Figure 8.16 below illustrates a *hashed page table* using chain-and-bucket hashing:

**Figure 8.19 - Hashed page table**



## 3. Inverted Page Tables

- Another approach is to use an *inverted page table*. Instead of a table listing all of the pages for a particular process, an inverted page table lists all of the pages currently loaded in memory, for all processes. ( I.e. there is one entry per *frame* instead of one entry per *page*. )

- Access to an inverted page table can be slow, as it may be necessary to search the entire table in order to find the desired page ( or to discover that it is not there. ) Hashing the table can help speedup the search process.

- Inverted page tables prohibit the normal method of implementing shared memory, which is to map multiple logical pages to a common physical frame. ( Because each frame is now mapped to one and only one process. )



**Figure 8.20 - Inverted page table**

# Virtual Memory

## Background

- Preceding sections talked about how to avoid memory fragmentation by breaking process memory requirements down into smaller bites ( pages ), and storing the pages non-contiguously in memory. However the entire process still had to be stored in memory somewhere.

- In practice, most real processes do not need all their pages, or at least not all at once, for several reasons:

  1. Error handling code is not needed unless that specific error occurs, some of which are quite rare.

  2. Arrays are often over-sized for worst-case scenarios, and only a small fraction of the arrays are actually used in practice.

  3. Certain features of certain programs are rarely used, such as the routine to balance the federal budget. :-

- The ability to load only the portions of processes that were actually needed ( and only *when* they were needed ) has several benefits:

  o Programs could be written for a much larger address space ( virtual memory space ) than physically exists on the computer.

  o Because each process is only using a fraction of their total address space, there is more memory left for other programs, improving CPU utilization and system throughput.

  o Less I/O is needed for swapping processes in and out of RAM, speeding things up.

**Figure 9.1 - Diagram showing virtual memory that is larger than physical memory**



- Figure 9.2 shows *virtual address space*, which is the programmers logical view of process memory storage. The actual physical layout is controlled by the process's page table.

- Note that the address space shown in Figure 9.2 is *sparse* - A great hole in the middle of the address space is never used, unless the stack and/or the heap grow to fill the hole.
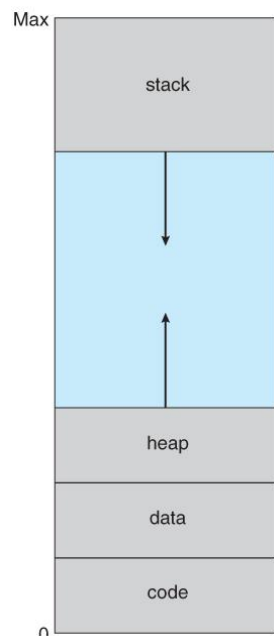


**Figure 9.2 - Virtual address space**

- Virtual memory also allows the sharing of files and memory by multiple processes, with several benefits:

142

o System libraries can be shared by mapping them into the virtual address space of more than one process.



o Processes can also share virtual memory by mapping the same block of memory to more than one process.

o Process pages can be shared during a fork( ) system call, eliminating the need to copy all of the pages of the original ( parent ) process.

**Figure 9.3 - Shared library using virtual memory**

## Demand Paging:



- The basic idea behind **demand paging** is that when a process is swapped in, its pages are not swapped in all at once. Rather they are swapped in only when the process needs them. (on demand. ) This is termed a **lazy swapper**, although a **pager** is a more accurate term.

**Figure 9.4 - Transfer of a paged memory to contiguous disk space**

**Basic Concepts:**

- The basic idea behind paging is that when a process is swapped in, the pager only loads into memory those pages that it expects the process to need (right away.)

- Pages that are not loaded into memory are marked as invalid in the page table, using the invalid bit. (The rest of the page table entry may either be blank or contain information about where to find the swapped-out page on the hard drive. )

- If the process only ever accesses pages that are loaded in memory (memory *resident*



3

pages), then the process runs exactly as if all the pages were loaded in to memory.

**Figure 9.5 - Page table when some pages are not in main memory.**

- On the other hand, if a page is needed that was not originally loaded up, then a *page fault trap* is generated, which must be handled in a series of steps:

  1. The memory address requested is first checked, to make sure it was a valid memory request.
  2. If the reference was invalid, the process is terminated. Otherwise, the page must be paged in.
  3. A free frame is located, possibly from a free-frame list.
  4. A disk operation is scheduled to bring in the necessary page from disk. ( This will usually block the process on an I/O wait, allowing some other process to use the CPU in the meantime. )
  5. When the I/O operation is complete, the process's page table is updated with the new frame number, and the invalid bit is changed to indicate that this is now a valid page reference.
  6. The instruction that caused the page fault must now be restarted from the beginning, ( as soon as this process gets another turn on the CPU. )



**Figure 9.6 - Steps in handling a page fault**

- In an extreme case, NO pages are swapped in for a process until they are requested by page faults. This is known as *pure demand paging.*

- In theory each instruction could generate multiple page faults. In practice this is very rare, due to *locality of reference*.

- The hardware necessary to support virtual memory is the same as for paging and swapping: A page table and secondary memory. (*Swap space,* whose allocation is discussed in chapter 12.)

- A crucial part of the process is that the instruction must be restarted from scratch once the desired page has been made available in memory. For most simple instructions this is not a major difficulty. However there are some architectures that allow a single instruction to modify a fairly large block of data, (which may span a page boundary), and if some of the data gets modified before the page fault occurs, this could cause problems. One solution is to access both ends of the

block before executing the instruction, guaranteeing that the necessary pages get paged in before the instruction begins.

**Performance of Demand Paging**

- Obviously there is some slowdown and performance hit whenever a page fault occurs and the system has to go get it from memory, but just how big a hit is it exactly?
- There are many steps that occur when servicing a page fault ( see book for full details ), and some of the steps are optional or variable. But just for the sake of discussion, suppose that a normal memory access requires 200 nanoseconds, and that servicing a page fault takes 8 milliseconds. (8,000,000 nanoseconds, or 40,000 times a normal memory access. ) With a *page fault rate* of p, ( on a scale from 0 to 1 ), the effective access time is now:

$$( 1 - p ) * ( 200 ) + p * 8000000 = 200 + 7{,}999{,}800 * p$$

Which *clearly* depends heavily on p! Even if only one access in 1000 causes a page fault, the effective access time drops from 200 nanoseconds to 8.2 microseconds, a slowdown of a factor of 40 times. In order to keep the slowdown less than 10%, the page fault rate must be less than 0.0000025, or one in 399,990 accesses.

- A subtlety is that swap space is faster to access than the regular file system, because it does not have to go through the whole directory structure. For this reason some systems will transfer an entire process from the file system to swap space before starting up the process, so that future paging all occurs from the (relatively) faster swap space.
- Some systems use demand paging directly from the file system for binary code ( which never changes and hence does not have to be stored on a page operation ), and to reserve the swap space for data segments that must be stored. This approach is used by both Solaris and BSD Unix.

**Copy-on-Write:**



- The idea behind a copy-on-write fork is that the pages for a parent process do not have to be actually copied for the child until one or the other of the processes changes the page. They can be simply shared between the two processes in the meantime, with a bit set that the page needs to be copied if it ever gets written to. This is a reasonable approach, since the child process usually issues an exec( ) system call immediately after the fork.

**Figure 9.7 - Before process 1 modifies page C.**

**Figure 9.8 - After process 1 modifies page C.**

- Obviously only pages that can be modified even need to be labeled as copy-on-write. Code segments can simply be shared.

- Pages used to satisfy copy-on-write duplications are typically allocated using *zero-fill-on-demand*, meaning that their previous contents are zeroed out before the copy proceeds.

- Some systems provide an alternative to the fork( ) system call called a *virtual memory fork, vfork( )*. In this case the parent is suspended, and the child uses the parent's memory pages. This is very fast for process creation, but requires that the child not modify any of the shared memory pages before performing the exec( ) system call. (In essence this addresses the question of which process executes first after a call to fork, the parent or the child. With vfork, the parent is suspended, allowing the child to execute first until it calls exec( ), sharing pages with the parent in the meantime.

## Page Replacement:

- In order to make the most use of virtual memory, we load several processes into memory at the same time. Since we only load the pages that are actually needed by each process at any given time, there is room to load many more processes than if we had to load in the entire process.

- However memory is also needed for other purposes (such as I/O buffering), and what happens if some process suddenly decides it needs more pages and there aren't any free frames available? There are several possible solutions to consider:

  1. Adjust the memory used by I/O buffering, etc., to free up some frames for user processes. The decision of how to allocate memory for I/O versus user processes is a complex one, yielding different policies on different systems. (Some allocate a fixed amount for I/O, and others let the I/O system contend for memory along with everything else.)

  2. Put the process requesting more pages into a wait queue until some free frames become available.

  3. Swap some process out of memory completely, freeing up its page frames.

  4. Find some page in memory that isn't being used right now, and swap that page only out to disk, freeing up a frame that can be allocated to the process requesting it. This is known as *page replacement*, and is the most common solution. There are

many different algorithms for page replacement, which is the subject of the remainder of this section.

**Figure 9.9 - Ned for page replacement.**

**Basic Page Replacement**

- The previously discussed page-fault processing assumed that there would be free frames available on the free-frame list. Now the page-fault handling must be modified to free up a frame if necessary, as follows:

    1. Find the location of the desired page on the disk, either in swap space or in the file system.
    2. Find a free frame:
        a. If there is a free frame, use it.
        b. If there is no free frame, use a page-replacement algorithm to select an existing frame to be replaced, known as the *victim frame*.
        c. Write the victim frame to disk. Change all related page tables to indicate that this page is no longer in memory.
    3. Read in the desired page and store it in the frame. Adjust all related page and frame tables to indicate the change.
    4. Restart the process that was waiting for this page.



**Figure 9.10 - Page replacement.**

- Note that step 3c adds an extra disk write to the page-fault handling, effectively doubling the time required to process a page fault. This can be alleviated somewhat by assigning a *modify bit,* or *dirty bit* to each page, indicating whether or not it has been changed since it was last loaded in from disk. If the dirty bit has not been set, then the page is unchanged, and does not need to be written out to disk. Otherwise the page write is required. It should come as no surprise that many page replacement strategies specifically look for pages that do not have their dirty bit set, and preferentially select clean pages as victim pages. It should also be obvious that unmodifiable code pages never get their dirty bits set.

- There are two major requirements to implement a successful demand paging system. We must develop a *frame-allocation algorithm* and a *page-replacement algorithm.* The former centers around how many frames are allocated to each process (and to other needs), and the latter deals with how to select a page for replacement when there are no free frames available.

- The overall goal in selecting and tuning these algorithms is to generate the fewest number of overall page faults. Because disk access is so slow relative to memory access, even slight improvements to these algorithms can yield large improvements in overall system performance.

- Algorithms are evaluated using a given string of memory accesses known as a *reference string,* which can be generated in one of ( at least ) three common ways:

  1. Randomly generated, either evenly distributed or with some distribution curve based on observed system behavior. This is the fastest and easiest approach, but may not reflect real performance well, as it ignores locality of reference.

  2. Specifically designed sequences. These are useful for illustrating the properties of comparative algorithms in published papers and textbooks, ( and also for homework and exam problems. :-) )

  3. Recorded memory references from a live system. This may be the best approach, but the amount of data collected can be enormous, on the order of a million addresses per second. The volume of collected data can be reduced by making two important observations:

  1. Only the page number that was accessed is relevant. The offset within that page does not affect paging operations.

  2. Successive accesses within the same page can be treated as a single page request, because all requests after the first are guaranteed to be page hits. (Since there are no intervening requests for other pages that could remove this page from the page table.)

  - So for example, if pages were of size 100 bytes, then the sequence of address requests ( 0100, 0432, 0101, 0612, 0634, 0688, 0132, 0038, 0420 ) would reduce to page requests ( 1, 4, 1, 6, 1, 0, 4 )

As the number of available frames increases, the number of page faults should decrease, as shown in Figure 9.11:



**Figure 9.11 - Graph of page faults versus number of frames.**

**2. FIFO Page Replacement**

- A simple and obvious page replacement strategy is *FIFO*, i.e. first-in-first-out.
- As new pages are brought in, they are added to the tail of a queue, and the page at the head of the queue is the next victim. In the following example, 20 page requests result in 15 page faults:

reference string

7  0  1  2  0  3  0  4  2  3  0  3  2  1  2  0  1  7  0  1



page frames

**Figure 9.12 - FIFO page-replacement algorithm.**

- Although FIFO is simple and easy, it is not always optimal, or even efficient.
- An interesting effect that can occur with FIFO is ***Belady's anomaly***, in which increasing the number of frames available can actually ***increase*** the number of page faults that occur! Consider, for example, the following chart based on the page sequence ( 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5 ) and a varying number of available frames. Obviously the maximum number of faults is 12 ( every request generates a fault ), and the minimum number is 5 ( each page loaded only once ), but in between there are some interesting results:



**Figure 9.13 - Page-fault curve for FIFO replacement on a reference string.**

## 9.4.3 Optimal Page Replacement

- The discovery of Belady's anomaly lead to the search for an ***optimal page-replacement algorithm***, which is simply that which yields the lowest of all possible page-faults, and which does not suffer from Belady's anomaly.
- Such an algorithm does exist, and is called ***OPT or MIN.*** This algorithm is simply "Replace the page that will not be used for the longest time in the future."
- For example, Figure 9.14 shows that by applying OPT to the same reference string used for the FIFO example, the minimum number of possible page faults is 9. Since 6 of the page-faults are unavoidable ( the first reference to each new page ), FIFO can be shown to require 3 times as many ( extra ) page faults as the optimal algorithm. ( Note: The book claims that only the first three page faults are required by all algorithms, indicating that FIFO is only twice as bad as OPT. )
- Unfortunately OPT cannot be implemented in practice, because it requires foretelling the future, but it makes a nice benchmark for the comparison and evaluation of real proposed new algorithms.

150

- In practice most page-replacement algorithms try to approximate OPT by predicting ( estimating ) in one fashion or another what page will not be used for the longest period of time. The basis of FIFO is the prediction that the page that was brought in the longest time ago is the one that will not be needed again for the longest future time, but as we shall see, there are many other prediction methods, all striving to match the performance of OPT.
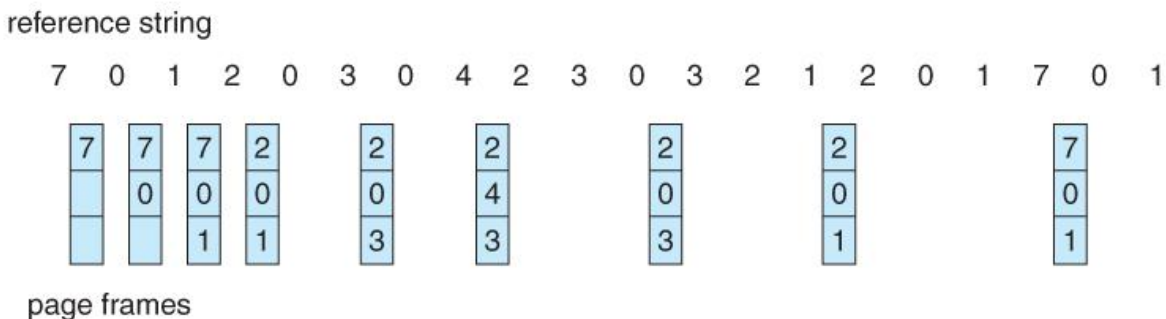


**Figure 9.14 - Optimal page-replacement algorithm**

### 9.4.4 LRU Page Replacement

- The prediction behind *LRU,* the *Least Recently Used,* algorithm is that the page that has not been used in the longest time is the one that will not be used again in the near future. ( Note the distinction between FIFO and LRU: The former looks at the oldest *load* time, and the latter looks at the oldest *use* time. )
- Some view LRU as analogous to OPT, except looking backwards in time instead of forwards. ( OPT has the interesting property that for any reference string S and its reverse R, OPT will generate the same number of page faults for S and for R. It turns out that LRU has this same property. )
- Figure 9.15 illustrates LRU for our sample string, yielding 12 page faults, ( as compared to 15 for FIFO and 9 for OPT. )



**Figure 9.15 - LRU page-replacement algorithm.**

- LRU is considered a good replacement policy, and is often used. The problem is how exactly to implement it. There are two simple approaches commonly used:
  1. **Counters.** Every memory access increments a counter, and the current value of this counter is stored in the page table entry for that page. Then finding the LRU page involves

simple searching the table for the page with the smallest counter value. Note that overflowing of the counter must be considered.

2. **Stack.** Another approach is to use a stack, and whenever a page is accessed, pull that page from the middle of the stack and place it on the top. The LRU page will always be at the bottom of the stack. Because this requires removing objects from the middle of the stack, a doubly linked list is the recommended data structure.

- Note that both implementations of LRU require hardware support, either for incrementing the counter or for managing the stack, as these operations must be performed for *every* memory access.

- Neither LRU or OPT exhibit Belady's anomaly. Both belong to a class of page-replacement algorithms called *stack algorithms,* which can never exhibit Belady's anomaly. A stack algorithm is one in which the pages kept in memory for a frame set of size N will always be a subset of the pages kept for a frame size of N + 1. In the case of LRU, ( and particularly the stack implementation thereof ), the top N pages of the stack will be the same for all frame set sizes of N or anything larger.



**Figure 9.16 - Use of a stack to record the most recent page references.**

### 9.4.5 LRU-Approximation Page Replacement

- Unfortunately full implementation of LRU requires hardware support, and few systems provide the full hardware support necessary.

- However many systems offer some degree of HW support, enough to approximate LRU fairly well. ( In the absence of ANY hardware support, FIFO might be the best available choice. )

- In particular, many systems provide a *reference bit* for every entry in a page table, which is set anytime that page is accessed. Initially all bits are set to zero, and they can also all be cleared at any time. One bit of precision is enough to distinguish pages that have been accessed since the last clear from those that have not, but does not provide any finer grain of detail.

### 9.4.5.1 Additional-Reference-Bits Algorithm

- Finer grain is possible by storing the most recent 8 reference bits for each page in an 8-bit byte in the page table entry, which is interpreted as an unsigned int.
  - At periodic intervals (clock interrupts ), the OS takes over, and right-shifts each of the reference bytes by one bit.
  - The high-order (leftmost ) bit is then filled in with the current value of the reference bit, and the reference bits are cleared.
  - At any given time, the page with the smallest value for the reference byte is the LRU page.
- Obviously the specific number of bits used and the frequency with which the reference byte is updated are adjustable, and are tuned to give the fastest performance on a given hardware platform.

**9.4.5.2 Second-Chance Algorithm**

- The *second chance algorithm* is essentially a FIFO, except the reference bit is used to give pages a second chance at staying in the page table.
  - When a page must be replaced, the page table is scanned in a FIFO ( circular queue ) manner.
  - If a page is found with its reference bit not set, then that page is selected as the next victim.
  - If, however, the next page in the FIFO **does** have its reference bit set, then it is given a second chance:
    - The reference bit is cleared, and the FIFO search continues.
    - If some other page is found that did not have its reference bit set, then that page will be selected as the victim, and this page ( the one being given the second chance ) will be allowed to stay in the page table.
    - If , however, there are no other pages that do not have their reference bit set, then this page will be selected as the victim when the FIFO search circles back around to this page on the second pass.
- If all reference bits in the table are set, then second chance degrades to FIFO, but also requires a complete search of the table for every page-replacement.
- As long as there are some pages whose reference bits are not set, then any page referenced frequently enough gets to stay in the page table indefinitely.



circular queue of pages          circular queue of pages

(a)                              (b)

- This algorithm is also known as the ***clock*** algorithm, from the hands of the clock moving around the circular queue.

**Figure 9.17 - Second-chance ( clock ) page-replacement algorithm.**

**Allocation of Frames:**

We said earlier that there were two important tasks in virtual memory management: a page-replacement strategy and a frame-allocation strategy. This section covers the second part of that pair.

**Minimum Number of Frames**

- The absolute minimum number of frames that a process must be allocated is dependent on system architecture, and corresponds to the worst-case scenario of the number of pages that could be touched by a single (machine) instruction.
- If an instruction (and its operands) spans a page boundary, then multiple pages could be needed just for the instruction fetch.
- Memory references in an instruction touch more pages, and if those memory locations can span page boundaries, then multiple pages could be needed for operand access also.
- The worst case involves indirect addressing, particularly where multiple levels of indirect addressing are allowed. Left unchecked, a pointer to a pointer to a pointer to a pointer to a . . . could theoretically touch every page in the virtual address space in a single machine instruction, requiring every virtual page be loaded into physical memory simultaneously. For this reason architectures place a limit ( say 16 ) on the number of levels of indirection allowed in an instruction, which is enforced with a counter initialized to the limit and decremented with every level of indirection in an instruction - If the counter reaches zero, then an "excessive indirection" trap occurs. This example would still require a minimum frame allocation of 17 per process.

**Allocation Algorithms**

- **Equal Allocation -** If there are m frames available and n processes to share them, each process gets m / n frames, and the leftovers are kept in a free-frame buffer pool.
- **Proportional Allocation -** Allocate the frames proportionally to the size of the process, relative to the total size of all processes. So if the size of process i is $S_i$, and S is the sum of all $S_i$, then the allocation for process $P_i$ is $a_i = m * S_i / S$.
- Variations on proportional allocation could consider priority of process rather than just their size.
- Obviously all allocations fluctuate over time as the number of available free frames, m, fluctuates, and all are also subject to the constraints of minimum allocation. ( If the minimum allocations cannot be met, then processes must either be swapped out or not allowed to start until more free frames become available. )

**Global versus Local Allocation**

- One big question is whether frame allocation ( page replacement ) occurs on a local or global level.

- With local replacement, the number of pages allocated to a process is fixed, and page replacement occurs only amongst the pages allocated to this process.
- With global replacement, any page may be a potential victim, whether it currently belongs to the process seeking a free frame or not.
- Local page replacement allows processes to better control their own page fault rates, and leads to more consistent performance of a given process over different system load levels.
- Global page replacement is overall more efficient, and is the more commonly used approach.

**9.5.4 Non-Uniform Memory Access**

- The above arguments all assume that all memory is equivalent, or at least has equivalent access times.
- This may not be the case in multiple-processor systems, especially where each CPU is physically located on a separate circuit board which also holds some portion of the overall system memory.
- In these latter systems, CPUs can access memory that is physically located on the same board much faster than the memory on the other boards.
- The basic solution is akin to processor affinity - At the same time that we try to schedule processes on the same CPU to minimize cache misses, we also try to allocate memory for those processes on the same boards, to minimize access times.
- The presence of threads complicates the picture, especially when the threads get loaded onto different processors.
- Solaris uses an *lgroup* as a solution, in a hierarchical fashion based on relative latency. For example, all processors and RAM on a single board would probably be in the same lgroup. Memory assignments are made within the same lgroup if possible, or to the next nearest lgroup otherwise. ( Where "nearest" is defined as having the lowest access time. )

# Thrashing

- If a process cannot maintain its minimum required number of frames, then it must be swapped out, freeing up frames for other processes. This is an intermediate level of CPU scheduling.
- But what about a process that can keep its minimum, but cannot keep all of the frames that it is currently using on a regular basis? In this case it is forced to page out pages that it will need again in the very near future, leading to large numbers of page faults.
- A process that is spending more time paging than executing is said to be *thrashing.*

**Cause of Thrashing**

- Early process scheduling schemes would control the level of multiprogramming allowed based on CPU utilization, adding in more processes when CPU utilization was low.

156

- The problem is that when memory filled up and processes started spending lots of time waiting for their pages to page in, then CPU utilization would lower, causing 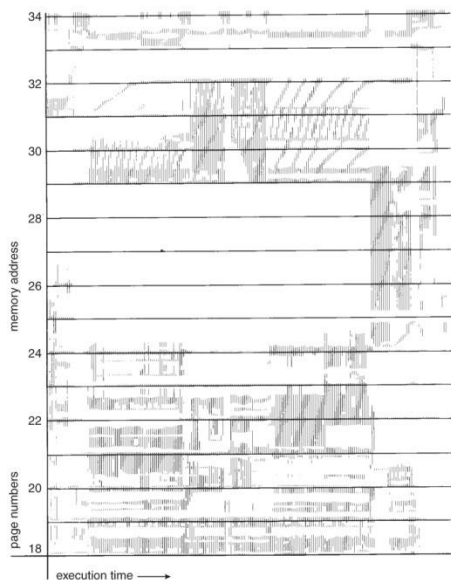the schedule to add in even more processes and exacerbating the problem! Eventually the system would essentially grind to a halt.



- Local page replacement policies can prevent one thrashing process from taking pages away from other processes, but it still tends to clog up the I/O queue, thereby slowing down any other process that needs to do even a little bit of paging ( or any other I/O for that matter. )

**Figure 9.18 - Thrashing**



- To prevent thrashing we must provide processes with as many frames as they really need "right now", but how do we know what that is?

- The *locality model* notes that processes typically access memory references in a given *locality,* making lots of references to the same general area of memory before moving periodically to a new locality, as shown in Figure 9.19 below. If we could just keep as many frames as are involved in the current locality, then page faulting would occur primarily on switches from one locality to another. ( E.g. when one function exits and another is called.)

**Figure 9.19 - Locality in a memory-reference pattern.**

**Working-Set Model**

- The *working set model* is based on the concept of locality, and defines a *working set window,* of length *delta.* Whatever pages are included in the most recent delta page references are said to be in the processes working set window, and comprise its current working set, as illustrated in Figure 9.20:

page reference table

$\ldots$ 2 6 1 5 7 7 7 7 5 1 6 2 3 4 1 2 3 4 4 4 3 4 3 4 4 4 1 3 2 3 4 4 4 3 4 4 4 $\ldots$



$WS(t_1) = \{1,2,5,6,7\}$          $WS(t_2) = \{3,4\}$

**Figure 9.20 - Working-set model.**

- The selection of delta is critical to the success of the working set model - If it is too small then it does not encompass all of the pages of the current locality, and if it is too large, then it encompasses pages that are no longer being frequently accessed.

- The total demand, D, is the sum of the sizes of the working sets for all processes. If D exceeds the total number of available frames, then at least one process is thrashing, because there are not enough frames available to satisfy its minimum working set. If D is significantly less than the currently available frames, then additional processes can be launched.

- The hard part of the working-set model is keeping track of what pages are in the current working set, since every reference adds one to the set and removes one older page. An approximation can be made using reference bits and a timer that goes off after a set interval of memory references:
    - For example, suppose that we set the timer to go off after every 5000 references ( by any process ), and we can store two additional historical reference bits in addition to the current reference bit.
    - Every time the timer goes off, the current reference bit is copied to one of the two historical bits, and then cleared.
    - If any of the three bits is set, then that page was referenced within the last 15,000 references, and is considered to be in that processes reference set.
    - Finer resolution can be achieved with more historical bits and a more frequent timer, at the expense of greater overhead.

**Page-Fault Frequency**

- A more direct approach is to recognize that what we really want to control is the page-fault rate, and to allocate frames based on this directly measurable value. If the page-fault rate exceeds a certain upper bound then that process needs more frames, and if it is below a given lower bound, then it can afford to give up some of its frames to other processes.
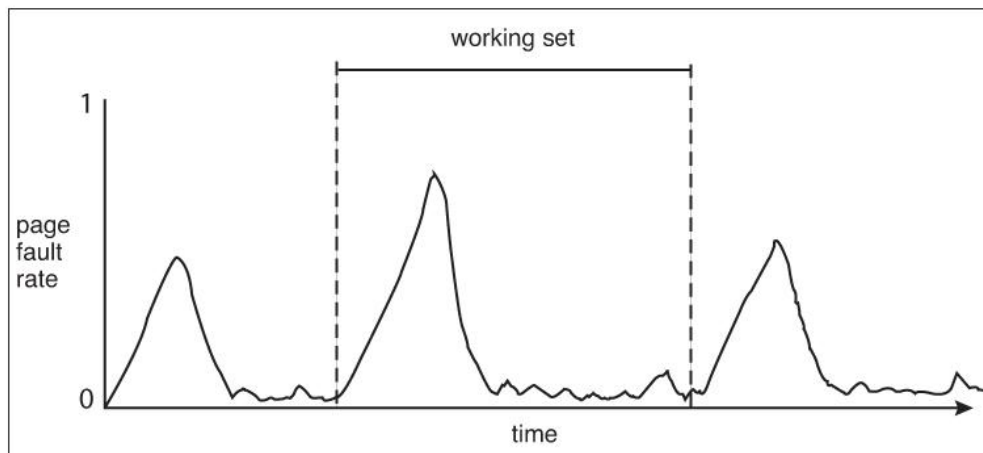
- ( I suppose a page-replacement strategy could be devised that would select victim frames based on the process with the lowest current page-fault frequency. )

**Figure 9.21 - Page-fault frequency.**

- Note that there is a direct relationship between the page-fault rate and the working-set, as a process moves from one locality to another:
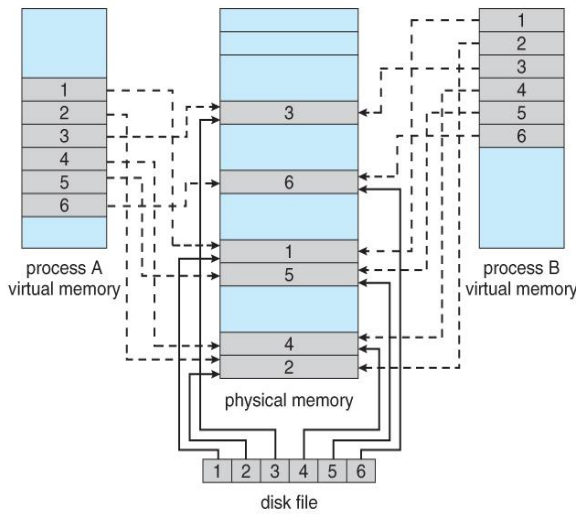
**Unnumbered side bar in Ninth Edition**



# Memory-Mapped Files

- Rather than accessing data files directly via the file system with every file access, data files can be paged into memory the same as process files, resulting in much faster accesses ( except of course when page-faults occur. ) This is known as *memory-mapping* a file.

### 9.7.1 Basic Mechanism

- Basically a file is mapped to an address range within a process's virtual address space, and then paged in as needed using the ordinary demand paging system.
- Note that file writes are made to the memory page frames, and are not immediately written out to disk. ( This is the purpose of the "flush( )" system call, which may also be needed for stdout in some cases. See the timekiller program for an example of this. )

- This is also why it is important to "close( )" a file when one is done writing to it - So that the data can be safely flushed out to disk and so that the memory frames can be freed up for other purposes.

- Some systems provide special system calls to memory map files and use direct disk access otherwise. Other systems map the file to process address space if the special system calls are used and map the file to kernel address space otherwise, but do memory mapping in either case.

- File sharing is made possible by mapping the same file to the address space of more than one process, as shown in Figure 9.23 below. Copy-on-write is supported, and mutual exclusion techniques ( chapter 6 ) may be needed to avoid synchronization problems.

**Figure 9.22** Memory-mapped files.

- Shared memory can be implemented via shared memory-mapped files (Windows), or it can be implemented through a separate process (Linux, UNIX. )

### 9.7.2 Shared Memory in the Win32 API

- Windows implements shared memory using shared memory-mapped files, involving three basic steps:

    1. Create a file, producing a HANDLE to the new file.
    2. Name the file as a shared object, producing a HANDLE to the shared object.
    3. Map the shared object to virtual memory address space, returning its base address as a void pointer (LPVOID).

- This is illustrated in Figures 9.24 to 9.26 (annotated.)

**Figure 9.23 - Shared memory in Windows using memory-mapped I/O.**

## Memory-Mapped I/O

**Consumer accesses named shared object, maps to memory, reads, and closes**

```
#include <windows.h>
#include <stdio.h>

int main(int argc, char *argv[])
{
  HANDLE hMapFile;
  LPVOID lpMapAddress;
                              Access shared memory
  hMapFile = OpenFileMapping(FILE_MAP_ALL_ACCESS, // R/W access
    FALSE, // no inheritance
    TEXT("SharedObject")); // name of mapped file object

  lpMapAddress = MapViewOfFile(hMapFile, // mapped object handle
    FILE_MAP_ALL_ACCESS, // read/write access
    0, // mapped view of entire file
    0,            Map to virtual memory space
    0);

  // read from shared memory
  printf("Read message %s", lpMapAddress);

  UnmapViewOfFile(lpMapAddress);     Read and close
  CloseHandle(hMapFile);
}
```

Figure 9.26  Consumer reading from shared memory using the Win32 API.

- All access to devices is done by writing into ( or reading from ) the device's registers. Normally this is done via special I/O instructions.
- For certain devices it makes sense to simply map the device's registers to addresses in the process's virtual address space, making device I/O as fast and simple as any other memory access. Video controller cards are a classic example of this.
- erial and parallel devices can also use memory mapped I/O, mapping the device registers to specific memory addresses known as *I/O Ports*, e.g. 0xF8. Transferring a series of bytes must be done one at a time, moving only as fast as the I/O device is prepared to process the data, through one of two mechanisms:
  - *Programmed I/O ( PIO )*, also known as *polling.* The CPU periodically checks the control bit on the device, to see if it is ready to handle another byte of data.
  - *Interrupt Driven.* The device generates an interrupt when it either has another byte of data to deliver or is ready to receive another byte.

# Allocating Kernel Memory

- Previous discussions have centered on process memory, which can be conveniently broken up into page-sized chunks, and the only fragmentation that occurs is the average half-page lost to internal fragmentation for each process (segment.)
- There is also additional memory allocated to the kernel, however, which cannot be so easily paged. Some of it is used for I/O buffering and direct access by devices, example, and must therefore be contiguous and not affected by paging. Other memory is used for internal kernel data structures of various sizes, and since kernel memory is often locked (restricted from being ever swapped out), management of this resource must be done carefully to avoid internal fragmentation or other waste. (I.e. you would like the kernel to consume as little memory as possible, leaving as much as possible for user processes. ) Accordingly there are several classic algorithms in place for allocating kernel memory structures.

### 9.8.1 Buddy System

- The Buddy *System* allocates memory using a *power of two allocator*.
- Under this scheme, memory is always allocated as a power of 2 ( 4K, 8K, 16K, etc ), rounding up to the next nearest power of two if necessary.
- If a block of the correct size is not currently available, then one is formed by splitting the next larger block in two, forming two matched buddies. ( And if that larger size is not available, then the next largest available size is split, and so on. )
- One nice feature of the buddy system is that if the address of a block is exclusively ORed with the size of the block, the resulting address is the address of the buddy of the same size, which allows for fast and easy *coalescing* of free blocks back into larger blocks.
  - Free lists are maintained for every size block.
  - If the necessary block size is not available upon request, a free block from the next largest size is split into two buddies of the desired size. ( Recursively splitting larger size blocks if necessary. )

- When a block is freed, its buddy's address is calculated, and the free list for that size block is checked to see if the buddy is also free. If it is, then the two buddies are coalesced into one larger free block, and the process is repeated with successively larger free lists.
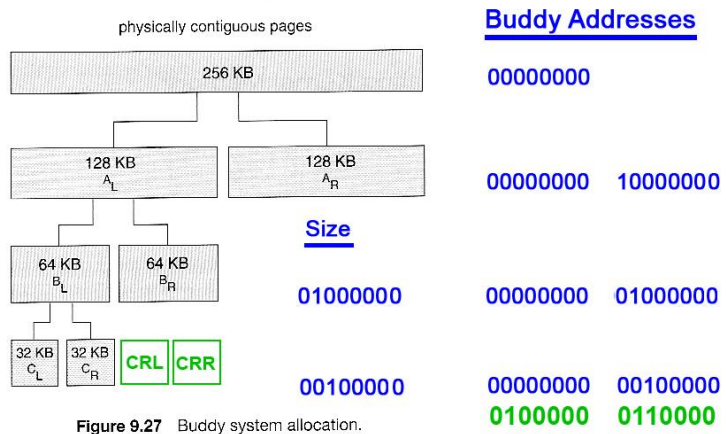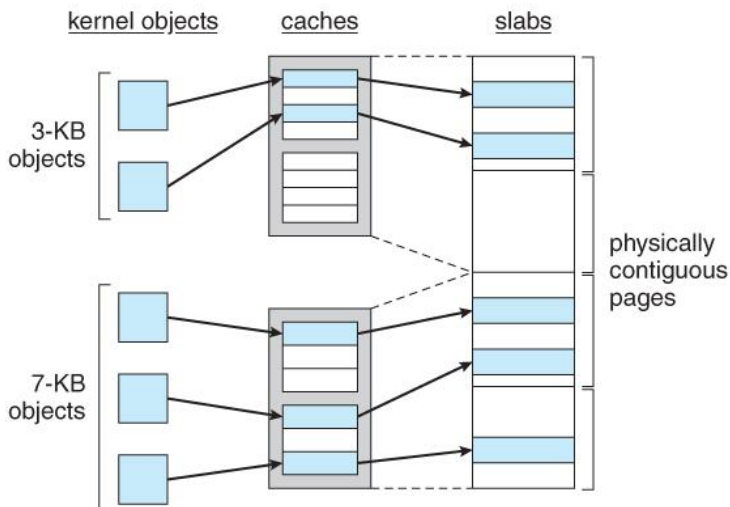- See the ( annotated ) Figure 9.27 below for an example.



**Figure 9.27**   Buddy system allocation.

**Figure 9.26**

### 9.8.2 Slab Allocation

- *Slab Allocation* allocates memory to the kernel in chunks called *slabs*, consisting of one or more contiguous pages. The kernel then creates separate caches for each type of data structure it might need from one or more slabs. Initially the caches are marked empty, and are marked full as they are used.
- New requests for space in the cache is first granted from empty or partially empty slabs, and if all slabs are full, then additional slabs are allocated.
- ( This essentially amounts to allocating space for arrays of structures, in large chunks suitable to the size of the structure being stored. For example if a particular structure were 512 bytes long, space for them would be allocated in groups of 8 using 4K pages. If the structure were 3K, then space for 4 of them could be allocated at one time in a slab of 12K using three 4K pages.
- Benefits of slab allocation include lack of internal fragmentation and fast allocation of space for individual structures
- Solaris uses slab allocation for the kernel and also for certain user-mode memory allocations. Linux used the buddy system prior to 2.2 and switched to slab allocation since then.
  - **New in 9th Edition: Linux SLOB and SLUB allocators replace SLAB**
    - SLOB, Simple List of Blocks, maintains 3 linked lists of free blocks - small, medium, and large - designed

for ( imbedded ) systems with limited amounts of memory.

- o SLUB modifies some implementation issues for better performance on systems with large numbers of processors.

**Figure 9.27 - Slab allocation.**

# **Other Considerations**

### 9.9.1 Prepaging

- The basic idea behind *prepaging* is to predict the pages that will be needed in the near future, and page them in before they are actually requested.
- If a process was swapped out and we know what its working set was at the time, then when we swap it back in we can go ahead and page back in the entire working set, before the page faults actually occur.
- With small ( data ) files we can go ahead and prepage all of the pages at one time.
- Prepaging can be of benefit if the prediction is good and the pages are needed eventually, but slows the system down if the prediction is wrong.

### 9.9.2 Page Size

- There are quite a few trade-offs of small versus large page sizes:
- Small pages waste less memory due to internal fragmentation.
- Large pages require smaller page tables.
- For disk access, the latency and seek times greatly outweigh the actual data transfer times. This makes it much faster to transfer one large page of data than two or more smaller pages containing the same amount of data.
- Smaller pages match locality better, because we are not bringing in data that is not really needed.
- Small pages generate more page faults, with attending overhead.
- The physical hardware may also play a part in determining page size.
- It is hard to determine an "optimal" page size for any given system. Current norms range from 4K to 4M, and tend towards larger page sizes as time passes.

### 9.9.3 TLB Reach

- *TLB Reach* is defined as the amount of memory that can be reached by the pages listed in the TLB.
- Ideally the working set would fit within the reach of the TLB.
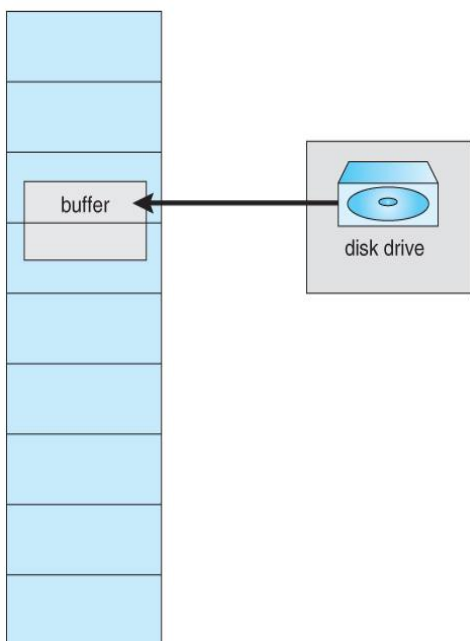
- Increasing the size of the TLB is an obvious way of increasing TLB reach, but TLB memory is very expensive and also draws lots of power.
- Increasing page sizes increases TLB reach, but also leads to increased fragmentation loss.
- Some systems provide multiple size pages to increase TLB reach while keeping fragmentation low.
- Multiple page sizes requires that the TLB be managed by software, not hardware.

### 9.9.4 Inverted Page Tables

- Inverted page tables store one entry for each frame instead of one entry for each virtual page. This reduces the memory requirement for the page table, but loses the information needed to implement virtual memory paging.
- A solution is to keep a separate page table for each process, for virtual memory management purposes. These are kept on disk, and only paged in when a page fault occurs. ( I.e. they are not referenced with every memory access the way a traditional page table would be. )

### 9.9.5 Program Structure

- Consider a pair of nested loops to access every element in a 1024 x 1024 two-dimensional array of 32-bit ints.
- Arrays in C are stored in row-major order, which means that each row of the array would occupy a page of memory.
- If the loops are nested so that the outer loop increments the row and the inner loop increments the column, then an entire row can be processed before the next page fault, yielding 1024 page faults total.
- On the other hand, if the loops are nested the other way, so that the program worked down the columns instead of across the rows, then every access would be to a different page, yielding a new page fault for each access, or over a million page faults all together.
- Be aware that different languages store their arrays differently. FORTRAN for example stores arrays in column-major format instead of row-major. This means that blind translation of code from one language to another may turn a fast program into a very slow one, strictly because of the extra page faults.

### 9.9.6 I/O Interlock and Page Locking

There are several occasions when it may be desirable to *lock* pages in memory, and not let them get paged out:

- Certain kernel operations cannot tolerate having their pages swapped out.
- If an I/O controller is doing direct-memory access, it would be wrong to change pages in the middle of the I/O operation.
- In a priority based scheduling system, low priority jobs may need to wait quite a while before getting their turn on the CPU, and there is a danger of their pages being paged out before they get a chance to use them even once after paging them in. In this situation pages may be locked when they are paged in, until the process that requested them gets at least one turn in the CPU.

**Figure 9.28 - The reason why frames used for I/O must be in memory.**

# Operating-System Examples ( Optional )

### 9.10.1 Windows

- Windows uses demand paging with *clustering,* meaning they page in multiple pages whenever a page fault occurs.
- The working set minimum and maximum are normally set at 50 and 345 pages respectively. ( Maximums can be exceeded in rare circumstances. )
- Free pages are maintained on a free list, with a minimum threshold indicating when there are enough free frames available.
- If a page fault occurs and the process is below their maximum, then additional pages are allocated. Otherwise some pages from this process must be replaced, using a local page replacement algorithm.
- If the amount of free frames falls below the allowable threshold, then *working set trimming* occurs, taking frames away from any processes which are above their minimum, until all are at their minimums. Then additional frames can be allocated to processes that need them.
- The algorithm for selecting victim frames depends on the type of processor:
  - On single processor 80x86 systems, a variation of the clock ( second chance ) algorithm is used.
  - On Alpha and multiprocessor systems, clearing the reference bits may require invalidating entries in the TLB on other processors, which is an expensive operation. In this case Windows uses a variation of FIFO.

### 9.10.2 Solaris

- Solaris maintains a list of free pages, and allocates one to a faulting thread whenever a fault occurs. It is therefore imperative that a minimum amount of free memory be kept on hand at all times.

- Solaris has a parameter, *lotsfree,* usually set at 1/64 of total physical memory. Solaris checks 4 times per second to see if the free memory falls below this threshhold, and if it does, then the *pageout* process is started.

- Pageout uses a variation of the clock ( second chance ) algorithm, with two hands rotating around through the frame table. The first hand clears the reference bits, and the second hand comes by afterwards and checks them. Any frame whose reference bit has not been reset before the second hand gets there gets paged out.

- The Pageout method is adjustable by the distance between the two hands, ( the *handspan* ), and the speed at which the hands move. For example, if the hands each check 100 frames per second, and the handspan is 1000 frames, then there would be a 10 second interval between the time when the leading hand clears the reference bits and the time when the trailing hand checks them.

- The speed of the hands is usually adjusted according to the amount of free memory, as shown below. *Slowscan* is usually set at 100 pages per second, and *fastscan* is usually set at the smaller of 1/2 of the total physical pages per second and 8192 pages per second.



**Figure 9.29 - Solaris page scanner.**

- Solaris also maintains a cache of pages that have been reclaimed but which have not yet been overwritten, as opposed to the free list which only holds pages whose current contents are invalid. If one of the pages from the cache is needed before it gets moved to the free list, then it can be quickly recovered.

- Normally pageout runs 4 times per second to check if memory has fallen below *lotsfree.* However if it falls below *desfree,* then pageout will run at 100 times per second in an attempt to keep at least desfree pages free. If it is unable to do this for a 30-second average, then Solaris begins swapping processes, starting preferably with processes that have been idle for a long time.

- If free memory falls below *minfree,* then pageout runs with every page fault.

- Recent releases of Solaris have enhanced the virtual memory management system, including recognizing pages from shared libraries, and protecting them from being paged out.

# STORAGE MANAGEMENT

**File-System Interface**

**File Concept:**

**File Attributes:**

- Different OSes keep track of different file attributes, including:
  - **Name** - Some systems give special significance to names, and particularly extensions ( .exe, .txt, etc. ), and some do not. Some extensions may be of significance to the OS ( .exe ), and others only to certain applications ( .jpg )
  - **Identifier** ( e.g. inode number )
  - **Type** - Text, executable, other binary, etc.
  - **Location** - on the hard drive.
  - **Size**
  - **Protection**
  - **Time & Date**
  - **User ID**

**File Operations:**

- The file ADT supports many common operations:
  - Creating a file
  - Writing a file
  - Reading a file
  - Repositioning within a file
  - Deleting a file
  - Truncating a file.
- Most Operating Systems require that files be *opened* before access and *closed* after all access is complete. Normally the programmer must open and close files explicitly, but some rare systems open the file automatically at first access. Information about currently open files is stored in an *open file table*, containing for example:
  - **File pointer** - records the current position in the file, for the next read or write access.
  - **File-open count** - How many times has the current file been opened (simultaneously by different processes) and not yet closed? When this counter reaches zero the file can be removed from the table.
  - **Disk location of the file.**
  - **Access rights**
- Some systems provide support for *file locking.*
  - A *shared lock* is for reading only.
  - A *exclusive lock* is for writing as well as reading.
  - An *advisory lock* is informational only, and not enforced. ( A "Keep Out" sign, which may be ignored. )
  - A *mandatory lock* is enforced. ( A truly locked door. )
  - UNIX used advisory locks, and Windows uses mandatory locks.

**File Types:**

| file type | usual extension | function |
|---|---|---|
| executable | exe, com, bin or none | ready-to-run machine-language program |
| object | obj, o | compiled, machine language, not linked |
| source code | c, cc, java, perl, asm | source code in various languages |
| batch | bat, sh | commands to the command interpreter |
| markup | xml, html, tex | textual data, documents |
| word processor | xml, rtf, docx | various word-processor formats |
| library | lib, a, so, dll | libraries of routines for programmers |
| print or view | gif, pdf, jpg | ASCII or binary file in a format for printing or viewing |
| archive | rar, zip, tar | related files grouped into one file, sometimes compressed, for archiving or storage |
| multimedia | mpeg, mov, mp3, mp4, avi | binary file containing audio or A/V information |

- Windows ( and some other systems ) use special file extensions to indicate the type of each file:

**Figure 11.3 - Common file types.**

- Macintosh stores a creator attribute for each file, according to the program that first created it with the create( ) system call.
- UNIX stores magic numbers at the beginning of certain files. ( Experiment with the "file" command, especially in directories such as /bin and /dev )

**File Structure:**

- Some files contain an internal structure, which may or may not be known to the OS.
- For the OS to support particular file formats increases the size and complexity of the OS.
- UNIX treats all files as sequences of bytes, with no further consideration of the internal structure. (With the exception of executable binary programs, which it must know how to load and find the first executable statement, etc. )
- Macintosh files have two *forks* - a *resource fork*, and a *data fork*. The resource fork contains information relating to the UI, such as icons and button images, and can be modified independently of the data fork, which contains the code or data as appropriate.

**Internal File Structure**

- Disk files are accessed in units of physical blocks, typically 512 bytes or some power-of-two multiple thereof. ( Larger physical disks use larger block sizes, to keep the range of block numbers within the range of a 32-bit integer. )
- Internally files are organized in units of logical units, which may be as small as a single byte, or may be a larger size corresponding to some data record or structure size.

- The number of logical units which fit into one physical block determines its *packing*, and has an impact on the amount of internal fragmentation ( wasted space ) that occurs.
- As a general rule, half a physical block is wasted for each file, and the larger the block sizes the more space is lost to internal fragmentation.

## Access Methods

### 1 Sequential Access:

- A sequential access file emulates magnetic tape operation, and generally supports a few operations:
  - read next - read a record and advance the tape to the next position.
  - write next - write a record and advance the tape to the next position.
  - rewind
  - skip n records - May or may not be supported. N may be limited to positive numbers, or may be limited to +/- 1.



**Figure 11.4 - Sequential-access file.**

### 2 Direct Access:

- Jump to any record and read that record. Operations supported include:
  - read n - read record number n. ( Note an argument is now required. )
  - write n - write record number n. ( Note an argument is now required. )
  - jump to record n - could be 0 or the end of file.
  - Query current record - used to return back to this record later.
  - Sequential access can be easily emulated using direct access. The inverse is complicated and inefficient.

| sequential access | implementation for direct access |
|---|---|
| reset | cp = 0; |
| read_next | read cp ;<br>cp = cp + 1; |
| write_next | write cp;<br>cp = cp + 1; |

**Figure 11.5 - Simulation of sequential access on a direct-access file.**

## 3 Other Access Methods:

- An indexed access scheme can be easily built on top of a direct access system. Very large files may require a multi-tiered indexing scheme, i.e. indexes of indexes.
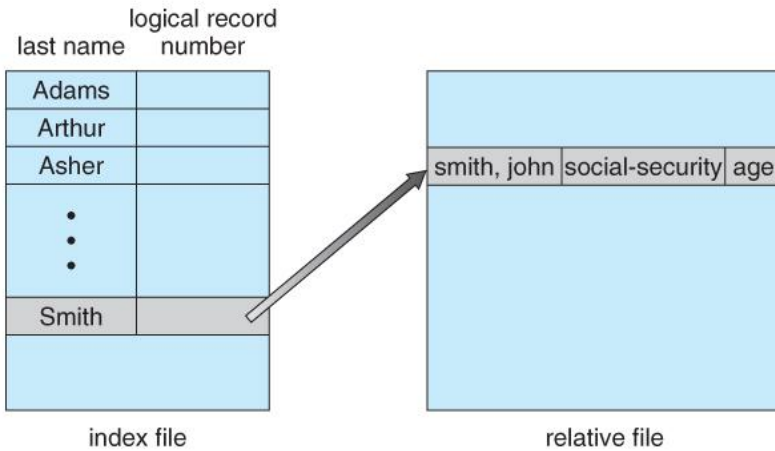


**Figure 11.6 - Example of index and relative files.**

## Directory Structure

### 1 Storage Structure

- A disk can be used in its entirety for a file system.
- Alternatively a physical disk can be broken up into multiple *partitions, slices, or mini-disks*, each of which becomes a virtual disk and can have its own filesystem. ( or be used for raw storage, swap space, etc. )
- Or, multiple physical disks can be combined into one *volume*, i.e. a larger virtual disk, with its own filesystem spanning the physical disks.

| / | ufs |
|---|---|
| /devices | devfs |
| /dev | dev |
| /system/contract | ctfs |
| /proc | proc |
| /etc/mnttab | mntfs |
| /etc/svc/volatile | tmpfs |
| /system/object | objfs |
| /lib/libc.so.1 | lofs |
| /dev/fd | fd |
| /var | ufs |
| /tmp | tmpfs |
| /var/run | tmpfs |
| /opt | ufs |
| /zpbge | zfs |
| /zpbge/backup | zfs |
| /export/home | zfs |
| /var/mail | zfs |
| /var/spool/mqueue | zfs |
| /zpbg | zfs |
| /zpbg/zones | zfs |

**Figure 11.8** Solaris file systems.



**Figure 11.7 - A typical file-system organization.**

## 2 Directory Overview

- Directory operations to be supported include:
    - Search for a file
    - Create a file - add to the directory
    - Delete a file - erase from the directory
    - List a directory - possibly ordered in different ways.
    - Rename a file - may change sorting order
    - Traverse the file system.

## 3. Single-Level Directory

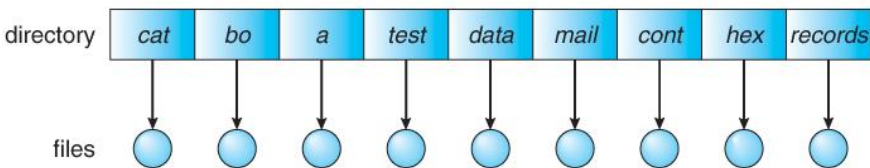- Simple to implement, but each file must have a unique name.



**Figure 11.9 - Single-level directory.**

## 4 Two-Level Directory

- Each user gets their own directory space.

- File names only need to be unique within a given user's directory.

- A master file directory is used to keep track of each users directory, and must be maintained when users are added to or removed from the system.

- A separate directory is generally needed for system ( executable ) files.

- Systems may or may not allow users to access other directories besides their own

    - If access to other directories is allowed, then provision must be made to specify the directory being accessed.

    - If access is denied, then special consideration must be made for users to run programs located in system directories. A **search path** is the list of directories in which to search for executable programs, and can be set uniquely for each user.



**Figure 11.10 - Two-level directory structure.**

## 5 Tree-Structured Directories

- An obvious extension to the two-tiered directory structure, and the one with which we are all most familiar.
- Each user / process has the concept of a *current directory* from which all ( relative ) searches take place.
- Files may be accessed using either absolute pathnames ( relative to the root of the tree ) or relative pathnames ( relative to the current directory. )
- Directories are stored the same as any other file in the system, except there is a bit that identifies them as directories, and they have some special structure that the OS understands.
- One question for consideration is whether or not to allow the removal of directories that are not empty - Windows requires that directories be emptied first, and UNIX provides an option for deleting entire sub-trees.



**Figure 11.11 - Tree-structured directory structure.**

## 6 Acyclic-Graph Directories

- When the same files need to be accessed in more than one place in the directory structure ( e.g. because they are being shared by more than one user / process ), it can be useful to provide an acyclic-graph structure. ( Note the *directed* arcs from parent to child. )
- UNIX provides two types of *links* for implementing the acyclic-graph structure. ( See "man ln" for more details. )
  - A *hard link* ( usually just called a link ) involves multiple directory entries that both refer to the same file. Hard links are only valid for ordinary files in the same filesystem.
  - A *symbolic link*, that involves a special file, containing information about where to find the linked file. Symbolic links may be used to link directories and/or files in other filesystems, as well as ordinary files in the current filesystem.
- Windows only supports symbolic links, termed *shortcuts.*

- Hard links require a *reference count*, or *link count* for each file, keeping track of how many directory entries are currently referring to this file. Whenever one of the references is removed the link count is reduced, and when it reaches zero, the disk space can be reclaimed.

- For symbolic links there is some question as to what to do with the symbolic links when the original file is moved or deleted:

  o One option is to find all the symbolic links and adjust them also.

  o Another is to leave the symbolic links dangling, and discover that they are no longer valid the next time they are used.

  o What if the original file is removed, and replaced with another file having the same name before the symbolic link is next used?



**Figure 11.12 - Acyclic-graph directory structure.**

**7 General Graph Directory**

- If cycles are allowed in the graphs, then several problems can arise:

  o Search algorithms can go into infinite loops. One solution is to not follow links in search algorithms. ( Or not to follow symbolic links, and to only allow symbolic links to refer to directories. )

  o Sub-trees can become disconnected from the rest of the tree and still not have their reference counts reduced to zero. Periodic garbage collection is required to detect and resolve this problem. ( chkdsk in DOS and fsck in UNIX search for these problems, among others, even though cycles are not supposed to be allowed in either system. Disconnected disk blocks that are not marked as free are added back to the file systems with made-up file names, and can usually be safely deleted. )
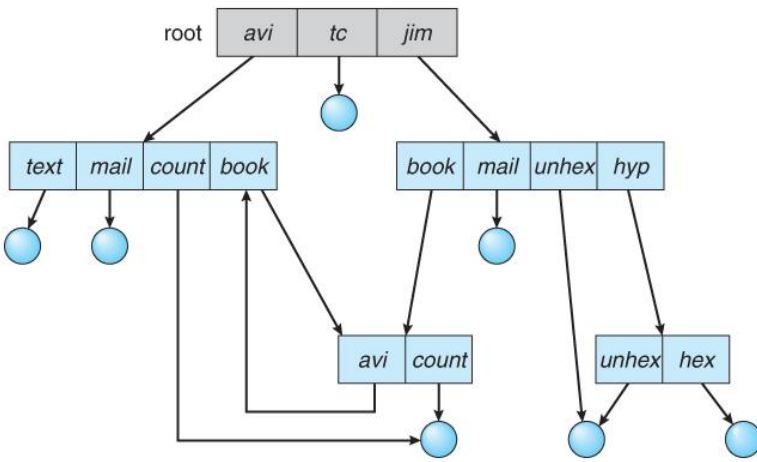
**Figure 11.13 - General graph directory.**

## File-System Mounting

- The basic idea behind mounting file systems is to combine multiple file systems into one large tree structure.

- The mount command is given a filesystem to mount and a *mount point* ( directory ) on which to attach it.

- Once a file system is mounted onto a mount point, any further references to that directory actually refer to the root of the mounted file system.

- Any files ( or sub-directories ) that had been stored in the mount point directory prior to mounting the new filesystem are now hidden by the mounted filesystem, and are no longer available. For this reason some systems only allow mounting onto empty directories.

- Filesystems can only be mounted by root, unless root has previously configured certain filesystems to be mountable onto certain pre-determined mount points. ( E.g. root may allow users to mount floppy filesystems to /mnt or something like it. ) Anyone can run the mount command to see what filesystems are currently mounted.

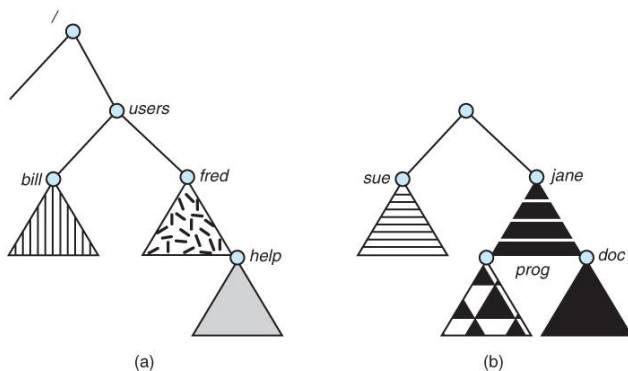- Filesystems may be mounted read-only, or have other restrictions imposed.



**Figure 11.14 - File system. (a) Existing system. (b) Unmounted volume.**
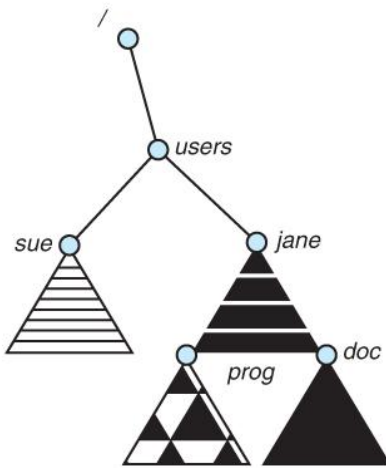
**Figure 11.15 - Mount point.**

- The traditional Windows OS runs an extended two-tier directory structure, where the first tier of the structure separates volumes by drive letters, and a tree structure is implemented below that level.

- Macintosh runs a similar system, where each new volume that is found is automatically mounted and added to the desktop when it is found.

- More recent Windows systems allow filesystems to be mounted to any directory in the filesystem, much like UNIX.

## File Sharing

### 1 Multiple Users

- On a multi-user system, more information needs to be stored for each file:
    - The owner ( user ) who owns the file, and who can control its access.
    - The group of other user IDs that may have some special access to the file.
    - What access rights are afforded to the owner ( **U**ser ), the **G**roup, and to the rest of the world ( the universe, a.k.a. **O**thers. )
    - Some systems have more complicated access control, allowing or denying specific accesses to specifically named users or groups.

### 2 Remote File Systems

- The advent of the Internet introduces issues for accessing files stored on remote computers
    - The original method was ftp, allowing individual files to be transported across systems as needed. Ftp can be either account and password controlled, or *anonymous*, not requiring any user name or password.
    - Various forms of *distributed file systems* allow remote file systems to be mounted onto a local directory structure, and accessed using normal file access commands.

175

( The actual files are still transported across the network as needed, possibly using ftp as the underlying transport mechanism. )

o The WWW has made it easy once again to access files on remote systems without mounting their filesystems, generally using ( anonymous ) ftp as the underlying file transport mechanism.

## 2.1 The Client-Server Model

- When one computer system remotely mounts a filesystem that is physically located on another system, the system which physically owns the files acts as a *server*, and the system which mounts them is the *client.*

- User IDs and group IDs must be consistent across both systems for the system to work properly. ( I.e. this is most applicable across multiple computers managed by the same organization, shared by a common group of users. )

- The same computer can be both a client and a server. ( E.g. cross-linked file systems. )

- There are a number of security concerns involved in this model:

    o Servers commonly restrict mount permission to certain trusted systems only. Spoofing ( a computer pretending to be a different computer ) is a potential security risk.

    o Servers may restrict remote access to read-only.

    o Servers restrict which filesystems may be remotely mounted. Generally the information within those subsystems is limited, relatively public, and protected by frequent backups.

- The NFS ( Network File System ) is a classic example of such a system.

## 2.2 Distributed Information Systems

- The *Domain Name System, DNS,* provides for a unique naming system across all of the Internet.

- Domain names are maintained by the *Network Information System, NIS*, which unfortunately has several security issues. NIS+ is a more secure version, but has not yet gained the same widespread acceptance as NIS.

- Microsoft's *Common Internet File System, CIFS*, establishes a *network login* for each user on a networked system with shared file access. Older Windows systems used *domains*, and newer systems ( XP, 2000 ), use *active directories.* User names must match across the network for this system to be valid.

176

- A newer approach is the *Lightweight Directory-Access Protocol, LDAP,* which provides a *secure single sign-on* for all users to access all resources on a network. This is a secure system which is gaining in popularity, and which has the maintenance advantage of combining authorization information in one central location.

**2.3 Failure Modes**

- When a local disk file is unavailable, the result is generally known immediately, and is generally non-recoverable. The only reasonable response is for the response to fail.
- However when a remote file is unavailable, there are many possible reasons, and whether or not it is unrecoverable is not readily apparent. Hence most remote access systems allow for blocking or delayed response, in the hopes that the remote system ( or the network ) will come back up eventually.

## Protection

- Files must be kept safe for reliability ( against accidental damage ), and protection ( against deliberate malicious access. ) The former is usually managed with backup copies. This section discusses the latter.
- One simple protection scheme is to remove all access to a file. However this makes the file unusable, so some sort of controlled access must be arranged.

**1 Types of Access**

- The following low-level operations are often controlled:
  o Read - View the contents of the file
  o Write - Change the contents of the file.
  o Execute - Load the file onto the CPU and follow the instructions contained therein.
  o Append - Add to the end of an existing file.
  o Delete - Remove a file from the system.
  o List -View the name and other attributes of files on the system.
- Higher-level operations, such as copy, can generally be performed through combinations of the above.

**2 Access Control**

- One approach is to have complicated *Access Control Lists, ACL,* which specify exactly what access is allowed or denied for specific users or groups.
    - The AFS uses this system for distributed access.
    - Control is very finely adjustable, but may be complicated, particularly when the specific users involved are unknown. ( AFS allows some wild cards, so for example all users on a certain remote system may be trusted, or a given username may be trusted when accessing from any remote system. )
- UNIX uses a set of 9 access control bits, in three groups of three. These correspond to R, W, and X permissions for each of the Owner, Group, and Others. ( See "man chmod" for full details. ) The RWX bits control the following privileges for ordinary files and directories:

| bit | Files | Directories |
|-----|-------|-------------|
| R | Read ( view ) file contents. | Read directory contents. Required to get a listing of the directory. |
| W | Write ( change ) file contents. | Change directory contents. Required to create or delete files. |
| X | Execute file contents as a program. | Access detailed directory information. Required to get a long listing, or to access any specific file in the directory. Note that if a user has X but not R permissions on a directory, they can still access specific files, but only if they already know the name of the file they are trying to access. |

- In addition there are some special bits that can also be applied:
    - The set user ID ( SUID ) bit and/or the set group ID ( SGID ) bits applied to executable files temporarily change the identity of whoever runs the program to match that of the owner / group of the executable program. This allows users running specific programs to have access to files ( *while running that program* ) to which they would normally be unable to access. Setting of these two bits is usually restricted to root, and must be done with caution, as it introduces a potential security leak.
    - The sticky bit on a directory modifies write permission, allowing users to only delete files for which they are the owner. This allows everyone to create files in /tmp, for example, but to only delete files which they have created, and not anyone else's.
    - The SUID, SGID, and sticky bits are indicated with an S, S, and T in the positions for execute permission for the user, group, and others, respectively. If the letter is

lower case, ( s, s, t ), then the corresponding execute permission is not also given. If it is upper case, ( S, S, T ), then the corresponding execute permission IS given.

- o The numeric form of chmod is needed to set these advanced bits.

```
-rw-rw-r--    1 pbg   staff      31200  Sep 3 08:30   intro.ps
drwx------    5 pbg   staff        512  Jul 8 09.33    private/
drwxrwxr-x    2 pbg   staff        512  Jul 8 09:35    doc/
drwxrwx---    2 jwg   student      512  Aug 3 14:13   student-proj/
-rw-r--r--    1 pbg   staff       9423  Feb 24 2012   program.c
-rwxr-xr-x    1 pbg   staff      20471  Feb 24 2012   program
drwx--x--x    4 tag   faculty      512  Jul 31 10:31   lib/
drwx------    3 pbg   staff       1024  Aug 29 06:52  mail/
drwxrwxrwx    3 pbg   staff        512  Jul 8 09:35    test/
```

**Sample permissions in a UNIX system.**

- Windows adjusts files access through a simple GUI:



**Figure 11.16 - Windows 7 access-control list management.**

### 11.6.3 Other Protection Approaches and Issues

- Some systems can apply passwords, either to individual files, or to specific sub-directories, or to the entire system. There is a trade-off between the number of passwords that must be maintained ( and remembered by the users ) and the amount of information that is vulnerable to a lost or forgotten password.

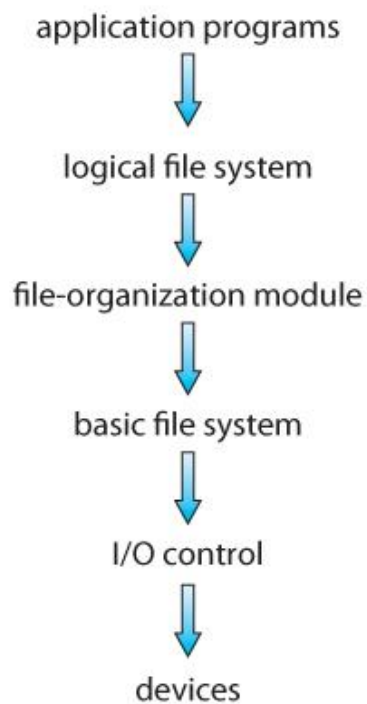- Older systems which did not originally have multi-user file access permissions ( DOS and older versions of Mac ) must now be *retrofitted* if they are to share files on a network.

- Access to a file requires access to all the files along its path as well. In a cyclic directory structure, users may have different access to the same file accessed through different paths.

- Sometimes just the knowledge of the existence of a file of a certain name is a security ( or privacy ) concern. Hence the distinction between the R and X bits on UNIX directories.

# File-System Implementation

## File-System Structure

- Hard disks have two important properties that make them suitable for secondary storage of files in file systems: (1) Blocks of data can be rewritten in place, and (2) they are direct access, allowing any block of data to be accessed with only ( relatively ) minor movements of the disk heads and rotational latency. ( See Chapter 12 )
- Disks are usually accessed in physical blocks, rather than a byte at a time. Block sizes may range from 512 bytes to 4K or larger.
- File systems organize storage on disk drives, and can be viewed as a layered design:
  - At the lowest layer are the physical devices, consisting of the magnetic media, motors & controls, and the electronics connected to them and controlling them. Modern disk put more and more of the electronic controls directly on the disk drive itself, leaving relatively little work for the disk controller card to perform.
  - *I/O Control* consists of *device drivers*, special software programs ( often written in assembly ) which communicate with the devices by reading and writing special codes directly to and from memory addresses corresponding to the controller card's registers. Each controller card ( device ) on a system has a different set of addresses ( registers, a.k.a. *ports* ) that it listens to, and a unique set of command codes and results codes that it understands.
  - The *basic file system* level works directly with the device drivers in terms of retrieving and storing raw blocks of data, without any consideration for what is in each block. Depending on the system, blocks may be referred to with a single block number, ( e.g. block # 234234 ), or with head-sector-cylinder combinations.
  - The *file organization module* knows about files and their logical blocks, and how they map to physical blocks on the disk. In addition to translating from logical to physical blocks, the file organization module also maintains the list of free blocks, and allocates free blocks to files as needed.
  - The *logical file system* deals with all of the meta data associated with a file ( UID, GID, mode, dates, etc ), i.e. everything about the file except the data itself. This level manages the directory structure and the mapping of file names to *file control blocks, FCBs*, which contain all of the meta data as well as block number information for finding the data on the disk.
- The layered approach to file systems means that much of the code can be used uniformly for a wide variety of different file systems, and only certain layers need to be filesystem specific. Common file systems in use include the UNIX file system, UFS, the Berkeley Fast File System, FFS, Windows systems FAT, FAT32, NTFS, CD-ROM systems ISO 9660, and for Linux the extended file systems ext2 and ext3 ( among 40 others supported. )

**Figure 12.1 - Layered file system.**

## File-System Implementation

### 1 Overview

- File systems store several important data structures on the disk:

- o A ***boot-control block***, ( per volume ) a.k.a. the ***boot block*** in UNIX or the ***partition boot sector*** in Windows contains information about how to boot the system off of this disk. This will generally be the first sector of the volume if there is a bootable system loaded on that volume, or the block will be left vacant otherwise.
- o A ***volume control block,*** ( per volume ) a.k.a. the ***master file table*** in UNIX or the ***superblock*** in Windows, which contains information such as the partition table, number of blocks on each filesystem, and pointers to free blocks and free FCB blocks.
- o A directory structure ( per file system ), containing file names and pointers to corresponding FCBs. UNIX uses inode numbers, and NTFS uses a ***master file table.***
- o The ***File Control Block, FCB,*** ( per file ) containing details about ownership, size, permissions, dates, etc. UNIX stores this information in inodes, and NTFS in the master file table as a relational database structure.

| file permissions |
| file dates (create, access, write) |
| file owner, group, ACL |
| file size |
| file data blocks or pointers to file data blocks |

**Figure 12.2 - A typical file-control block.**

- • There are also several key data structures stored in memory:
  - o An in-memory mount table.
  - o An in-memory directory cache of recently accessed directory information.
  - o *A **system-wide open file table***, containing a copy of the FCB for every currently open file in the system, as well as some other related information.
  - o *A **per-process open file table,*** containing a pointer to the system open file table as well as some other information. ( For example the current file position pointer may be either here or in the system file table, depending on the implementation and whether the file is being shared or not. )
- • Figure 12.3 illustrates some of the interactions of file system components when files are created and/or used:
  - o When a new file is created, a new FCB is allocated and filled out with important information regarding the new file. The appropriate directory is modified with the new file name and FCB information.
  - o When a file is accessed during a program, the open( ) system call reads in the FCB information from disk, and stores it in the system-wide open file table. An entry is added to the per-process open file table referencing the system-wide table, and an index into the per-process table is returned by the open( ) system call. UNIX refers to this index as a ***file descriptor***, and Windows refers to it as a ***file handle***.
  - o If another process already has a file open when a new request comes in for the same file, and it is sharable, then a counter in the system-wide table is incremented and the per-process table is adjusted to point to the existing entry in the system-wide table.
  - o When a file is closed, the per-process table entry is freed, and the counter in the system-wide table is decremented. If that counter reaches zero, then the system

wide table is also freed. Any data currently stored in memory cache for this file is written out to disk if necessary.
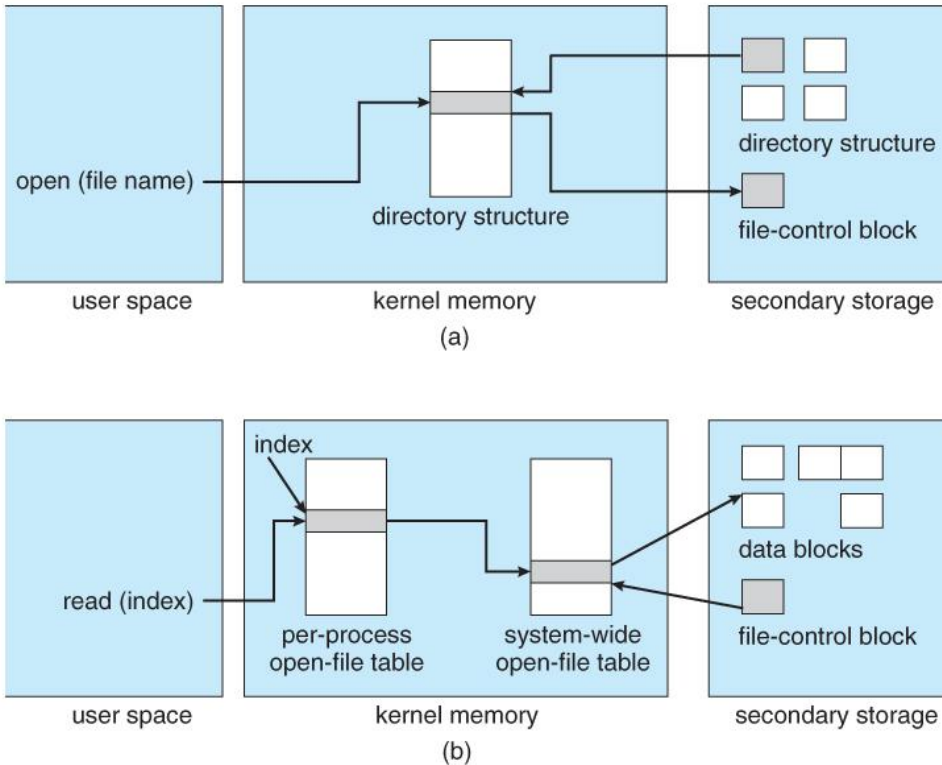


**Figure 12.3 - In-memory file-system structures. (a) File open. (b) File read.**

## 2 Partitions and Mounting

- Physical disks are commonly divided into smaller units called partitions. They can also be combined into larger units, but that is most commonly done for RAID installations and is left for later chapters.
- Partitions can either be used as raw devices ( with no structure imposed upon them ), or they can be formatted to hold a filesystem ( i.e. populated with FCBs and initial directory structures as appropriate. ) Raw partitions are generally used for swap space, and may also be used for certain programs such as databases that choose to manage their own disk storage system. Partitions containing filesystems can generally only be accessed using the file system structure by ordinary users, but can often be accessed as a raw device also by root.
- The boot block is accessed as part of a raw partition, by the boot program prior to any operating system being loaded. Modern boot programs understand multiple OSes and filesystem formats, and can give the user a choice of which of several available systems to boot.
- The *root partition* contains the OS kernel and at least the key portions of the OS needed to complete the boot process. At boot time the root partition is mounted, and control is transferred from the boot program to the kernel found there. ( Older systems required that the root partition lie completely within the first 1024 cylinders of the disk, because that was as far as the boot program could reach. Once the kernel had control, then it could access partitions beyond the 1024 cylinder boundary. )
- Continuing with the boot process, additional filesystems get mounted, adding their information into the appropriate mount table structure. As a part of the mounting process the file systems may be checked for errors or inconsistencies, either because they are flagged as not having been closed properly the last time they were used, or just for general

183

principals. Filesystems may be mounted either automatically or manually. In UNIX a mount point is indicated by setting a flag in the in-memory copy of the inode, so all future references to that inode get re-directed to the root directory of the mounted filesystem.

## 3 Virtual File Systems

- *Virtual File Systems, VFS*, provide a common interface to multiple different filesystem types. In addition, it provides for a unique identifier ( vnode ) for files across the entire space, including across all filesystems of different types. ( UNIX inodes are unique only across a single filesystem, and certainly do not carry across networked file systems. )
- The VFS in Linux is based upon four key object types:
    o The *inode* object, representing an individual file
    o The *file* object, representing an open file.
    o The *superblock* object, representing a filesystem.
    o The *dentry* object, representing a directory entry.
- Linux VFS provides a set of common functionalities for each filesystem, using function pointers accessed through a table. The same functionality is accessed through the same table position for all filesystem types, though the actual functions pointed to by the pointers may be filesystem-specific. See /usr/include/linux/fs.h for full details. Common operations provided include open( ), read( ), write( ), and mmap( ).
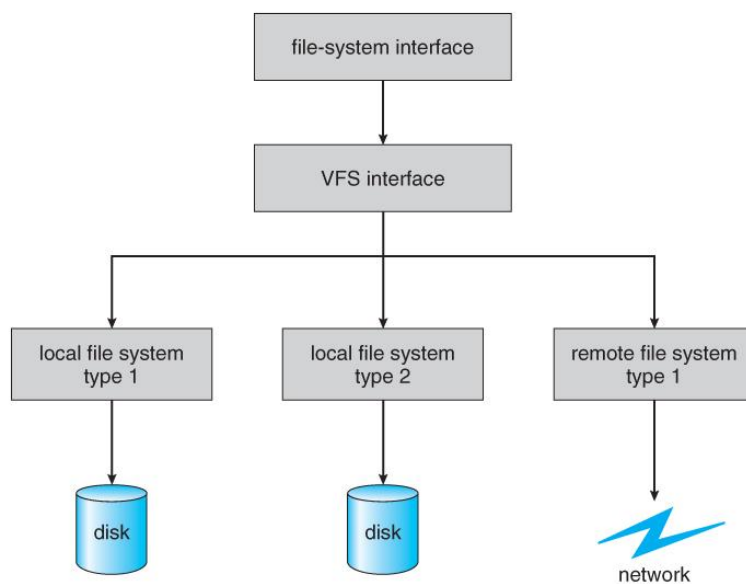


**Figure 12.4 - Schematic view of a virtual file system.**

## 3 Directory Implementation

- Directories need to be fast to search, insert, and delete, with a minimum of wasted disk space.

### 1 Linear List

- A linear list is the simplest and easiest directory structure to set up, but it does have some drawbacks.
- Finding a file ( or verifying one does not already exist upon creation ) requires a linear search.
- Deletions can be done by moving all entries, flagging an entry as deleted, or by moving the last entry into the newly vacant position.

- Sorting the list makes searches faster, at the expense of more complex insertions and deletions.
- A linked list makes insertions and deletions into a sorted list easier, with overhead for the links.
- More complex data structures, such as B-trees, could also be considered.

**2 Hash Table**

- A hash table can also be used to speed up searches.
- Hash tables are generally implemented *in addition to* a linear or other structure

## Allocation Methods

- There are three major methods of storing files on disks: contiguous, linked, and indexed.

**1 Contiguous Allocation**

- *Contiguous Allocation* requires that all blocks of a file be kept together contiguously.
- Performance is very fast, because reading successive blocks of the same file generally requires no movement of the disk heads, or at most one small step to the next adjacent cylinder.
- Storage allocation involves the same issues discussed earlier for the allocation of contiguous blocks of memory ( first fit, best fit, fragmentation problems, etc. ) The distinction is that the high time penalty required for moving the disk heads from spot to spot may now justify the benefits of keeping files contiguously when possible.
- ( Even file systems that do not by default store files contiguously can benefit from certain utilities that compact the disk and make all files contiguous in the process. )
- Problems can arise when files grow, or if the exact size of a file is unknown at creation time:
  o Over-estimation of the file's final size increases external fragmentation and wastes disk space.
  o Under-estimation may require that a file be moved or a process aborted if the file grows beyond its originally allocated space.
  o If a file grows slowly over a long time period and the total final space must be allocated initially, then a lot of space becomes unusable before the file fills the space.
- A variation is to allocate file space in large contiguous chunks, called *extents.* When a file outgrows its original extent, then an additional one is allocated. ( For example an extent may be the size of a complete track or even cylinder, aligned on an appropriate track or cylinder boundary. ) The high-performance files system Veritas uses extents to optimize performance.
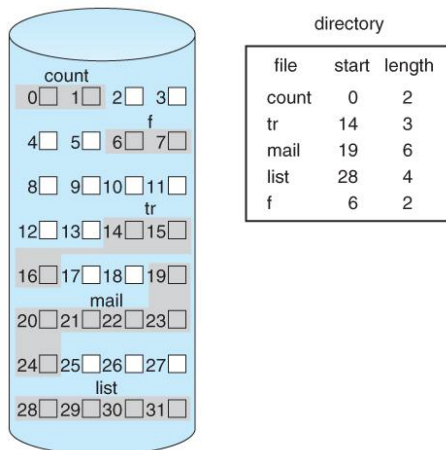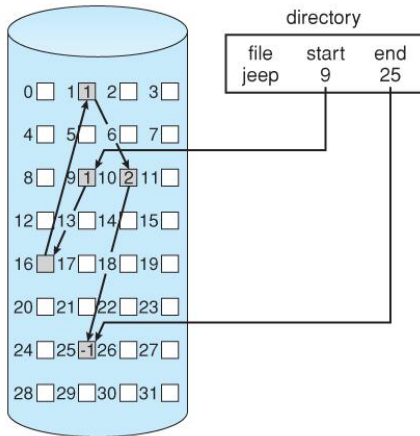


| file | start | length |
|---|---|---|
| count | 0 | 2 |
| tr | 14 | 3 |
| mail | 19 | 6 |
| list | 28 | 4 |
| f | 6 | 2 |

185

**Figure 12.5 - Contiguous allocation of disk space.**

## 2 Linked Allocation

- Disk files can be stored as linked lists, with the expense of the storage space consumed by each link. ( E.g. a block may be 508 bytes instead of 512. )
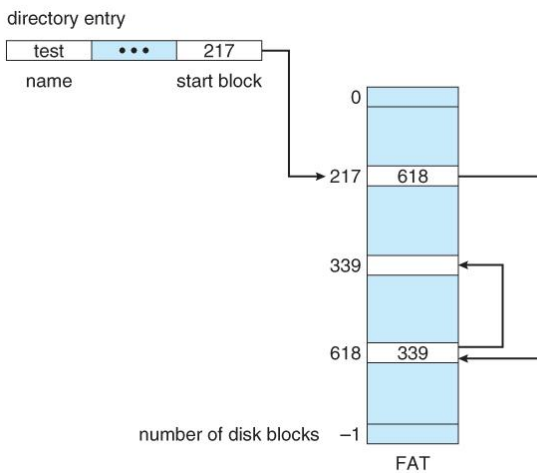


- Linked allocation involves no external fragmentation, does not require pre-known file sizes, and allows files to grow dynamically at any time.
- Unfortunately linked allocation is only efficient for sequential access files, as random access requires starting at the beginning of the list for each new location access.
- Allocating *clusters* of blocks reduces the space wasted by pointers, at the cost of internal fragmentation.
- Another big problem with linked allocation is reliability if a pointer is lost or damaged. Doubly linked lists provide some protection, at the cost of additional overhead and wasted space.
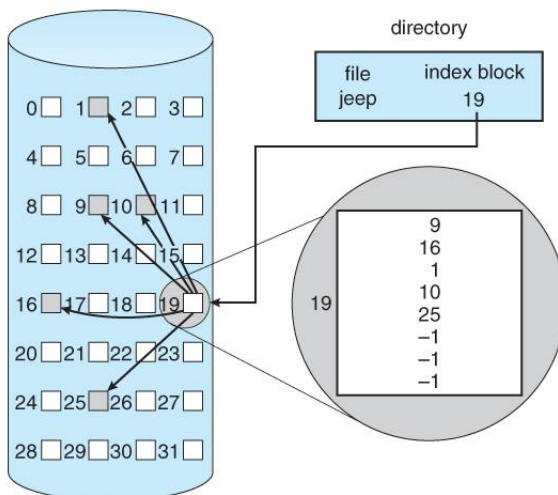
**Figure 12.6 - Linked allocation of disk space.**



- The *File Allocation Table, FAT,* used by DOS is a variation of linked allocation, where all the links are stored in a separate table at the beginning of the disk. The benefit of this approach is that the FAT table can be cached in memory, greatly improving random access speeds.

**Figure 12.7 File-allocation table.**

## 3 Indexed Allocation



- *Indexed Allocation* combines all of the indexes for accessing each file into a common block ( for that file ), as opposed to spreading them all over the disk or storing them in a FAT table.

**Figure 12.8 - Indexed allocation of disk space.**

- Some disk space is wasted ( relative to linked lists or FAT tables ) because an entire index block must be allocated for each file, regardless of

186

how many data blocks the file contains. This leads to questions of how big the index block should be, and how it should be implemented. There are several approaches:

- o **Linked Scheme -** An index block is one disk block, which can be read and written in a single disk operation. The first index block contains some header information, the first N block addresses, and if necessary a pointer to additional linked index blocks.
- o **Multi-Level Index -** The first index block contains a set of pointers to secondary index blocks, which in turn contain pointers to the actual data blocks.
- o **Combined Scheme -** This is the scheme used in UNIX inodes, in which the first 12 or so data block pointers are stored directly in the inode, and then singly, doubly, and triply indirect pointers provide access to more data blocks as needed. ( See below. ) The advantage of this scheme is that for small files ( which many are ), the data blocks are readily accessible ( up to 48K with 4K block sizes ); files up to about 4144K ( using 4K blocks ) are accessible with only a single indirect block ( which can be cached ), and huge files are still accessible using a relatively small number of disk accesses ( larger in theory than can be addressed by a 32-bit address, which is why some systems have moved to 64-bit file pointers. )
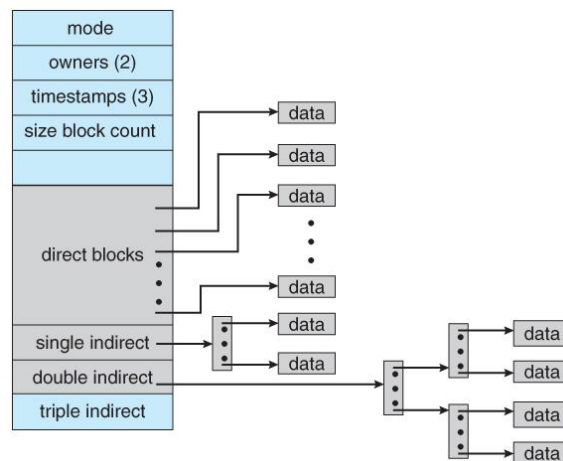


**Figure 12.9 - The UNIX inode.**

## 4 Performance

- The optimal allocation method is different for sequential access files than for random access files, and is also different for small files than for large files.
- Some systems support more than one allocation method, which may require specifying how the file is to be used ( sequential or random access ) at the time it is allocated. Such systems also provide conversion utilities.
- Some systems have been known to use contiguous access for small files, and automatically switch to an indexed scheme when file sizes surpass a certain threshold.
- And of course some systems adjust their allocation schemes ( e.g. block sizes ) to best match the characteristics of the hardware for optimum performance.

## Free-Space Management

- Another important aspect of disk management is keeping track of and allocating free space.

## 1 Bit Vector

- One simple approach is to use a *bit vector*, in which each bit represents a disk block, set to 1 if free or 0 if allocated.
- Fast algorithms exist for quickly finding contiguous blocks of a given size
- The down side is that a 40GB disk requires over 5MB just to store the bitmap. ( For example. )

## 2 Linked List

- A linked list can also be used to keep track of all free blocks.
- Traversing the list and/or finding a contiguous block of a given size are not easy, but fortunately are not frequently needed operations. Generally the system just adds and removes single blocks from the beginning of the list.
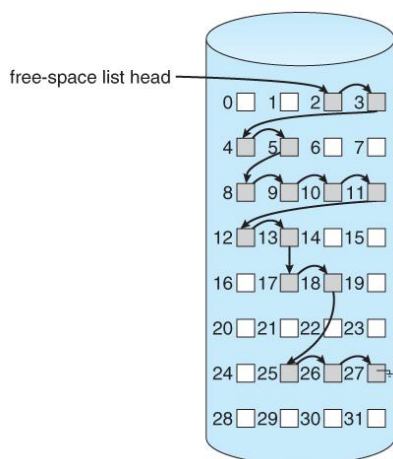- The FAT table keeps track of the free list as just one more linked list on the table.



**Figure 12.10 - Linked free-space list on disk.**

## 3 Grouping

- A variation on linked list free lists is to use links of blocks of indices of free blocks. If a block holds up to N addresses, then the first block in the linked-list contains up to N-1 addresses of free blocks and a pointer to the next block of free addresses.

## 4 Counting

- When there are multiple contiguous blocks of free space then the system can keep track of the starting address of the group and the number of contiguous free blocks. As long as the average length of a contiguous group of free blocks is greater than two this offers a savings in space needed for the free list. ( Similar to compression techniques used for graphics images when a group of pixels all the same color is encountered. )

## 5 Space Maps

- Sun's ZFS file system was designed for HUGE numbers and sizes of files, directories, and even file systems.
- The resulting data structures could be VERY inefficient if not implemented carefully. For example, freeing up a 1 GB file on a 1 TB file system could involve updating thousands of blocks of free list bit maps if the file was spread across the disk.
- ZFS uses a combination of techniques, starting with dividing the disk up into ( hundreds of ) *metaslabs* of a manageable size, each having their own space map.

- Free blocks are managed using the counting technique, but rather than write the information to a table, it is recorded in a log-structured transaction record. Adjacent free blocks are also coalesced into a larger single free block.
- An in-memory space map is constructed using a balanced tree data structure, constructed from the log data.
- The combination of the in-memory tree and the on-disk log provide for very fast and efficient management of these very large files and free blocks.

## Efficiency and Performance

### 1 Efficiency

- UNIX pre-allocates inodes, which occupies space even before any files are created.
- UNIX also distributes inodes across the disk, and tries to store data files near their inode, to reduce the distance of disk seeks between the inodes and the data.
- Some systems use variable size clusters depending on the file size.
- The more data that is stored in a directory ( e.g. last access time ), the more often the directory blocks have to be re-written.
- As technology advances, addressing schemes have had to grow as well.
  - Sun's ZFS file system uses 128-bit pointers, which should theoretically never need to be expanded. ( The mass required to store $2^{128}$ bytes with atomic storage would be at least 272 trillion kilograms! )
- Kernel table sizes used to be fixed, and could only be changed by rebuilding the kernels. Modern tables are dynamically allocated, but that requires more complicated algorithms for accessing them.

### 2 Performance

- Disk controllers generally include on-board caching. When a seek is requested, the heads are moved into place, and then an entire track is read, starting from whatever sector is currently under the heads ( reducing latency. ) The requested sector is returned and the unrequested portion of the track is cached in the disk's electronics.
- Some OSes cache disk blocks they expect to need again in a *buffer cache.*
- A *page cache* connected to the virtual memory system is actually more efficient as memory addresses do not need to be converted to disk block addresses and back again.
- Some systems ( Solaris, Linux, Windows 2000, NT, XP ) use page caching for both process pages and file data in a *unified virtual memory.*
- Figures 11.11 and 11.12 show the advantages of the *unified buffer cache* found in some versions of UNIX and Linux - Data does not need to be stored twice, and problems of inconsistent buffer information are avoided.
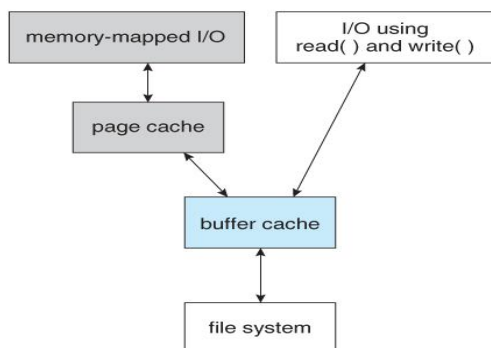


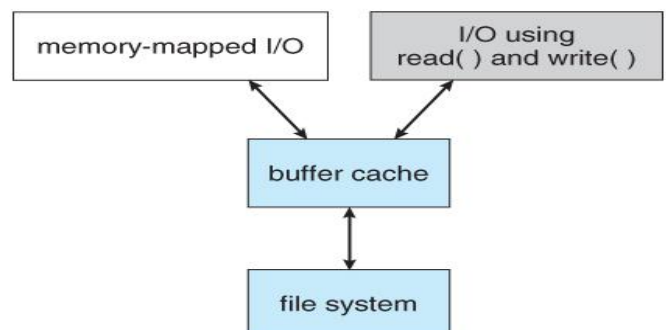**Figure 12.11 - I/O without a unified buffer cache.**

**Figure 12.12 - I/O using a unified buffer cache.**

- Page replacement strategies can be complicated with a unified cache, as one needs to decide whether to replace process or file pages, and how many pages to guarantee to each category of pages. Solaris, for example, has gone through many variations, resulting in *priority paging* giving process pages priority over file I/O pages, and setting limits so that neither can knock the other completely out of memory.
- Another issue affecting performance is the question of whether to implement *synchronous writes* or *asynchronous writes.* Synchronous writes occur in the order in which the disk subsystem receives them, without caching; Asynchronous writes are cached, allowing the disk subsystem to schedule writes in a more efficient order ( See Chapter 12. ) Metadata writes are often done synchronously. Some systems support flags to the open call requiring that writes be synchronous, for example for the benefit of database systems that require their writes be performed in a required order.
- The type of file access can also have an impact on optimal page replacement policies. For example, LRU is not necessarily a good policy for sequential access files. For these types of files progression normally goes in a forward direction only, and the most recently used page will not be needed again until after the file has been rewound and re-read from the beginning, ( if it is ever needed at all. ) On the other hand, we can expect to need the next page in the file fairly soon. For this reason sequential access files often take advantage of two special policies:
    - *Free-behind* frees up a page as soon as the next page in the file is requested, with the assumption that we are now done with the old page and won't need it again for a long time.
    - *Read-ahead* reads the requested page and several subsequent pages at the same time, with the assumption that those pages will be needed in the near future. This is similar to the track caching that is already performed by the disk controller, except it saves the future latency of transferring data from the disk controller memory into motherboard main memory.
- The caching system and asynchronous writes speed up disk writes considerably, because the disk subsystem can schedule physical writes to the disk to minimize head movement and disk seek times. ( See Chapter 12. ) Reads, on the other hand, must be done more synchronously in spite of the caching system, with the result that disk writes can counter-intuitively be much faster on average than disk reads.

# Mass-Storage Structure

## References:

1. Abraham Silberschatz, Greg Gagne, and Peter Baer Galvin, "Operating System Concepts, Ninth Edition ", Chapter 10 ( Was chapter 12 )

## 10.1 Overview of Mass-Storage Structure

### 10.1.1 Magnetic Disks

- Traditional magnetic disks have the following basic structure:
    - One or more **platters** in the form of disks covered with magnetic media. **Hard disk** platters are made of rigid metal, while "**floppy**" disks are made of more flexible plastic.
    - Each platter has two working **surfaces.** Older hard disk drives would sometimes not use the very top or bottom surface of a stack of platters, as these surfaces were more susceptible to potential damage.
    - Each working surface is divided into a number of concentric rings called **tracks.** The collection of all tracks that are the same distance from the edge of the platter, ( i.e. all tracks immediately above one another in the following diagram ) is called a **cylinder**.
    - Each track is further divided into **sectors,** traditionally containing 512 bytes of data each, although some modern disks occasionally use larger sector sizes. ( Sectors also include a header and a trailer, including checksum information among other things. Larger sector sizes reduce the fraction of the disk consumed by headers and trailers, but increase internal fragmentation and the amount of disk that must be marked bad in the case of errors. )
    - The data on a hard drive is read by read-write **heads.** The standard configuration ( shown below ) uses one head per surface, each on a separate **arm**, and controlled by a common **arm assembly** which moves all heads simultaneously from one cylinder to another. ( Other configurations, including independent read-write heads, may speed up disk access, but involve serious technical difficulties. )
    - The storage capacity of a traditional disk drive is equal to the number of heads ( i.e. the number of working surfaces ), times the number of tracks per surface, times the number of sectors per track, times the number of bytes per sector. A particular physical block of data is specified by providing the head-sector-cylinder number at which it is located.
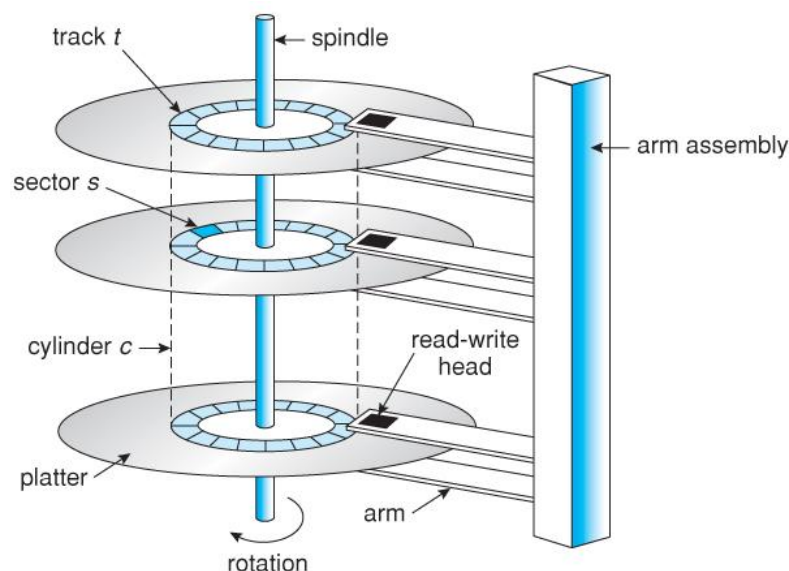


Figure 10.1 - Moving-head disk mechanism.

- In operation the disk rotates at high speed, such as 7200 rpm ( 120 revolutions per second. ) The rate at which data can be transferred from the disk to the computer is composed of several steps:
    - The *positioning time*, a.k.a. the *seek time* or *random access time* is the time required to move the heads from one cylinder to another, and for the heads to settle down after the move. This is typically the slowest step in the process and the predominant bottleneck to overall transfer rates.
    - The *rotational latency* is the amount of time required for the desired sector to rotate around and come under the read-write head.This can range anywhere from zero to one full revolution, and on the average will equal one-half revolution. This is another physical step and is usually the second slowest step behind seek time. ( For a disk rotating at 7200 rpm, the average rotational latency would be 1/2 revolution / 120 revolutions per second, or just over 4 milliseconds, a long time by computer standards.
    - The *transfer rate*, which is the time required to move the data electronically from the disk to the computer. ( Some authors may also use the term transfer rate to refer to the overall transfer rate, including seek time and rotational latency as well as the electronic data transfer rate. )
- Disk heads "fly" over the surface on a very thin cushion of air. If they should accidentally contact the disk, then a *head crash* occurs, which may or may not permanently damage the disk or even destroy it completely. For this reason it is normal to *park* the disk heads when turning a computer off, which means to move the heads off the disk or to an area of the disk where there is no data stored.
- Floppy disks are normally *removable*. Hard drives can also be removable, and some are even *hot-swappable*, meaning they can be removed while the computer is running, and a new hard drive inserted in their place.
- Disk drives are connected to the computer via a cable known as the *I/O Bus.* Some of the common interface formats include Enhanced Integrated Drive Electronics, EIDE; Advanced Technology Attachment, ATA; Serial ATA, SATA, Universal Serial Bus, USB; Fiber Channel, FC, and Small Computer Systems Interface, SCSI.
- The *host controller* is at the computer end of the I/O bus, and the *disk controller* is built into the disk itself. The CPU issues commands to the host controller via I/O ports. Data is transferred between the magnetic surface and onboard *cache* by the disk controller, and then the data is transferred from that cache to the host controller and the motherboard memory at electronic speeds.

### 10.1.2 Solid-State Disks - New

- As technologies improve and economics change, old technologies are often used in different ways. One example of this is the increasing used of *solid state disks, or SSDs.*
- SSDs use memory technology as a small fast hard disk. Specific implementations may use either flash memory or DRAM chips protected by a battery to sustain the information through power cycles.
- Because SSDs have no moving parts they are much faster than traditional hard drives, and certain problems such as the scheduling of disk accesses simply do not apply.
- However SSDs also have their weaknesses: They are more expensive than hard drives, generally not as large, and may have shorter life spans.
- SSDs are especially useful as a high-speed cache of hard-disk information that must be accessed quickly. One example is to store filesystem meta-data, e.g. directory and inode information, that must be accessed quickly and often. Another variation is a boot disk containing the OS and some application executables, but no vital user data. SSDs are also used in laptops to make them smaller, faster, and lighter.
- Because SSDs are so much faster than traditional hard disks, the throughput of the bus can become a limiting factor, causing some SSDs to be connected directly to the system PCI bus for example.

### 10.1.3 Magnetic Tapes - was 12.1.2

- Magnetic tapes were once used for common secondary storage before the days of hard disk drives, but today are used primarily for backups.
- Accessing a particular spot on a magnetic tape can be slow, but once reading or writing commences, access speeds are comparable to disk drives.
- Capacities of tape drives can range from 20 to 200 GB, and compression can double that capacity.

## 10.2 Disk Structure

- The traditional head-sector-cylinder, HSC numbers are mapped to linear block addresses by numbering the first sector on the first head on the outermost track as sector 0. Numbering proceeds with the rest of the sectors on that same track, and then the rest of the tracks on the same cylinder before proceeding through the rest of the cylinders to the center of the disk. In modern practice these linear block addresses are used in place of the HSC numbers for a variety of reasons:
    1. The linear length of tracks near the outer edge of the disk is much longer than for those tracks located near the center, and therefore it is possible to squeeze many more sectors onto outer tracks than onto inner ones.
    2. All disks have some bad sectors, and therefore disks maintain a few spare sectors that can be used in place of the bad ones. The mapping of spare sectors to bad sectors in managed internally to the disk controller.
    3. Modern hard drives can have thousands of cylinders, and hundreds of sectors per track on their outermost tracks. These numbers exceed the range of HSC numbers for many ( older ) operating systems, and therefore disks can be configured for any convenient combination of HSC values that falls within the total number of sectors physically on the drive.
- There is a limit to how closely packed individual bits can be placed on a physical media, but that limit is growing increasingly more packed as technological advances are made.
- Modern disks pack many more sectors into outer cylinders than inner ones, using one of two approaches:
    o With **Constant Linear Velocity, CLV,** the density of bits is uniform from cylinder to cylinder. Because there are more sectors in outer cylinders, the disk spins slower when reading those cylinders, causing the rate of bits passing under the read-write head to remain constant. This is the approach used by modern CDs and DVDs.
    o With **Constant Angular Velocity, CAV,** the disk rotates at a constant angular speed, with the bit density decreasing on outer cylinders. ( These disks would have a constant number of sectors per track on all cylinders. )

## 10.3 Disk Attachment

Disk drives can be attached either directly to a particular host ( a local disk ) or to a network.

### 10.3.1 Host-Attached Storage

- Local disks are accessed through I/O Ports as described earlier.
- The most common interfaces are IDE or ATA, each of which allow up to two drives per host controller.
- SATA is similar with simpler cabling.
- High end workstations or other systems in need of larger number of disks typically use SCSI disks:
    o The SCSI standard supports up to 16 **targets** on each SCSI bus, one of which is generally the host adapter and the other 15 of which can be disk or tape drives.
    o A SCSI target is usually a single drive, but the standard also supports up to 8 **units** within each target. These would generally be used for accessing individual disks within a RAID array. ( See below. )
    o The SCSI standard also supports multiple host adapters in a single computer, i.e. multiple SCSI busses.
    o Modern advancements in SCSI include "fast" and "wide" versions, as well as SCSI-2.
    o SCSI cables may be either 50 or 68 conductors. SCSI devices may be external as well as internal.
    o See wikipedia for more information on the SCSI interface.
- FC is a high-speed serial architecture that can operate over optical fiber or four-conductor copper wires, and has two variants:
    o A large switched fabric having a 24-bit address space. This variant allows for multiple devices and multiple hosts to interconnect, forming the basis for the **storage-area networks, SANs,** to be discussed in a future section.
    o The **arbitrated loop, FC-AL,** that can address up to 126 devices ( drives and controllers. )

- Network attached storage connects storage devices to computers using a remote procedure call, RPC, interface, typically with something like NFS filesystem mounts. This is convenient for allowing several computers in a group common access and naming conventions for shared storage.
- NAS can be implemented using SCSI cabling, or *ISCSI* uses Internet protocols and standard network connections, allowing long-distance remote access to shared files.
- NAS allows computers to easily share data storage, but tends to be less efficient than standard host-attached storage.



**Figure 10.2 - Network-attached storage.**

## 10.3.3 Storage-Area Network

- A *Storage-Area Network, SAN,* connects computers and storage devices in a network, using storage protocols instead of network protocols.
- One advantage of this is that storage access does not tie up regular networking bandwidth.
- SAN is very flexible and dynamic, allowing hosts and devices to attach and detach on the fly.
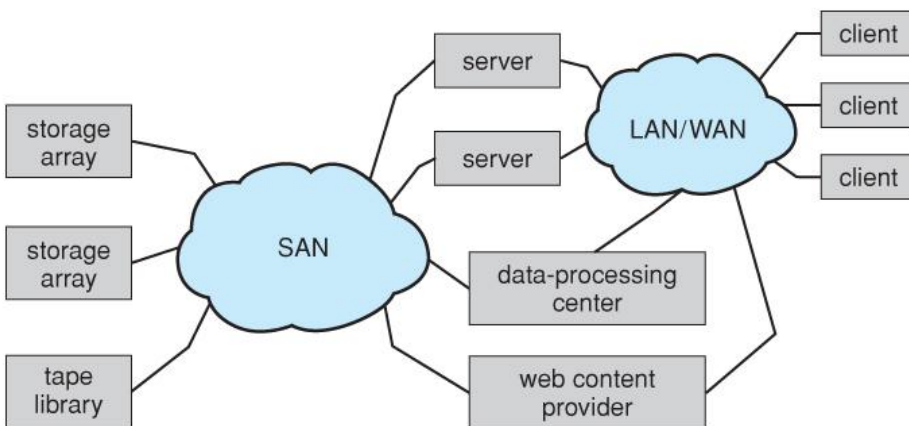- SAN is also controllable, allowing restricted access to certain hosts and devices.



**Figure 10.3 - Storage-area network.**

## 10.4 Disk Scheduling

- As mentioned earlier, disk transfer speeds are limited primarily by *seek times* and *rotational latency.* When multiple requests are to be processed there is also some inherent delay in waiting for other requests to be processed.
- *Bandwidth* is measured by the amount of data transferred divided by the total amount of time from the first request being made to the last transfer being completed, ( for a series of disk requests. )
- Both bandwidth and access time can be improved by processing requests in a good order.

- Disk requests include the disk address, memory address, number of sectors to transfer, and whether the request is for reading or writing.

### 10.4.1 FCFS Scheduling

- *First-Come First-Serve* is simple and intrinsically fair, but not very efficient. Consider in the following sequence the wild swing from cylinder 122 to 14 and then back to 124:
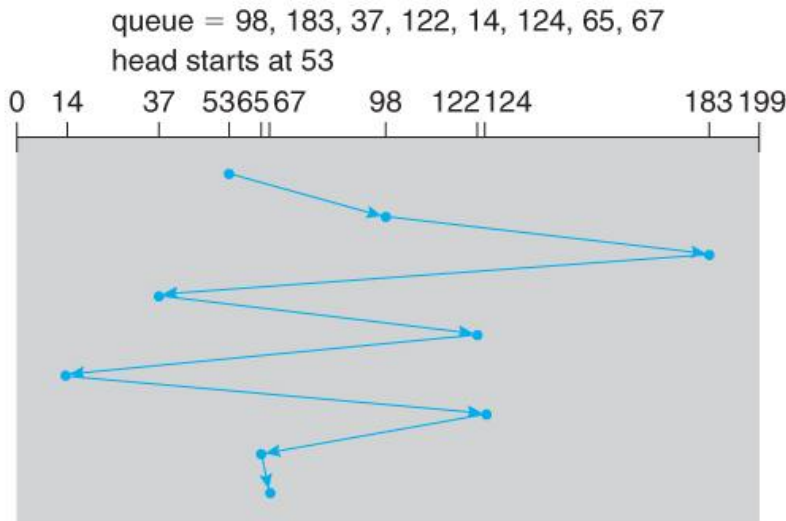


**Figure 10.4 - FCFS disk scheduling.**

### 10.4.2 SSTF Scheduling

- *Shortest Seek Time First* scheduling is more efficient, but may lead to starvation if a constant stream of requests arrives for the same general area of the disk.
- SSTF reduces the total head movement to 236 cylinders, down from 640 required for the same set of requests under FCFS. Note, however that the distance could be reduced still further to 208 by starting with 37 and then 14 first before processing the rest of the requests.
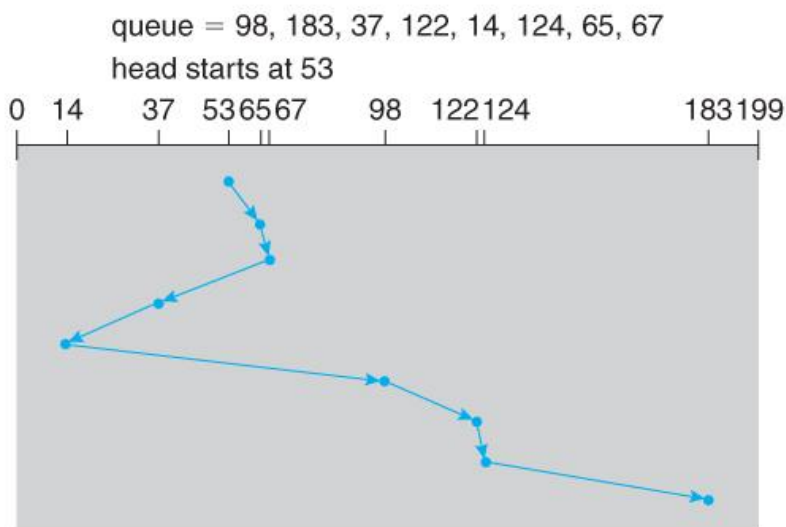


**Figure 10.5 - SSTF disk scheduling.**

### 10.4.3 SCAN Scheduling

- The *SCAN* algorithm, a.k.a. the *elevator* algorithm moves back and forth from one end of the disk to the other, similarly to an elevator processing requests in a tall building.
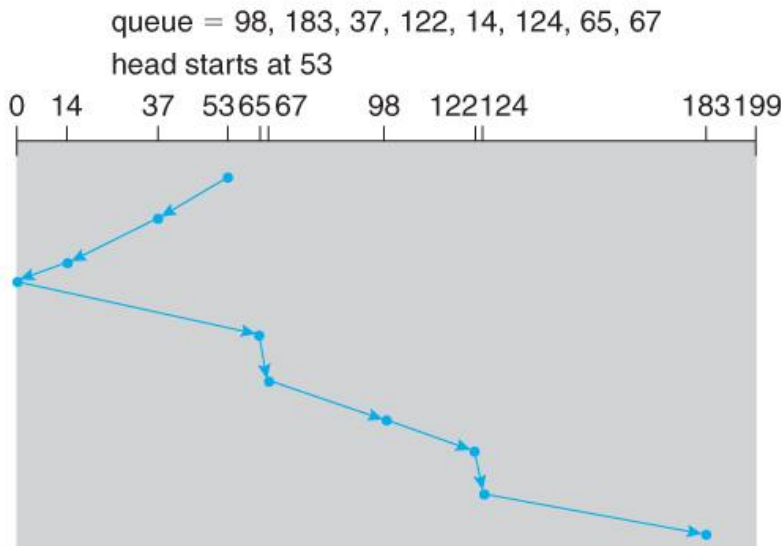


**Figure 10.6 - SCAN disk scheduling.**

- Under the SCAN algorithm, If a request arrives just ahead of the moving head then it will be processed right away, but if it arrives just after the head has passed, then it will have to wait for the head to pass going the other way on the return trip. This leads to a fairly wide variation in access times which can be improved upon.
- Consider, for example, when the head reaches the high end of the disk: Requests with high cylinder numbers just missed the passing head, which means they are all fairly recent requests, whereas requests with low numbers may have been waiting for a much longer time. Making the return scan from high to low then ends up accessing recent requests first and making older requests wait that much longer.

### 10.4.4 C-SCAN Scheduling

- The *Circular-SCAN* algorithm improves upon SCAN by treating all requests in a circular queue fashion - Once the head reaches the end of the disk, it returns to the other end without processing any requests, and then starts again from the beginning of the disk:
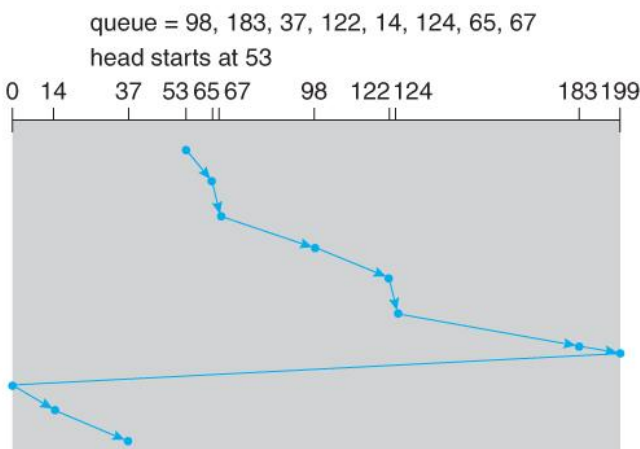


**Figure 10.7 - C-SCAN disk scheduling.**

- **LOOK** scheduling improves upon SCAN by looking ahead at the queue of pending requests, and not moving the heads any farther towards the end of the disk than is necessary. The following diagram illustrates the circular form of LOOK:
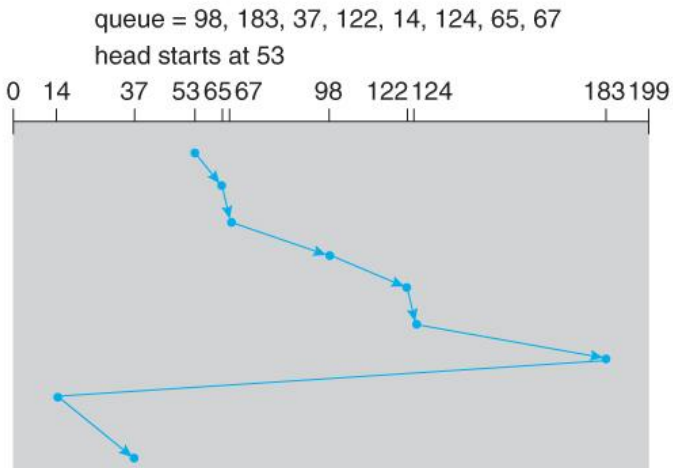


**Figure 10.8 - C-LOOK disk scheduling.**

*10.4.6 Selection of a Disk-Scheduling Algorithm*

- With very low loads all algorithms are equal, since there will normally only be one request to process at a time.
- For slightly larger loads, SSTF offers better performance than FCFS, but may lead to starvation when loads become heavy enough.
- For busier systems, SCAN and LOOK algorithms eliminate starvation problems.
- The actual optimal algorithm may be something even more complex than those discussed here, but the incremental improvements are generally not worth the additional overhead.
- Some improvement to overall filesystem access times can be made by intelligent placement of directory and/or inode information. If those structures are placed in the middle of the disk instead of at the beginning of the disk, then the maximum distance from those structures to data blocks is reduced to only one-half of the disk size. If those structures can be further distributed and furthermore have their data blocks stored as close as possible to the corresponding directory structures, then that reduces still further the overall time to find the disk block numbers and then access the corresponding data blocks.
- On modern disks the rotational latency can be almost as significant as the seek time, however it is not within the OSes control to account for that, because modern disks do not reveal their internal sector mapping schemes, ( particularly when bad blocks have been remapped to spare sectors. )
  - Some disk manufacturers provide for disk scheduling algorithms directly on their disk controllers, ( which do know the actual geometry of the disk as well as any remapping ), so that if a series of requests are sent from the computer to the controller then those requests can be processed in an optimal order.
  - Unfortunately there are some considerations that the OS must take into account that are beyond the abilities of the on-board disk-scheduling algorithms, such as priorities of some requests over others, or the need to process certain requests in a particular order. For this reason OSes may elect to spoon-feed requests to the disk controller one at a time in certain situations.

# 10.5 Disk Management

*105.1 Disk Formatting*

- Before a disk can be used, it has to be *low-level formatted*, which means laying down all of the headers and trailers marking the beginning and ends of each sector. Included in the header and trailer are the linear

sector numbers, and *error-correcting codes, ECC,* which allow damaged sectors to not only be detected, but in many cases for the damaged data to be recovered ( depending on the extent of the damage. ) Sector sizes are traditionally 512 bytes, but may be larger, particularly in larger drives.

- ECC calculation is performed with every disk read or write, and if damage is detected but the data is recoverable, then a *soft error* has occurred. Soft errors are generally handled by the on-board disk controller, and never seen by the OS. ( See below. )
- Once the disk is low-level formatted, the next step is to partition the drive into one or more separate partitions. This step must be completed even if the disk is to be used as a single large partition, so that the partition table can be written to the beginning of the disk.
- After partitioning, then the filesystems must be *logically formatted,* which involves laying down the master directory information ( FAT table or inode structure ), initializing free lists, and creating at least the root directory of the filesystem. ( Disk partitions which are to be used as raw devices are not logically formatted. This saves the overhead and disk space of the filesystem structure, but requires that the application program manage its own disk storage requirements. )

## 10.5.2 Boot Block

- Computer ROM contains a *bootstrap* program ( OS independent ) with just enough code to find the first sector on the first hard drive on the first controller, load that sector into memory, and transfer control over to it. ( The ROM bootstrap program may look in floppy and/or CD drives before accessing the hard drive, and is smart enough to recognize whether it has found valid boot code or not. )
- The first sector on the hard drive is known as the *Master Boot Record, MBR,* and contains a very small amount of code in addition to the *partition table.* The partition table documents how the disk is partitioned into logical disks, and indicates specifically which partition is the *active* or *boot* partition.
- The boot program then looks to the active partition to find an operating system, possibly loading up a slightly larger / more advanced boot program along the way.
- In a *dual-boot* ( or larger multi-boot ) system, the user may be given a choice of which operating system to boot, with a default action to be taken in the event of no response within some time frame.
- Once the kernel is found by the boot program, it is loaded into memory and then control is transferred over to the OS. The kernel will normally continue the boot process by initializing all important kernel data structures, launching important system services ( e.g. network daemons, sched, init, etc. ), and finally providing one or more login prompts. Boot options at this stage may include *single-user* a.k.a. *maintenance* or *safe* modes, in which very few system services are started - These modes are designed for system administrators to repair problems or otherwise maintain the system.
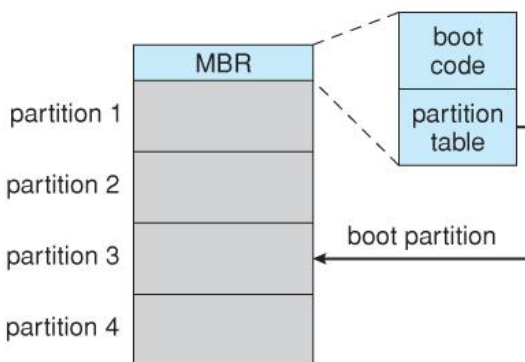


**Figure 10.9 - Booting from disk in Windows 2000.**

## 10.5.3 Bad Blocks

- No disk can be manufactured to 100% perfection, and all physical objects wear out over time. For these reasons all disks are shipped with a few bad blocks, and additional blocks can be expected to go bad slowly over time. If a large number of blocks go bad then the entire disk will need to be replaced, but a few here and there can be handled through other means.
- In the old days, bad blocks had to be checked for manually. Formatting of the disk or running certain disk-analysis tools would identify bad blocks, and attempt to read the data off of them one last time through

repeated tries. Then the bad blocks would be mapped out and taken out of future service. Sometimes the data could be recovered, and sometimes it was lost forever. ( Disk analysis tools could be either destructive or non-destructive. )

- Modern disk controllers make much better use of the error-correcting codes, so that bad blocks can be detected earlier and the data usually recovered. ( Recall that blocks are tested with every write as well as with every read, so often errors can be detected before the write operation is complete, and the data simply written to a different sector instead. )
- Note that re-mapping of sectors from their normal linear progression can throw off the disk scheduling optimization of the OS, especially if the replacement sector is physically far away from the sector it is replacing. For this reason most disks normally keep a few spare sectors on each cylinder, as well as at least one spare cylinder. Whenever possible a bad sector will be mapped to another sector on the same cylinder, or at least a cylinder as close as possible. *Sector slipping* may also be performed, in which all sectors between the bad sector and the replacement sector are moved down by one, so that the linear progression of sector numbers can be maintained.
- If the data on a bad block cannot be recovered, then a ***hard error*** has occurred., which requires replacing the file(s) from backups, or rebuilding them from scratch.

## 10.6 Swap-Space Management

- Modern systems typically swap out pages as needed, rather than swapping out entire processes. Hence the swapping system is part of the virtual memory management system.
- Managing swap space is obviously an important task for modern OSes.

### 10.6.1 Swap-Space Use

- The amount of swap space needed by an OS varies greatly according to how it is used. Some systems require an amount equal to physical RAM; some want a multiple of that; some want an amount equal to the amount by which virtual memory exceeds physical RAM, and some systems use little or none at all!
- Some systems support multiple swap spaces on separate disks in order to speed up the virtual memory system.

### 10.6.2 Swap-Space Location

Swap space can be physically located in one of two locations:

- As a large file which is part of the regular filesystem. This is easy to implement, but inefficient. Not only must the swap space be accessed through the directory system, the file is also subject to fragmentation issues. Caching the block location helps in finding the physical blocks, but that is not a complete fix.
- As a raw partition, possibly on a separate or little-used disk. This allows the OS more control over swap space management, which is usually faster and more efficient. Fragmentation of swap space is generally not a big issue, as the space is re-initialized every time the system is rebooted. The downside of keeping swap space on a raw partition is that it can only be grown by repartitioning the hard drive.

### 12.6.3 Swap-Space Management: An Example

- Historically OSes swapped out entire processes as needed. Modern systems swap out only individual pages, and only as needed. ( For example process code blocks and other blocks that have not been changed since they were originally loaded are normally just freed from the virtual memory system rather than copying them to swap space, because it is faster to go find them again in the filesystem and read them back in from there than to write them out to swap space and then read them back. )
- In the mapping system shown below for Linux systems, a map of swap space is kept in memory, where each entry corresponds to a 4K block in the swap space. Zeros indicate free slots and non-zeros refer to how many processes have a mapping to that particular block ( >1 for shared pages only. )
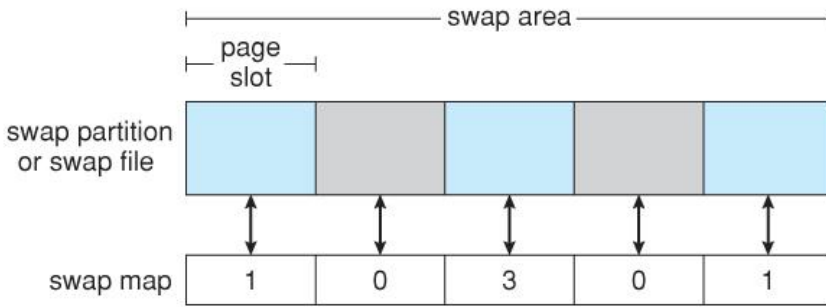
**Figure 10.10 - The data structures for swapping on Linux systems.**

## 10.7 RAID Structure

- The general idea behind RAID is to employ a group of hard drives together with some form of duplication, either to increase reliability or to speed up operations, ( or sometimes both. )
- **RAID** originally stood for **Redundant Array of Inexpensive Disks**, and was designed to use a bunch of cheap small disks in place of one or two larger more expensive ones. Today RAID systems employ large possibly expensive disks as their components, switching the definition to **Independent** disks.

### 10.7.1 Improvement of Reliability via Redundancy

- The more disks a system has, the greater the likelihood that one of them will go bad at any given time. Hence increasing disks on a system actually *decreases* the **Mean Time To Failure, MTTF** of the system.
- If, however, the same data was copied onto multiple disks, then the data would not be lost unless **both** ( or all ) copies of the data were damaged simultaneously, which is a **MUCH** lower probability than for a single disk going bad. More specifically, the second disk would have to go bad before the first disk was repaired, which brings the **Mean Time To Repair** into play. For example if two disks were involved, each with a MTTF of 100,000 hours and a MTTR of 10 hours, then the **Mean Time to Data Loss** would be 500 * 10^6 hours, or 57,000 years!
- This is the basic idea behind disk *mirroring*, in which a system contains identical data on two or more disks.
  - Note that a power failure during a write operation could cause both disks to contain corrupt data, if both disks were writing simultaneously at the time of the power failure. One solution is to write to the two disks in series, so that they will not both become corrupted ( at least not in the same way ) by a power failure. And alternate solution involves non-volatile RAM as a write cache, which is not lost in the event of a power failure and which is protected by error-correcting codes.

### 10.7.2 Improvement in Performance via Parallelism

- There is also a performance benefit to mirroring, particularly with respect to reads. Since every block of data is duplicated on multiple disks, read operations can be satisfied from any available copy, and multiple disks can be reading different data blocks simultaneously in parallel. ( Writes could possibly be sped up as well through careful scheduling algorithms, but it would be complicated in practice. )
- Another way of improving disk access time is with *striping*, which basically means spreading data out across multiple disks that can be accessed simultaneously.
  - With *bit-level striping* the bits of each byte are striped across multiple disks. For example if 8 disks were involved, then each 8-bit byte would be read in parallel by 8 heads on separate disks. A single disk read would access 8 * 512 bytes = 4K worth of data in the time normally required to read 512 bytes. Similarly if 4 disks were involved, then two bits of each byte could be stored on each disk, for 2K worth of disk access per read or write operation.
  - *Block-level striping* spreads a filesystem across multiple disks on a block-by-block basis, so if block N were located on disk 0, then block N + 1 would be on disk 1, and so on. This is particularly useful when filesystems are accessed in *clusters* of physical blocks. Other striping possibilities exist, with block-level striping being the most common.

### 10.7.3 RAID Levels

- Mirroring provides reliability but is expensive; Striping improves performance, but does not improve reliability. Accordingly there are a number of different schemes that combine the principals of mirroring and striping in different ways, in order to balance reliability versus performance versus cost. These are described by different *RAID levels*, as follows: ( In the diagram that follows, "C" indicates a copy, and "P" indicates parity, i.e. checksum bits. )
    1. *Raid Level 0 -* This level includes striping only, with no mirroring.
    2. *Raid Level 1 -* This level includes mirroring only, no striping.
    3. *Raid Level 2 -* This level stores error-correcting codes on additional disks, allowing for any damaged data to be reconstructed by subtraction from the remaining undamaged data. Note that this scheme requires only three extra disks to protect 4 disks worth of data, as opposed to full mirroring. ( The number of disks required is a function of the error-correcting algorithms, and the means by which the particular bad bit(s) is(are) identified. )
    4. *Raid Level 3 -* This level is similar to level 2, except that it takes advantage of the fact that each disk is still doing its own error-detection, so that when an error occurs, there is no question about which disk in the array has the bad data. As a result a single parity bit is all that is needed to recover the lost data from an array of disks. Level 3 also includes striping, which improves performance. The downside with the parity approach is that every disk must take part in every disk access, and the parity bits must be constantly calculated and checked, reducing performance. Hardware-level parity calculations and NVRAM cache can help with both of those issues. In practice level 3 is greatly preferred over level 2.
    5. *Raid Level 4 -* This level is similar to level 3, employing block-level striping instead of bit-level striping. The benefits are that multiple blocks can be read independently, and changes to a block only require writing two blocks ( data and parity ) rather than involving all disks. Note that new disks can be added seamlessly to the system provided they are initialized to all zeros, as this does not affect the parity results.
    6. *Raid Level 5 -* This level is similar to level 4, except the parity blocks are distributed over all disks, thereby more evenly balancing the load on the system. For any given block on the disk(s), one of the disks will hold the parity information for that block and the other N-1 disks will hold the data. Note that the same disk cannot hold both data and parity for the same block, as both would be lost in the event of a disk crash.
    7. *Raid Level 6 -* This level extends raid level 5 by storing multiple bits of error-recovery codes, ( such as the *Reed-Solomon codes* ), for each bit position of data, rather than a single parity bit. In the example shown below 2 bits of ECC are stored for every 4 bits of data, allowing data recovery in the face of up to two simultaneous disk failures. Note that this still involves only 50% increase in storage needs, as opposed to 100% for simple mirroring which could only tolerate a single disk failure.
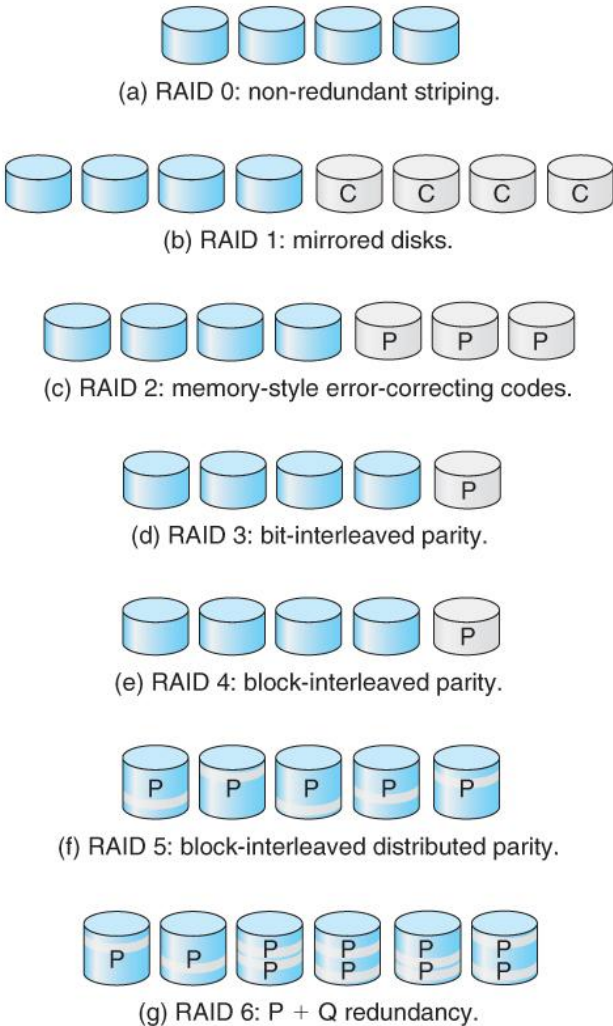
**Figure 10.11 - RAID levels.**

- There are also two RAID levels which combine RAID levels 0 and 1 ( striping and mirroring ) in different combinations, designed to provide both performance and reliability at the expense of increased cost.
    - **RAID level 0 + 1** disks are first striped, and then the striped disks mirrored to another set. This level generally provides better performance than RAID level 5.
    - **RAID level 1 + 0** mirrors disks in pairs, and then stripes the mirrored pairs. The storage capacity, performance, etc. are all the same, but there is an advantage to this approach in the event of multiple disk failures, as illustrated below:.
        - In diagram (a) below, the 8 disks have been divided into two sets of four, each of which is striped, and then one stripe set is used to mirror the other set.
            - If a single disk fails, it wipes out the entire stripe set, but the system can keep on functioning using the remaining set.
            - However if a second disk from the other stripe set now fails, then the entire system is lost, as a result of two disk failures.
        - In diagram (b), the same 8 disks are divided into four sets of two, each of which is mirrored, and then the file system is striped across the four sets of mirrored disks.
            - If a single disk fails, then that mirror set is reduced to a single disk, but the system rolls on, and the other three mirror sets continue mirroring.
            - Now if a second disk fails, ( that is not the mirror of the already failed disk ), then another one of the mirror sets is reduced to a single disk, but the system can continue without data loss.
            - In fact the second arrangement could handle as many as four simultaneously failed disks, as long as no two of them were from the same mirror pair.
    - See the wikipedia article on nested raid levels for more information.

202

a) RAID 0 + 1 with a single disk failure.



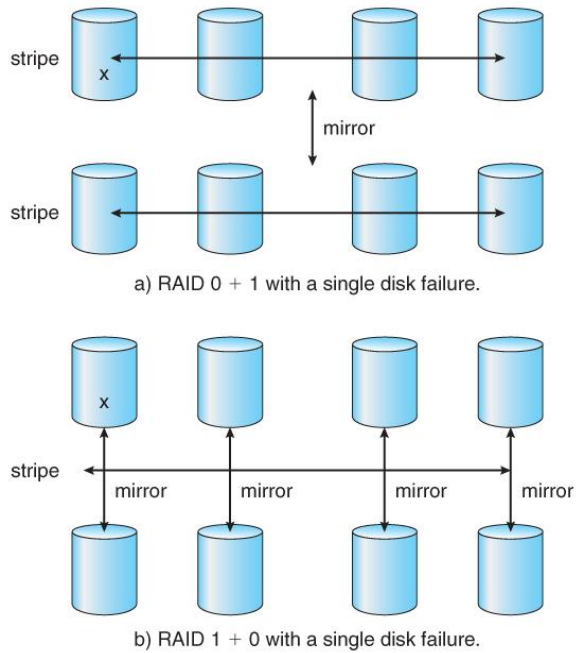b) RAID 1 + 0 with a single disk failure.

**Figure 10.12 - RAID 0 + 1 and 1 + 0**

## 10.7.4 Selecting a RAID Level

- Trade-offs in selecting the optimal RAID level for a particular application include cost, volume of data, need for reliability, need for performance, and rebuild time, the latter of which can affect the likelihood that a second disk will fail while the first failed disk is being rebuilt.
- Other decisions include how many disks are involved in a RAID set and how many disks to protect with a single parity bit. More disks in the set increases performance but increases cost. Protecting more disks per parity bit saves cost, but increases the likelihood that a second disk will fail before the first bad disk is repaired.

## 10.7.5 Extensions

- RAID concepts have been extended to tape drives ( e.g. striping tapes for faster backups or parity checking tapes for reliability ), and for broadcasting of data.

## 10.7.6 Problems with RAID

- RAID protects against physical errors, but not against any number of bugs or other errors that could write erroneous data.
- ZFS adds an extra level of protection by including data block checksums in all inodes along with the pointers to the data blocks. If data are mirrored and one copy has the correct checksum and the other does not, then the data with the bad checksum will be replaced with a copy of the data with the good checksum. This increases reliability greatly over RAID alone, at a cost of a performance hit that is acceptable because ZFS is so fast to begin with.
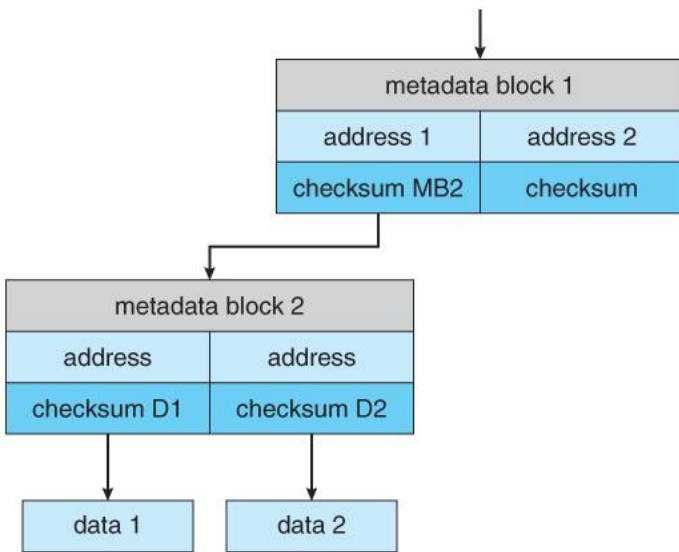
**Figure 10.13 - ZFS checksums all metadata and data.**

- Another problem with traditional filesystems is that the sizes are fixed, and relatively difficult to change. Where RAID sets are involved it becomes even harder to adjust filesystem sizes, because a filesystem cannot span across multiple filesystems.
- ZFS solves these problems by pooling RAID sets, and by dynamically allocating space to filesystems as needed. Filesystem sizes can be limited by quotas, and space can also be reserved to guarantee that a filesystem will be able to grow later, but these parameters can be changed at any time by the filesystem's owner. Otherwise filesystems grow and shrink dynamically as needed.
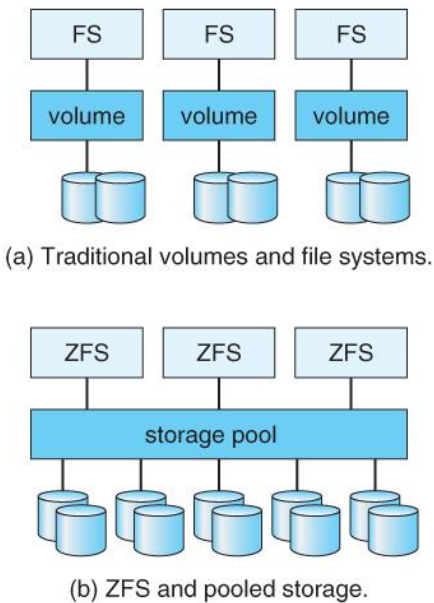


(a) Traditional volumes and file systems.



(b) ZFS and pooled storage.

**Figure 10.14 - (a) Traditional volumes and file systems. (b) a ZFS pool and file systems.**

## 10.8 Stable-Storage Implementation ( Optional )

- The concept of stable storage ( first presented in chapter 6 ) involves a storage medium in which data is *never* lost, even in the face of equipment failure in the middle of a write operation.
- To implement this requires two ( or more ) copies of the data, with separate failure modes.
- An attempted disk write results in one of three possible outcomes:
    1. The data is successfully and completely written.
    2. The data is partially written, but not completely. The last block written may be garbled.
    3. No writing takes place at all.

- Whenever an equipment failure occurs during a write, the system must detect it, and return the system back to a consistent state. To do this requires two physical blocks for every logical block, and the following procedure:
    1. Write the data to the first physical block.
    2. After step 1 had completed, then write the data to the second physical block.
    3. Declare the operation complete only after both physical writes have completed successfully.
- During recovery the pair of blocks is examined.
    o If both blocks are identical and there is no sign of damage, then no further action is necessary.
    o If one block contains a detectable error but the other does not, then the damaged block is replaced with the good copy. ( This will either undo the operation or complete the operation, depending on which block is damaged and which is undamaged. )
    o If neither block shows damage but the data in the blocks differ, then replace the data in the first block with the data in the second block. ( Undo the operation. )

Because the sequence of operations described above is slow, stable storage usually includes NVRAM as a cache, and declares a write operation complete once it has been written to the NVRAM.

## 10.9 Summary

---

## Was 12.9 Tertiary-Storage Structure - Optional, Omitted from Ninth Edition

- Primary storage refers to computer memory chips; Secondary storage refers to fixed-disk storage systems ( hard drives ); And *Tertiary Storage* refers to *removable media,* such as tape drives, CDs, DVDs, and to a lesser extend floppies, thumb drives, and other detachable devices.
- Tertiary storage is typically characterized by large capacity, low cost per MB, and slow access times, although there are exceptions in any of these categories.
- Tertiary storage is typically used for backups and for long-term archival storage of completed work. Another common use for tertiary storage is to swap large little-used files ( or groups of files ) off of the hard drive, and then swap them back in as needed in a fashion similar to secondary storage providing swap space for primary storage. ( Review The Paging Game, note 5 ).

### 12.9.1 Tertiary-Storage Devices

### 12.9.1.1 Removable Disks

- Removable magnetic disks ( e.g. floppies ) can be nearly as fast as hard drives, but are at greater risk for damage due to scratches. Variations of removable magnetic disks up to a GB or more in capacity have been developed. ( Hot-swappable hard drives? )
- A *magneto-optical* disk uses a magnetic disk covered in a clear plastic coating that protects the surface.
    o The heads sit a considerable distance away from the magnetic surface, and as a result do not have enough magnetic strength to switch bits *at normal room temperature.*
    o For writing, a laser is used to heat up a specific spot on the disk, to a temperature at which the weak magnetic field of the write head is able to flip the bits.
    o For reading, a laser is shined at the disk, and the *Kerr effect* causes the polarization of the light to become rotated either clockwise or counter-clockwise depending on the orientation of the magnetic field.
- *Optical disks* do not use magnetism at all, but instead use special materials that can be altered ( by lasers ) to have relatively light or dark spots.
    o For example the *phase-change disk* has a material that can be frozen into either a crystalline or an amorphous state, the latter of which is less transparent and reflects less light when a laser is bounced off a reflective surface under the material.
        ▪ Three powers of lasers are used with phase-change disks: (1) a low power laser is used to read the disk, without effecting the materials. (2) A medium power erases the disk, by melting and re-freezing the medium into a crystalline state, and (3) a high power writes to the disk by melting the medium and re-freezing it into the amorphous state.

- - The most common examples of these disks are *re-writable* CD-RWs and DVD-RWs.
  - An alternative to the disks described above are *Write-Once Read-Many, WORM* drives.
    - o The original version of WORM drives involved a thin layer of aluminum sandwiched between two protective layers of glass or plastic.
      - Holes were burned in the aluminum to write bits.
      - Because the holes could not be filled back in, there was no way to re-write to the disk. ( Although data could be erased by burning more holes. )
    - o WORM drives have important legal ramifications for data that must be stored for a very long time and must be provable in court as unaltered since it was originally written. ( Such as long-term storage of medical records. )
    - o Modern CD-R and DVD-R disks are examples of WORM drives that use organic polymer inks instead of an aluminum layer.
  - Read-only disks are similar to WORM disks, except the bits are pressed onto the disk at the factory, rather than being burned on one by one. ( See http://en.wikipedia.org/wiki/CD_manufacturing#Premastering for more information on CD manufacturing techniques. )

### 12.9.1.2 Tapes

- Tape drives typically cost more than disk drives, but the cost per MB of the tapes themselves is lower.
- Tapes are typically used today for backups, and for enormous volumes of data stored by certain scientific establishments. ( E.g. NASA's archive of space probe and satellite imagery, which is currently being downloaded from numerous sources faster than anyone can actually look at it. )
- Robotic tape changers move tapes from drives to archival tape libraries upon demand.
- ( Never underestimate the bandwidth of a station wagon full of tapes rolling down the highway! )

### 12.9.1.3 Future Technology

- *Solid State Disks, SSDs,* are becoming more and more popular.
- *Holographic storage* uses laser light to store images in a 3-D structure, and the entire data structure can be transferred in a single flash of laser light.
- *Micro-Electronic Mechanical Systems, MEMS,* employs the technology used for computer chip fabrication to create VERY tiny little machines. One example packs 10,000 read-write heads within a square centimeter of space, and as media are passed over it, all 10,000 heads can read data in parallel.

### 12.9.2 Operating-System Support

- The OS must provide support for tertiary storage as removable media, including the support to transfer data between different systems.

### 12.9.2.1 Application Interface

- File systems are typically not stored on tapes. ( It might be technically possible, but it is impractical. )
- Tapes are also not low-level formatted, and do not use fixed-length blocks. Rather data is written to tapes in variable length blocks as needed.
- Tapes are normally accessed as raw devices, requiring each application to determine how the data is to be stored and read back. Issues such as header contents and ASCII versus binary encoding ( and byte-ordering ) are generally application specific.
- Basic operations supported for tapes include locate( ), read( ), write( ), and read_position( ).
- ( Because of variable length writes ), writing to a tape erases all data that follows that point on the tape.
  - o Writing to a tape places the End of Tape ( EOT ) marker at the end of the data written.
  - o It is not possible to locate( ) to any spot past the EOT marker.

### 12.9.2.2 File Naming

- File naming conventions for removable media are not entirely uniquely specific, nor are they necessarily consistent between different systems. ( Two removable disks may contain files with the same name, and there is no clear way for the naming system to distinguish between them. )
- Fortunately music CDs have a common format, readable by all systems. Data CDs and DVDs have only a few format choices, making it easy for a system to support all known formats.

### 12.9.2.3 Hierarchical Storage Management

- Hierarchical storage involves extending file systems out onto tertiary storage, swapping files from hard drives to tapes in much the same manner as data blocks are swapped from memory to hard drives.
- A placeholder is generally left on the hard drive, storing information about the particular tape ( or other removable media ) on which the file has been swapped out to.
- A robotic system transfers data to and from tertiary storage as needed, generally automatically upon demand of the file(s) involved.

### 12.9.3 Performance Issues

### 12.9.3.1 Speed

- *Sustained Bandwidth* is the rate of data transfer during a large file transfer, once the proper tape is loaded and the file located.
- *Effective Bandwidth* is the effective overall rate of data transfer, including any overhead necessary to load the proper tape and find the file on the tape.
- *Access Latency* is all of the accumulated waiting time before a file can be actually read from tape. This includes the time it takes to find the file on the tape, the time to load the tape from the tape library, and the time spent waiting in the queue for the tape drive to become available.
- Clearly tertiary storage access is much slower than secondary access, although removable disks ( e.g. a CD jukebox ) have somewhat faster access than a tape library.
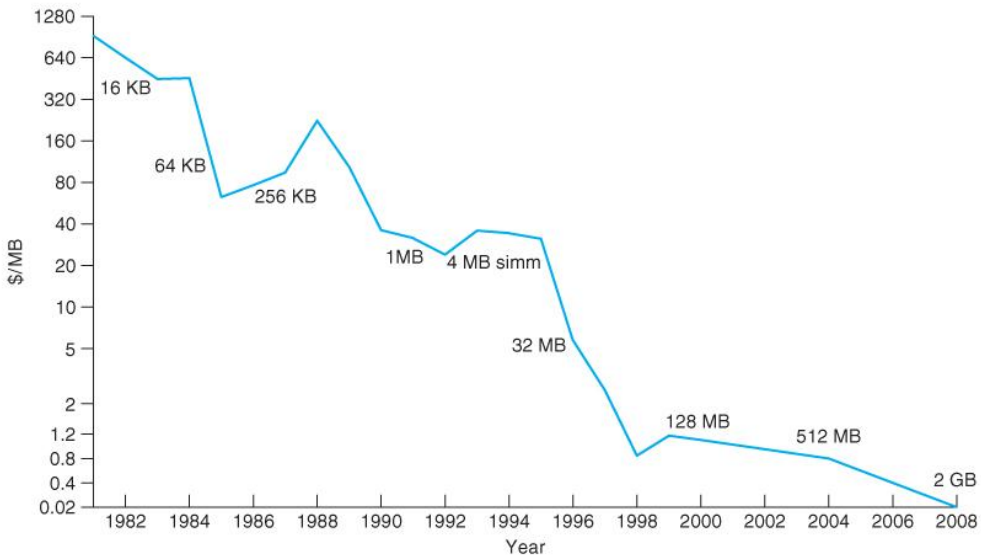
### 12.9.3.1 Reliability

- Fixed hard drives are generally more reliable than removable drives, because they are less susceptible to the environment.
- Optical disks are generally more reliable than magnetic media.
- A fixed hard drive crash can destroy all data, whereas an optical drive or tape drive failure will often not harm the data media, ( and certainly can't damage any media not in the drive at the time of the failure. )
- Tape drives are mechanical devices, and can wear out tapes over time, ( as the tape head is generally in much closer physical contact with the tape than disk heads are with platters. )
  - Some drives may only be able to read tapes a few times whereas other drives may be able to re-use the same tapes millions of times.
  - Backup tapes should be read after writing, to verify that the backup tape is readable. ( Unfortunately that may have been the LAST time that particular tape was readable, and the only way to be sure is to read it again, . . . )
  - Long-term tape storage can cause degradation, as magnetic fields "drift" from one layer of tape to the adjacent layers. Periodic fast-forwarding and rewinding of tapes can help, by changing which section of tape lays against which other layers.
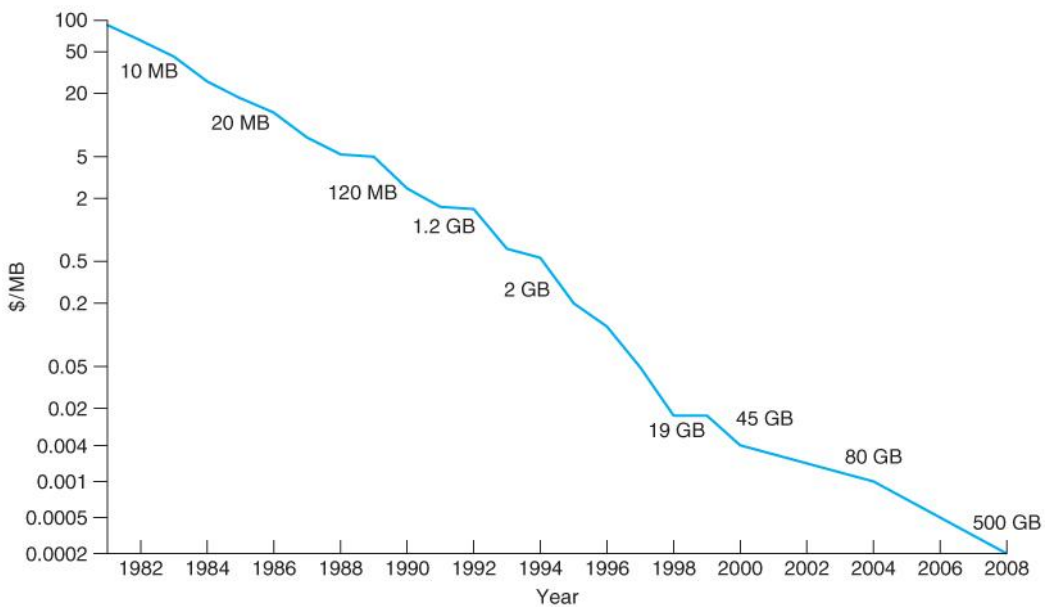
### 12.9.3.3 Cost

- The cost per megabyte for removable media is its strongest selling feature, particularly as the amount of storage involved ( i.e. the number of tapes, CDs, etc ) increases.
- However the cost per megabyte for hard drives has dropped more rapidly over the years than the cost of removable media, such that the currently most cost-effective backup solution for many systems is simply an additional ( external ) hard drive.
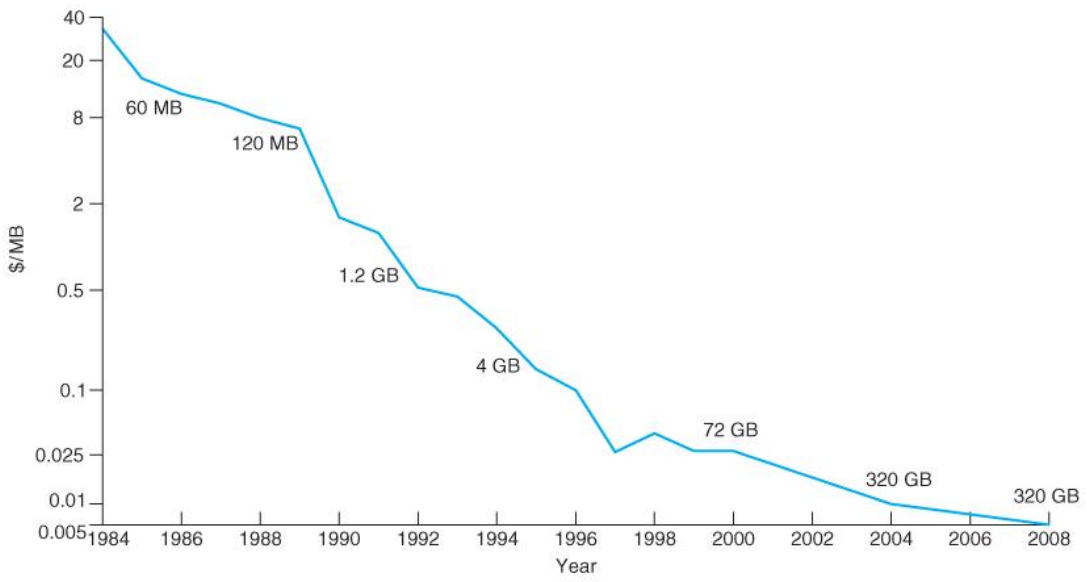
- ( One good use for old unwanted PCs is to put them on a network as a backup server and/or print server. The downside to this backup solution is that the backups are stored on-site with the original data, and a fire, flood, or burglary could wipe out both the original data and the backups. )



**Old Figure 12.15 - Price per megabyte of DRAM, from 1981 to 2008**



**Old Figure 12.16 - Price per megabyte of magnetic hard disk, from 1981 to 2008.**

**Old Figure 12.17 - Price per megabyte of a tape drive, from 1984 to 2008.**

# I/O Systems

**References:**

1. Abraham Silberschatz, Greg Gagne, and Peter Baer Galvin, "Operating System Concepts, Eighth Edition ", Chapter 13

## 13.1 Overview

- Management of I/O devices is a very important part of the operating system - so important and so varied that entire I/O subsystems are devoted to its operation. ( Consider the range of devices on a modern computer, from mice, keyboards, disk drives, display adapters, USB devices, network connections, audio I/O, printers, special devices for the handicapped, and many special-purpose peripherals. )
- I/O Subsystems must contend with two ( conflicting? ) trends: (1) The gravitation towards standard interfaces for a wide range of devices, making it easier to add newly developed devices to existing systems, and (2) the development of entirely new types of devices, for which the existing standard interfaces are not always easy to apply.
- *Device drivers* are modules that can be plugged into an OS to handle a particular device or category of similar devices.

## 13.2 I/O Hardware

- I/O devices can be roughly categorized as storage, communications, user-interface, and other
- Devices communicate with the computer via signals sent over wires or through the air.
- Devices connect with the computer via *ports*, e.g. a serial or parallel port.
- A common set of wires connecting multiple devices is termed a *bus.*
  - Buses include rigid protocols for the types of messages that can be sent across the bus and the procedures for resolving contention issues.
  - Figure 13.1 below illustrates three of the four bus types commonly found in a modern PC:
    1. The *PCI bus* connects high-speed high-bandwidth devices to the memory subsystem ( and the CPU. )
    2. The *expansion bus* connects slower low-bandwidth devices, which typically deliver data one character at a time ( with buffering. )
    3. The *SCSI bus* connects a number of SCSI devices to a common SCSI controller.
    4. A *daisy-chain bus,* ( not shown) is when a string of devices is connected to each other like beads on a chain, and only one of the devices is directly connected to the host.
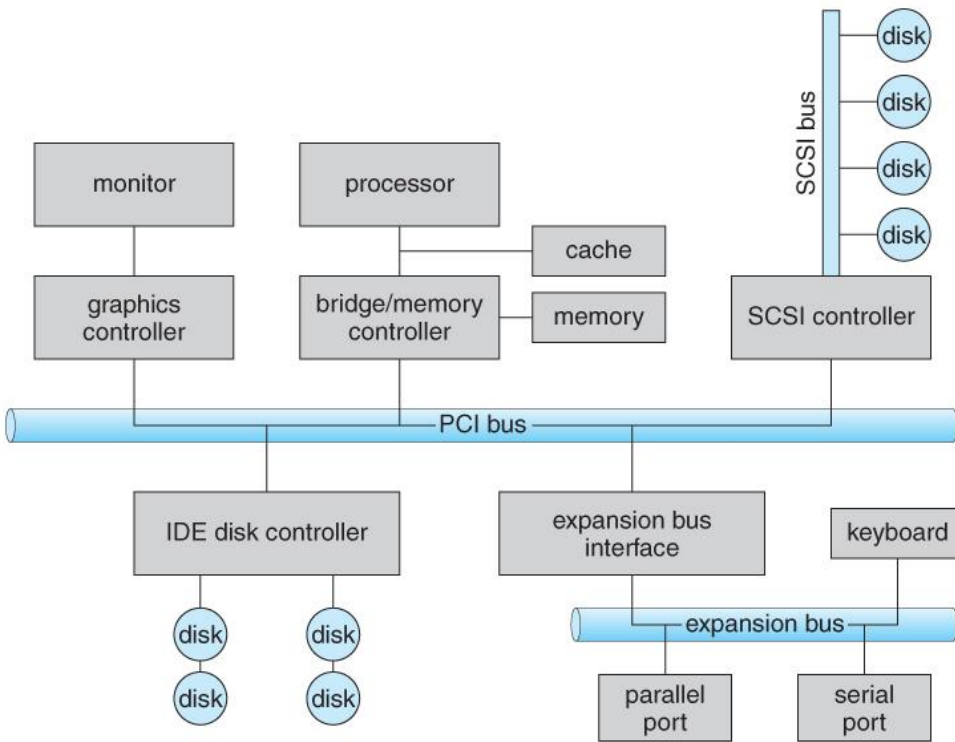
**Figure 13.1 - A typical PC bus structure.**

- One way of communicating with devices is through *registers* associated with each port. Registers may be one to four bytes in size, and may typically include ( a subset of ) the following four:
  1. The *data-in register* is read by the host to get input from the device.
  2. The *data-out register* is written by the host to send output.
  3. The *status register* has bits read by the host to ascertain the status of the device, such as idle, ready for input, busy, error, transaction complete, etc.
  4. The *control register* has bits written by the host to issue commands or to change settings of the device such as parity checking, word length, or full- versus half-duplex operation.
- Figure 13.2 shows some of the most common I/O port address ranges.

| I/O address range (hexadecimal) | device |
|---|---|
| 000–00F | DMA controller |
| 020–021 | interrupt controller |
| 040–043 | timer |
| 200–20F | game controller |
| 2F8–2FF | serial port (secondary) |
| 320–32F | hard-disk controller |
| 378–37F | parallel port |
| 3D0–3DF | graphics controller |
| 3F0–3F7 | diskette-drive controller |
| 3F8–3FF | serial port (primary) |

**Figure 13.2 - Device I/O port locations on PCs ( partial ).**

- Another technique for communicating with devices is *memory-mapped I/O.*
  o In this case a certain portion of the processor's address space is mapped to the device, and communications occur by reading and writing directly to/from those memory areas.

- Memory-mapped I/O is suitable for devices which must move large quantities of data quickly, such as graphics cards.
- Memory-mapped I/O can be used either instead of or more often in combination with traditional registers. For example, graphics cards still use registers for control information such as setting the video mode.
- A potential problem exists with memory-mapped I/O, if a process is allowed to write directly to the address space used by a memory-mapped I/O device.
- ( Note: Memory-mapped I/O is not the same thing as direct memory access, DMA. See section 13.2.3 below. )

## 13.2.1 Polling

- One simple means of device *handshaking* involves polling:
    1. The host repeatedly checks the *busy bit* on the device until it becomes clear.
    2. The host writes a byte of data into the data-out register, and sets the *write bit* in the command register ( in either order. )
    3. The host sets the *command ready bit* in the command register to notify the device of the pending command.
    4. When the device controller sees the command-ready bit set, it first sets the busy bit.
    5. Then the device controller reads the command register, sees the write bit set, reads the byte of data from the data-out register, and outputs the byte of data.
    6. The device controller then clears the *error bit* in the status register, the command-ready bit, and finally clears the busy bit, signaling the completion of the operation.
- Polling can be very fast and efficient, if both the device and the controller are fast and if there is significant data to transfer. It becomes inefficient, however, if the host must wait a long time in the busy loop waiting for the device, or if frequent checks need to be made for data that is infrequently there.

## 13.2.2 Interrupts

- Interrupts allow devices to notify the CPU when they have data to transfer or when an operation is complete, allowing the CPU to perform other duties when no I/O transfers need its immediate attention.
- The CPU has an *interrupt-request line* that is sensed after every instruction.
    - A device's controller *raises* an interrupt by asserting a signal on the interrupt request line.
    - The CPU then performs a state save, and transfers control to the *interrupt handler* routine at a fixed address in memory. ( The CPU *catches* the interrupt and *dispatches* the interrupt handler. )
    - The interrupt handler determines the cause of the interrupt, performs the necessary processing, performs a state restore, and executes a *return from interrupt* instruction to return control to the CPU. ( The interrupt handler *clears* the interrupt by servicing the device. )
        - ( Note that the state restored does not need to be the same state as the one that was saved when the interrupt went off. See below for an example involving time-slicing. )
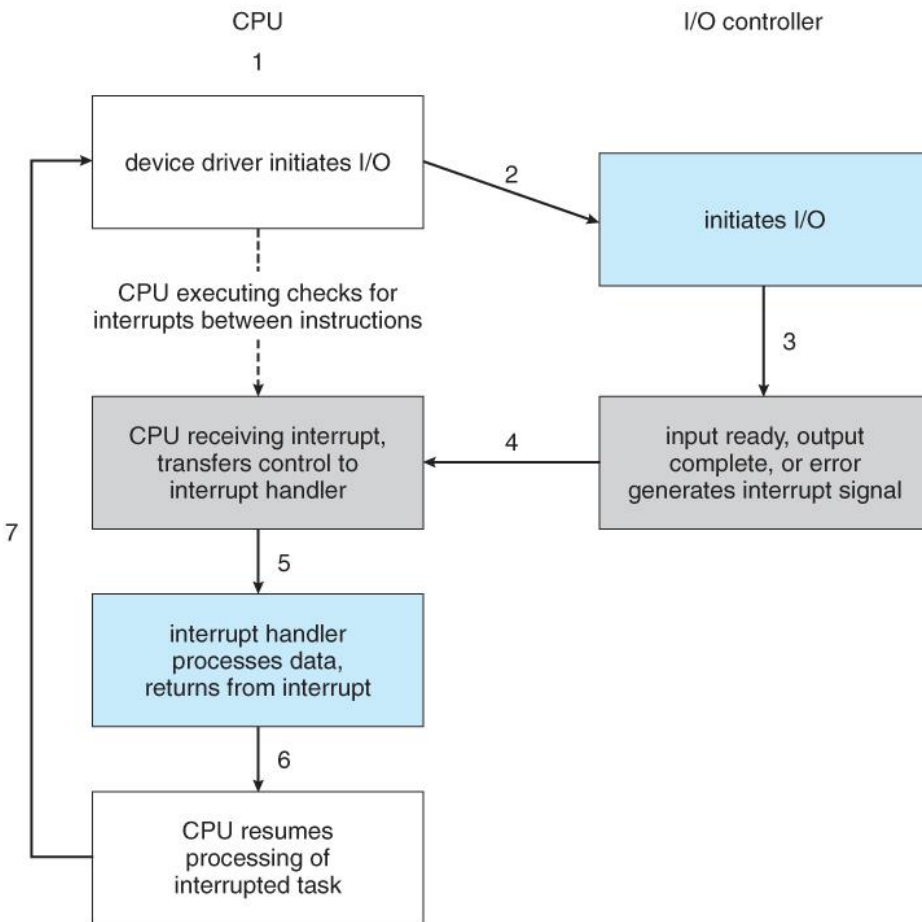- Figure 13.3 illustrates the interrupt-driven I/O procedure:

**Figure 13.3 - Interrupt-driven I/O cycle.**

- The above description is adequate for simple interrupt-driven I/O, but there are three needs in modern computing which complicate the picture:
  1. The need to defer interrupt handling during critical processing,
  2. The need to determine **which** interrupt handler to invoke, without having to poll all devices to see which one needs attention, and
  3. The need for multi-level interrupts, so the system can differentiate between high- and low-priority interrupts for proper response.
- These issues are handled in modern computer architectures with **interrupt-controller** hardware.
  - Most CPUs now have two interrupt-request lines: One that is **non-maskable** for critical error conditions and one that is **maskable,** that the CPU can temporarily ignore during critical processing.
  - The interrupt mechanism accepts an **address,** which is usually one of a small set of numbers for an offset into a table called the **interrupt vector.** This table ( usually located at physical address zero ? ) holds the addresses of routines prepared to process specific interrupts.
  - The number of possible interrupt handlers still exceeds the range of defined interrupt numbers, so multiple handlers can be **interrupt chained**. Effectively the addresses held in the interrupt vectors are the head pointers for linked-lists of interrupt handlers.
  - Figure 13.4 shows the Intel Pentium interrupt vector. Interrupts 0 to 31 are non-maskable and reserved for serious hardware and other errors. Maskable interrupts, including normal device I/O interrupts begin at interrupt 32.
  - Modern interrupt hardware also supports **interrupt priority levels**, allowing systems to mask off only lower-priority interrupts while servicing a high-priority

interrupt, or conversely to allow a high-priority signal to interrupt the processing of a low-priority one.

| vector number | description |
|---|---|
| 0 | divide error |
| 1 | debug exception |
| 2 | null interrupt |
| 3 | breakpoint |
| 4 | INTO-detected overflow |
| 5 | bound range exception |
| 6 | invalid opcode |
| 7 | device not available |
| 8 | double fault |
| 9 | coprocessor segment overrun (reserved) |
| 10 | invalid task state segment |
| 11 | segment not present |
| 12 | stack fault |
| 13 | general protection |
| 14 | page fault |
| 15 | (Intel reserved, do not use) |
| 16 | floating-point error |
| 17 | alignment check |
| 18 | machine check |
| 19–31 | (Intel reserved, do not use) |
| 32–255 | maskable interrupts |

**Figure 13.4 - Intel Pentium processor event-vector table.**

- At boot time the system determines which devices are present, and loads the appropriate handler addresses into the interrupt table.
- During operation, devices signal errors or the completion of commands via interrupts.
- Exceptions, such as dividing by zero, invalid memory accesses, or attempts to access kernel mode instructions can be signaled via interrupts.
- Time slicing and context switches can also be implemented using the interrupt mechanism.
    - The scheduler sets a hardware timer before transferring control over to a user process.
    - When the timer raises the interrupt request line, the CPU performs a state-save, and transfers control over to the proper interrupt handler, which in turn runs the scheduler.
    - The scheduler does a state-restore of a *different* process before resetting the timer and issuing the return-from-interrupt instruction.
- A similar example involves the paging system for virtual memory - A page fault causes an interrupt, which in turn issues an I/O request and a context switch as described above, moving the interrupted process into the wait queue and selecting a different process to run. When the I/O request has completed ( i.e. when the requested page has been loaded up into physical memory ), then the device interrupts, and the interrupt handler moves the process from the wait queue into the ready queue, ( or depending on scheduling algorithms and policies, may go ahead and context switch it back onto the CPU. )
- System calls are implemented via *software interrupts,* a.k.a. *traps.* When a ( library ) program needs work performed in kernel mode, it sets command information and possibly data addresses in certain registers, and then raises a software interrupt. ( E.g. 21 hex in DOS. ) The system does a state save and then calls on the proper interrupt handler to process the request in kernel mode. Software interrupts generally have low priority, as they are not as urgent as devices with limited buffering space.

- Interrupts are also used to control kernel operations, and to schedule activities for optimal performance. For example, the completion of a disk read operation involves **two** interrupts:
    - A high-priority interrupt acknowledges the device completion, and issues the next disk request so that the hardware does not sit idle.
    - A lower-priority interrupt transfers the data from the kernel memory space to the user space, and then transfers the process from the waiting queue to the ready queue.
- The Solaris OS uses a multi-threaded kernel and priority threads to assign different threads to different interrupt handlers. This allows for the "simultaneous" handling of multiple interrupts, and the assurance that high-priority interrupts will take precedence over low-priority ones and over user processes.

### 13.2.3 Direct Memory Access

- For devices that transfer large quantities of data ( such as disk controllers ), it is wasteful to tie up the CPU transferring data in and out of registers one byte at a time.
- Instead this work can be off-loaded to a special processor, known as the ***Direct Memory Access, DMA, Controller.***
- The host issues a command to the DMA controller, indicating the location where the data is located, the location where the data is to be transferred to, and the number of bytes of data to transfer. The DMA controller handles the data transfer, and then interrupts the CPU when the transfer is complete.
- A simple DMA controller is a standard component in modern PCs, and many ***bus-mastering*** I/O cards contain their own DMA hardware.
- Handshaking between DMA controllers and their devices is accomplished through two wires called the DMA-request and DMA-acknowledge wires.
- While the DMA transfer is going on the CPU does not have access to the PCI bus ( including main memory ), but it does have access to its internal registers and primary and secondary caches.
- DMA can be done in terms of either physical addresses or virtual addresses that are mapped to physical addresses. The latter approach is known as ***Direct Virtual Memory Access, DVMA,*** and allows direct data transfer from one memory-mapped device to another without using the main memory chips.
- Direct DMA access by user processes can speed up operations, but is generally forbidden by modern systems for security and protection reasons. ( I.e. DMA is a kernel-mode operation. )
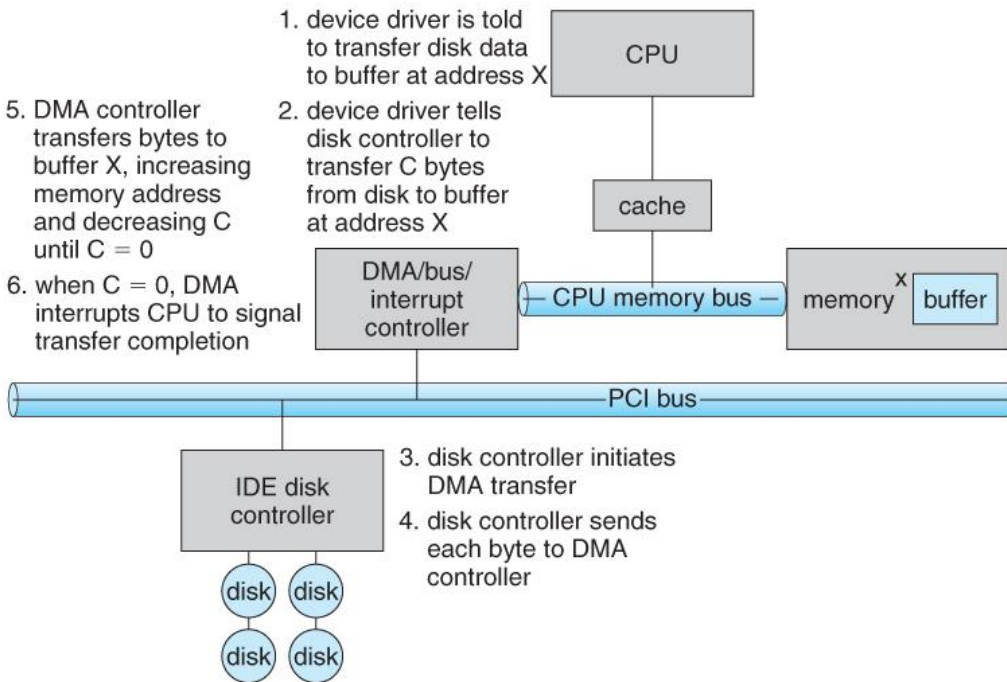- Figure 13.5 below illustrates the DMA process.

**Figure 13.5 - Steps in a DMA transfer.**

**13.2.4 I/O Hardware Summary**

## 13.3 Application I/O Interface

- User application access to a wide variety of different devices is accomplished through layering, and through encapsulating all of the device-specific code into *device drivers*, while application layers are presented with a common interface for all ( or at least large general categories of ) devices.
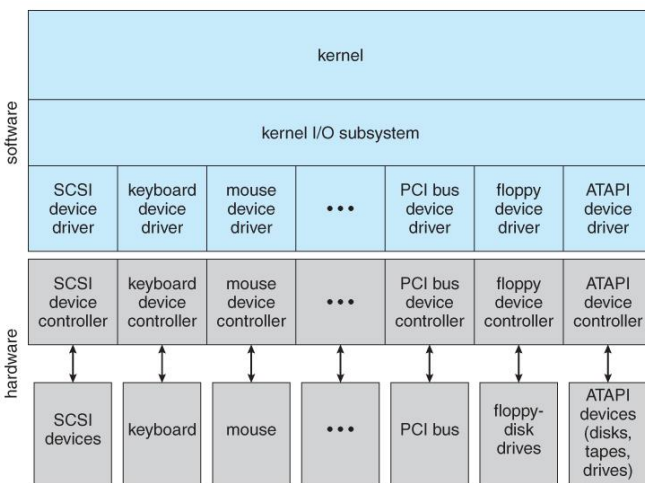


**Figure 13.6 - A kernel I/O structure.**

- Devices differ on many different dimensions, as outlined in Figure 13.7:

| aspect | variation | example |
|---|---|---|
| data-transfer mode | character<br>block | terminal<br>disk |
| access method | sequential<br>random | modem<br>CD-ROM |
| transfer schedule | synchronous<br>asynchronous | tape<br>keyboard |
| sharing | dedicated<br>sharable | tape<br>keyboard |
| device speed | latency<br>seek time<br>transfer rate<br>delay between operations | |
| I/O direction | read only<br>write only<br>read–write | CD-ROM<br>graphics controller<br>disk |

**Figure 13.7 - Characteristics of I/O devices.**

- Most devices can be characterized as either block I/O, character I/O, memory mapped file access, or network sockets. A few devices are special, such as time-of-day clock and the system timer.
- Most OSes also have an *escape,* or *back door,* which allows applications to send commands directly to device drivers if needed. In UNIX this is the *ioctl( )* system call ( I/O Control ). Ioctl( ) takes three arguments - The file descriptor for the device driver being accessed, an integer indicating the desired function to be performed, and an address used for communicating or transferring additional information.

### 13.3.1 Block and Character Devices

- *Block devices* are accessed a block at a time, and are indicated by a "b" as the first character in a long listing on UNIX systems. Operations supported include read( ), write( ), and seek( ).
    - o Accessing blocks on a hard drive directly ( without going through the filesystem structure ) is called *raw I/O*, and can speed up certain operations by bypassing the buffering and locking normally conducted by the OS. ( It then becomes the application's responsibility to manage those issues. )
    - o A new alternative is *direct I/O,* which uses the normal filesystem access, but which disables buffering and locking operations.
- Memory-mapped file I/O can be layered on top of block-device drivers.
    - o Rather than reading in the entire file, it is mapped to a range of memory addresses, and then paged into memory as needed using the virtual memory system.
    - o Access to the file is then accomplished through normal memory accesses, rather than through read( ) and write( ) system calls. This approach is commonly used for executable program code.
- *Character devices* are accessed one byte at a time, and are indicated by a "c" in UNIX long listings. Supported operations include get( ) and put( ), with more advanced functionality such as reading an entire line supported by higher-level library routines.

### 13.3.2 Network Devices

- Because network access is inherently different from local disk access, most systems provide a separate interface for network devices.
- One common and popular interface is the *socket* interface, which acts like a cable or pipeline connecting two networked entities. Data can be put into the socket at one end, and

read out sequentially at the other end. Sockets are normally full-duplex, allowing for bi-directional data transfer.
- The select( ) system call allows servers ( or other applications ) to identify sockets which have data waiting, without having to poll all available sockets.

### 13.3.3 Clocks and Timers

- Three types of time services are commonly needed in modern systems:
  - Get the current time of day.
  - Get the elapsed time ( system or wall clock ) since a previous event.
  - Set a timer to trigger event X at time T.
- Unfortunately time operations are not standard across all systems.
- A *programmable interrupt timer, PIT* can be used to trigger operations and to measure elapsed time. It can be set to trigger an interrupt at a specific future time, or to trigger interrupts periodically on a regular basis.
  - The scheduler uses a PIT to trigger interrupts for ending time slices.
  - The disk system may use a PIT to schedule periodic maintenance cleanup, such as flushing buffers to disk.
  - Networks use PIT to abort or repeat operations that are taking too long to complete. I.e. resending packets if an acknowledgement is not received before the timer goes off.
  - More timers than actually exist can be simulated by maintaining an ordered list of timer events, and setting the physical timer to go off when the next scheduled event should occur.
- On most systems the system clock is implemented by counting interrupts generated by the PIT. Unfortunately this is limited in its resolution to the interrupt frequency of the PIT, and may be subject to some drift over time. An alternate approach is to provide direct access to a high frequency hardware counter, which provides much higher resolution and accuracy, but which does not support interrupts.

### 13.3.4 Blocking and Non-blocking I/O

- With *blocking I/O* a process is moved to the wait queue when an I/O request is made, and moved back to the ready queue when the request completes, allowing other processes to run in the meantime.
- With *non-blocking I/O* the I/O request returns immediately, whether the requested I/O operation has ( completely ) occurred or not. This allows the process to check for available data without getting hung completely if it is not there.
- One approach for programmers to implement non-blocking I/O is to have a multi-threaded application, in which one thread makes blocking I/O calls ( say to read a keyboard or mouse ), while other threads continue to update the screen or perform other tasks.
- A subtle variation of the non-blocking I/O is the *asynchronous I/O,* in which the I/O request returns immediately allowing the process to continue on with other tasks, and then the process is notified ( via changing a process variable, or a software interrupt, or a callback function ) when the I/O operation has completed and the data is available for use. ( The regular non-blocking I/O returns immediately with whatever results are available, but does not complete the operation and notify the process later. )
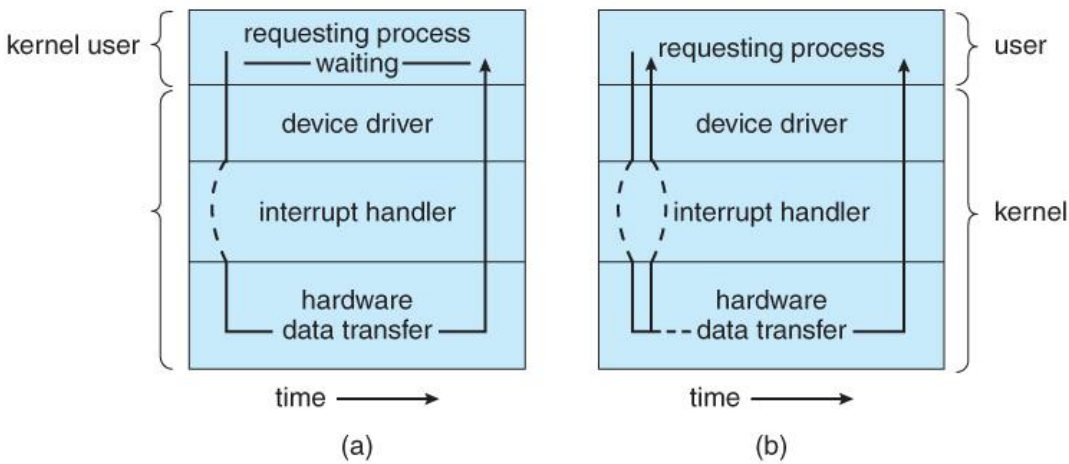
**Figure 13.8 - Two I/O methods: (a) synchronous and (b) asynchronous.**

### 13.3.5 Vectored I/O ( NEW )

## 13.4 Kernel I/O Subsystem

### 13.4.1 I/O Scheduling

- Scheduling I/O requests can greatly improve overall efficiency. Priorities can also play a part in request scheduling.
- The classic example is the scheduling of disk accesses, as discussed in detail in chapter 12.
- Buffering and caching can also help, and can allow for more flexible scheduling options.
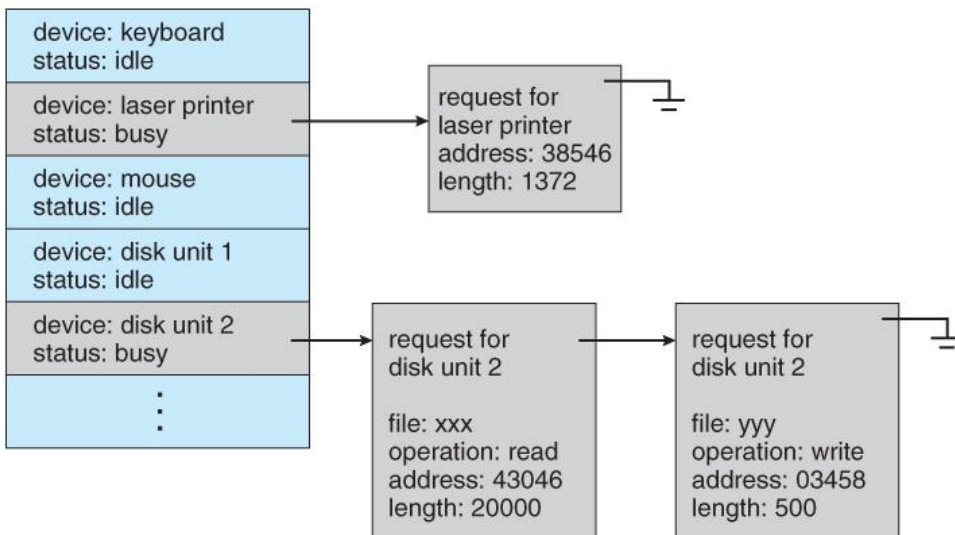- On systems with many devices, separate request queues are often kept for each device:



**Figure 13.9 - Device-status table.**

### 13.4.2 Buffering

- Buffering of I/O is performed for ( at least ) 3 major reasons:
    1. Speed differences between two devices. ( See Figure 13.10 below. ) A slow device may write data into a buffer, and when the buffer is full, the entire buffer is sent to the fast device all at once. So that the slow device still has somewhere to write while this is going on, a second buffer is used, and the two buffers alternate as each becomes full. This is known as ***double buffering.*** ( Double buffering is often

219

used in ( animated ) graphics, so that one screen image can be generated in a buffer while the other ( completed ) buffer is displayed on the screen. This prevents the user from ever seeing any half-finished screen images. )

2. Data transfer size differences. Buffers are used in particular in networking systems to break messages up into smaller packets for transfer, and then for re-assembly at the receiving side.

3. To support *copy semantics.* For example, when an application makes a request for a disk write, the data is copied from the user's memory area into a kernel buffer. Now the application can change their copy of the data, but the data which eventually gets written out to disk is the version of the data at the time the write request was made.
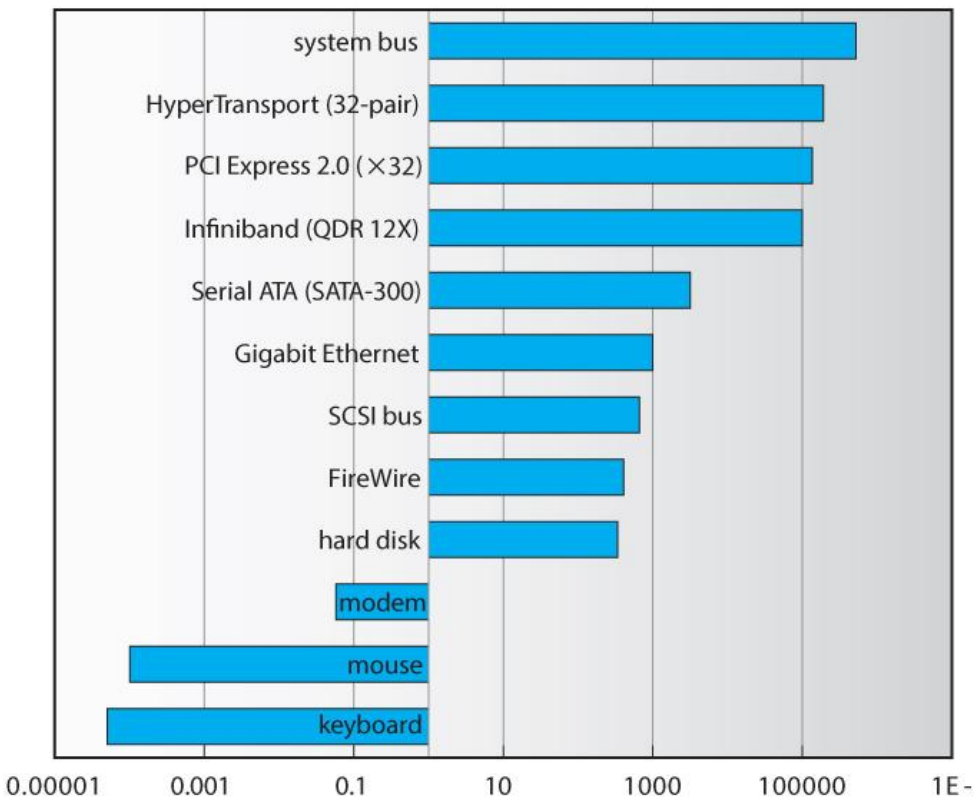


**Figure 13.10 - Sun Enterprise 6000 device-transfer rates ( logarithmic ).**

### 13.4.3 Caching

- Caching involves keeping a *copy* of data in a faster-access location than where the data is normally stored.
- Buffering and caching are very similar, except that a buffer may hold the only copy of a given data item, whereas a cache is just a duplicate copy of some other data stored elsewhere.
- Buffering and caching go hand-in-hand, and often the same storage space may be used for both purposes. For example, after a buffer is written to disk, then the copy in memory can be used as a cached copy, (until that buffer is needed for other purposes. )

### 13.4.4 Spooling and Device Reservation

- A *spool ( Simultaneous Peripheral Operations On-Line )* buffers data for ( peripheral ) devices such as printers that cannot support interleaved data streams.

220

- If multiple processes want to print at the same time, they each send their print data to files stored in the spool directory. When each file is closed, then the application sees that print job as complete, and the print scheduler sends each file to the appropriate printer one at a time.
- Support is provided for viewing the spool queues, removing jobs from the queues, moving jobs from one queue to another queue, and in some cases changing the priorities of jobs in the queues.
- Spool queues can be general ( any laser printer ) or specific ( printer number 42. )
- OSes can also provide support for processes to request / get exclusive access to a particular device, and/or to wait until a device becomes available.

### 13.4.5 Error Handling

- I/O requests can fail for many reasons, either transient ( buffers overflow ) or permanent ( disk crash ).
- I/O requests usually return an error bit ( or more ) indicating the problem. UNIX systems also set the global variable *errno* to one of a hundred or so well-defined values to indicate the specific error that has occurred. ( See errno.h for a complete listing, or man errno. )
- Some devices, such as SCSI devices, are capable of providing much more detailed information about errors, and even keep an on-board error log that can be requested by the host.

### 13.4.6 I/O Protection

- The I/O system must protect against either accidental or deliberate erroneous I/O.
- User applications are not allowed to perform I/O in user mode - All I/O requests are handled through system calls that must be performed in kernel mode.
- Memory mapped areas and I/O ports must be protected by the memory management system, **but** access to these areas cannot be totally denied to user programs. ( Video games and some other applications need to be able to write directly to video memory for optimal performance for example. ) Instead the memory protection system restricts access so that only one process at a time can access particular parts of memory, such as the portion of the screen memory corresponding to a particular window.
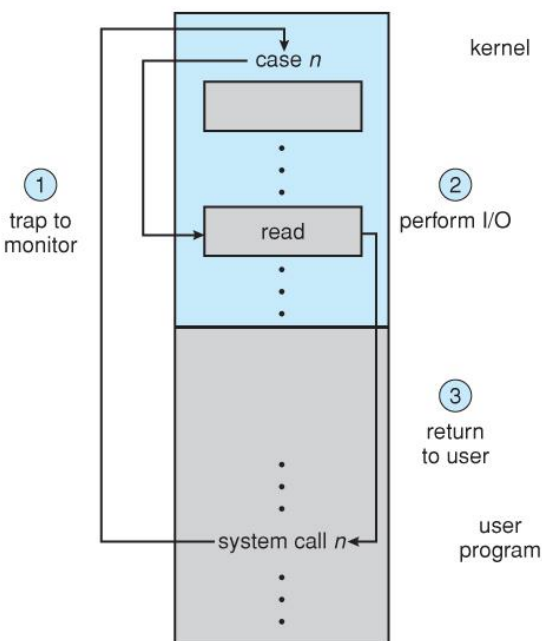


**Figure 13.11 - Use of a system call to perform I/O.**

### 13.4.7 Kernel Data Structures

- The kernel maintains a number of important data structures pertaining to the I/O system, such as the open file table.
- These structures are object-oriented, and flexible to allow access to a wide variety of I/O devices through a common interface. ( See Figure 13.12 below. )
- Windows NT carries the object-orientation one step further, implementing I/O as a message-passing system from the source through various intermediaries to the device.
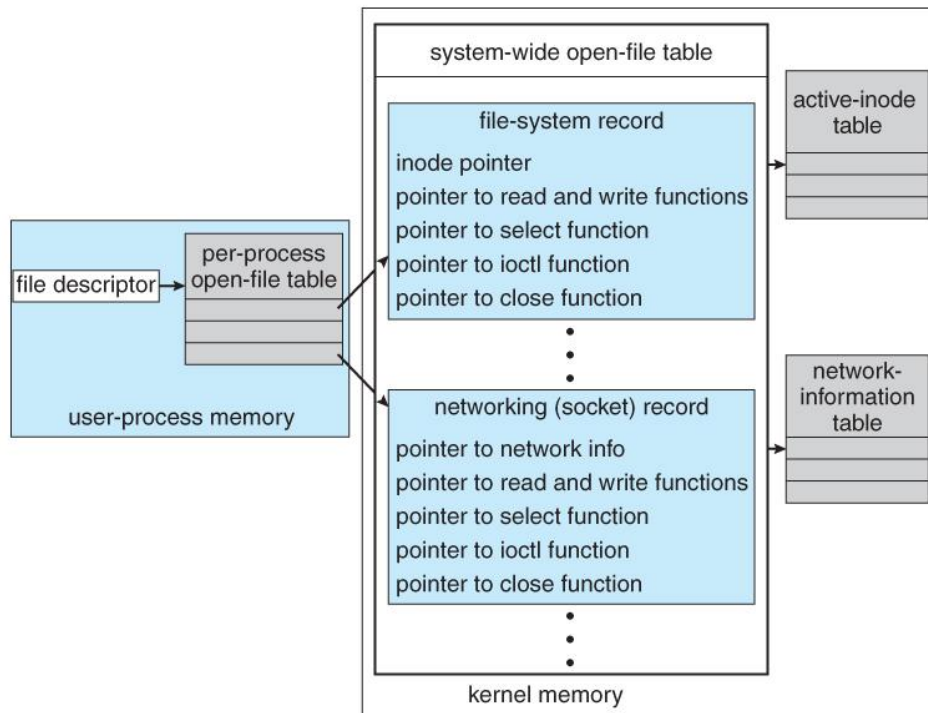


**Figure 13.12 - UNIX I/O kernel structure.**

### 13.4.6 Kernel I/O Subsystem Summary

## 13.5 Transforming I/O Requests to Hardware Operations

- Users request data using file names, which must ultimately be mapped to specific blocks of data from a specific device managed by a specific device driver.
- DOS uses the colon separator to specify a particular device ( e.g. C:, LPT:, etc. )
- UNIX uses a *mount table* to map filename prefixes ( e.g. /usr ) to specific mounted devices. Where multiple entries in the mount table match different prefixes of the filename the one that matches the longest prefix is chosen. ( e.g. /usr/home instead of /usr where both exist in the mount table and both match the desired file. )
- UNIX uses special *device files,* usually located in /dev, to represent and access physical devices directly.
    - o Each device file has a major and minor number associated with it, stored and displayed where the file size would normally go.
    - o The major number is an index into a table of device drivers, and indicates which device driver handles this device. ( E.g. the disk drive handler. )
    - o The minor number is a parameter passed to the device driver, and indicates which specific device is to be accessed, out of the many which may be handled by a particular device driver. ( e.g. a particular disk drive or partition. )

- A series of lookup tables and mappings makes the access of different devices flexible, and somewhat transparent to users.
- Figure 13.13 illustrates the steps taken to process a ( blocking ) read request:
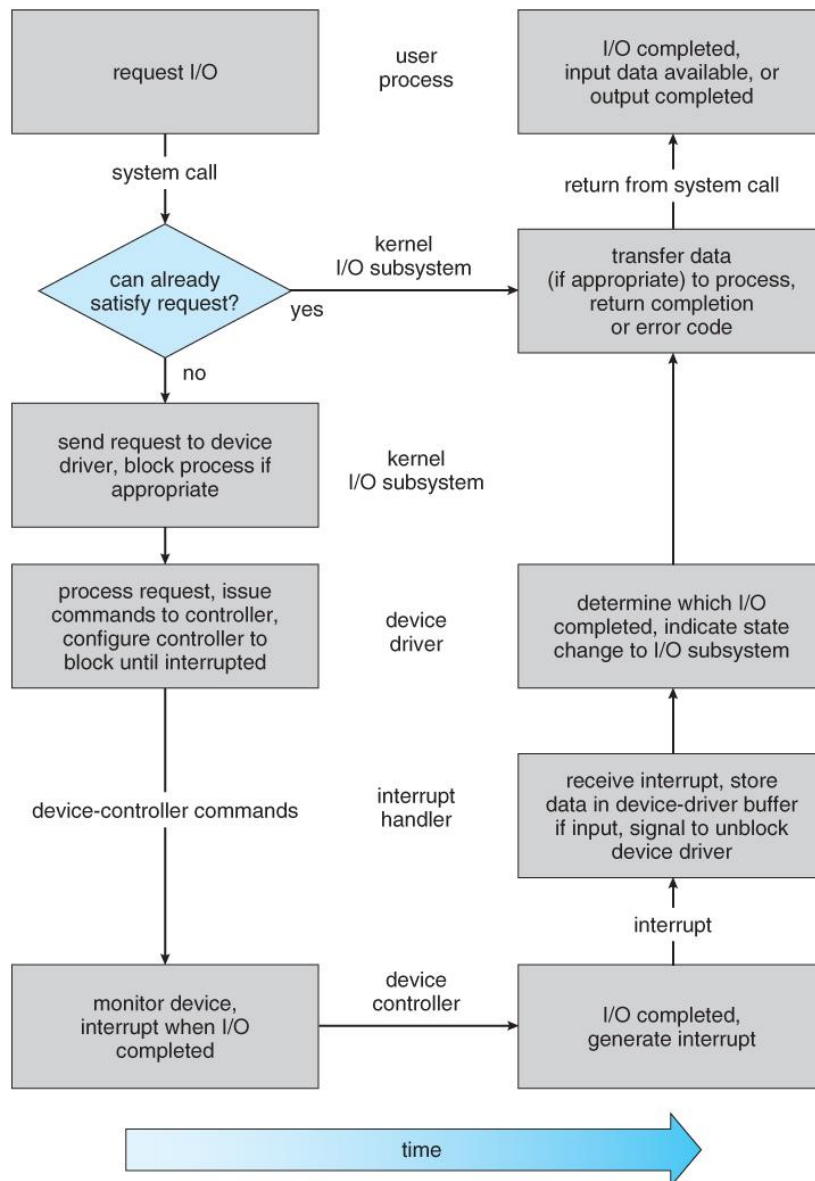


**Figure 13.13 - The life cycle of an I/O request.**

## 13.6 STREAMS ( Optional )

- The *streams* mechanism in UNIX provides a bi-directional pipeline between a user process and a device driver, onto which additional modules can be added.
- The user process interacts with the *stream head.*
- The device driver interacts with the *device end.*
- Zero or more *stream modules* can be pushed onto the stream, using ioctl( ). These modules may filter and/or modify the data as it passes through the stream.
- Each module has a *read queue* and a *write queue.*
- *Flow control* can be optionally supported, in which case each module will buffer data until the adjacent module is ready to receive it. Without flow control, data is passed along as soon as it is ready.
- User processes communicate with the stream head using either read( ) and write( ) ( or putmsg( ) and getmsg( ) for message passing. )

223

- Streams I/O is asynchronous ( non-blocking ), except for the interface between the user process and the stream head.
- The device driver **must** respond to interrupts from its device - If the adjacent module is not prepared to accept data and the device driver's buffers are all full, then data is typically dropped.
- Streams are widely used in UNIX, and are the preferred approach for device drivers. For example, UNIX implements sockets using streams.
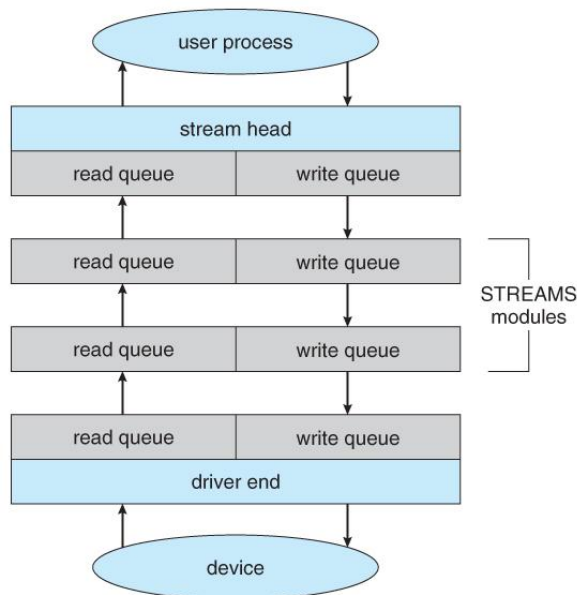


**Figure 13.14 - The SREAMS structure.**

## 13.7 Performance ( Optional )

- The I/O system is a major factor in overall system performance, and can place heavy loads on other major components of the system ( interrupt handling, process switching, memory access, bus contention, and CPU load for device drivers just to name a few. )
- Interrupt handling can be relatively expensive ( slow ), which causes programmed I/O to be faster than interrupt-driven I/O when the time spent busy waiting is not excessive.
- Network traffic can also put a heavy load on the system. Consider for example the sequence of events that occur when a single character is typed in a telnet session, as shown in figure 13.15. ( And the fact that a similar set of events must happen in reverse to echo back the character that was typed. ) Sun uses in-kernel threads for the telnet daemon, increasing the supportable number of simultaneous telnet sessions from the hundreds to the thousands.
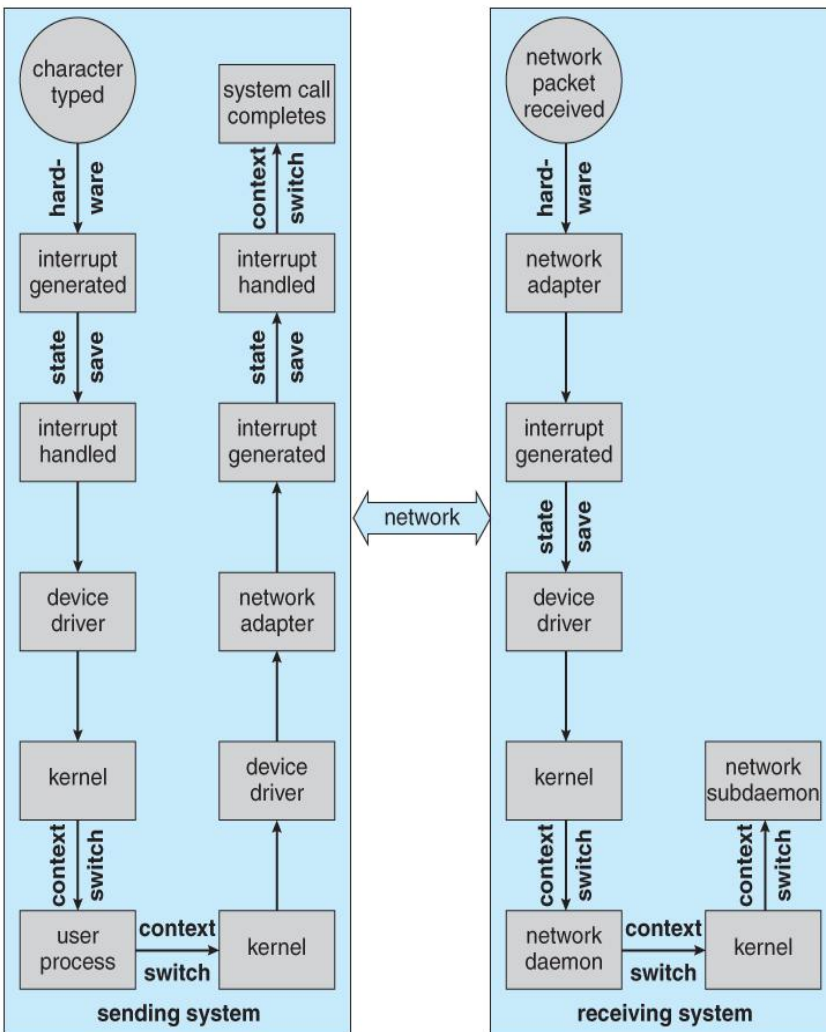
**Figure 13.15 - Intercomputer communications.**

- Other systems use ***front-end processors*** to off-load some of the work of I/O processing from the CPU. For example a ***terminal concentrator*** can multiplex with hundreds of terminals on a single port on a large computer.
- Several principles can be employed to increase the overall efficiency of I/O processing:
  1. Reduce the number of context switches.
  2. Reduce the number of times data must be copied.
  3. Reduce interrupt frequency, using large transfers, buffering, and polling where appropriate.
  4. Increase concurrency using DMA.
  5. Move processing primitives into hardware, allowing their operation to be concurrent with CPU and bus operations.
  6. Balance CPU, memory, bus, and I/O operations, so a bottleneck in one does not idle all the others.
- The development of new I/O algorithms often follows a progression from application level code to on-board hardware implementation, as shown in Figure 13.16. Lower-level implementations are faster and more efficient, but higher-level ones are more flexible and easier to modify. Hardware-level functionality may also be harder for higher-level authorities ( e.g. the kernel ) to control.
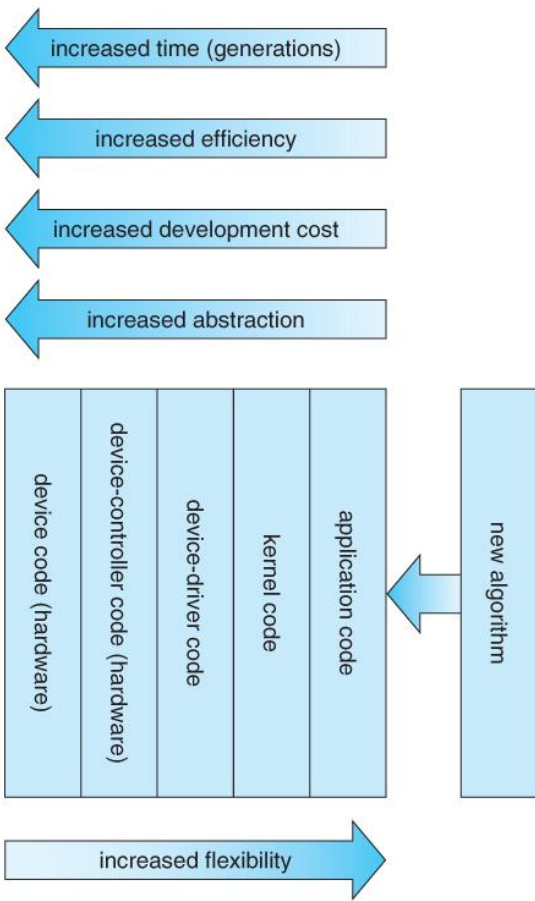
**Figure 13.16 - Device functionality progression.**

## 13.8 Summary

**Course material**

# UNIT V CASE STUDIES

## 5.1 The Linux System

☐ An operating system is a program that acts as an interface between the user and the computer hardware and controls the execution of all kinds of programs. The Linux open source operating system, or Linux OS, is a freely distributable, cross-platform operating system based on UNIX.

☐ The Linux consist of a kernel and some system programs. There are also some application programs for doing work. The kernel is the heart of the operating system which provides a set of tools that are used by system calls.

☐ The defining component of Linux is the Linux kernel, an operating system kernel first released on 5 October 1991 by *Linus Torvalds.*

☐ A Linux-based system is a modular Unix-like operating system. It derives much of its basic design from principles established in UNIX. Such a system uses a monolithic kernel which handles process control, networking, and peripheral and file system access.

## 5.2 Important features of Linux Operating System

☐ **Portable** - Portability means software can work on different types of hardware in same way. Linux kernel and application programs supports their installation on any kind of hardware platform.

☐ **Open Source** - Linux source code is freely available and it is community based development project.

☐ **Multi-User** & **Multiprogramming** - Linux is a multiuser system where multiple users can access system resources like memory/ ram/ application programs at same time. Linux is a multiprogramming system means multiple applications can run at same time.

☐ **Hierarchical File System** - Linux provides a standard file structure in which system files/ user files are arranged.

☐ **Shell** - Linux provides a special interpreter program which can be used to execute commands of the operating system.

☐ **Security** - Linux provides user security using authentication features like password protection/ controlled access to specific files/ encryption of data.

## 5.3 Components of Linux System

Linux Operating System has primarily three components

☐ **Kernel** - Kernel is the core part of Linux. It is responsible for all major activities of this operating system. It is consists of various modules and it interacts directly with the underlying hardware. Kernel provides the required abstraction to hide low level hardware details to system or application programs.

☐ **System Library** - System libraries are special functions or programs using which application programs or system utilities accesses Kernel's features. These libraries implements most of the functionalities of the operating system and do not requires kernel module's code access rights.

☐ **System Utility** - System Utility programs are responsible to do specialized, individual level tasks Installed components of a Linux system include the following:

☐ A **bootloader** is a program that loads the Linux kernel into the computer's main memory, by being executed by the computer when it is turned on and after the firmware initialization is performed.

☐ An **init** program is the first process launched by the Linux kernel, and is at the root of the process tree.

☐ **Software libraries**, which contain code that can be used by running processes. The most commonly used software library on Linux systems, the GNU C Library (glibc), C standard library and Widget toolkits.

☐ **User interface programs** such as command shells or windowing environments. The user interface, also known as the shell, is either a command-line interface (CLI), a graphical user interface (GUI), or through controls attached to the associated hardware.

## 5.4 Architecture

Linux System Architecture is consists of following layers

1. **Hardware layer** - Hardware consists of all peripheral devices (RAM/ HDD/ CPU etc).

2. **Kernel** - Core component of Operating System, interacts directly with hardware, provides low level

services to upper layer components.

3. **Shell** - An interface to kernel, hiding complexity of kernel's functions from users. Takes commands from user and executes kernel's functions.

4. **Utilities** - Utility programs giving user most of the functionalities of an operating systems.

## 5.5 Modes of operation

☐ **Kernel Mode:**

☐ Kernel component code executes in a special privileged mode called *kernel mode* with full access to all resources of the computer.

☐ This code represents a single process, executes in single address space and do not require any context switch and hence is very efficient and fast.

☐ Kernel runs each processes and provides system services to processes, provides protected access to hardware to processes.

☐ **User Mode:**

☐ The system programs use the tools provided by the kernel to implement the various services required from an operating system. System programs, and all other programs, run `on top of the kernel', in what is called the user mode.

☐ Support code which is not required to run in kernel mode is in System Library.

☐ User programs and other system programs work in User Mode which has no access to system hardware and kernel code.

☐ User programs/ utilities use System libraries to access Kernel functions to get system's low level tasks.

## 5.6 Major Services provided by LINUX System

### 1. Initialization (init)

The single most important service in a LINUX system is provided by **init** program. The *init* is started as the first process of every LINUX system, as the last thing the kernel does when it boots. When init starts, it continues the boot process by doing various startup chores (checking and mounting file systems, starting daemons, etc).

### 2. Logins from terminals (getty)

Logins from terminals (via serial lines) and the console are provided by the **getty** program. **init** starts a separate instance of **getty** for each terminal upon which logins are to be allowed. Getty reads the username and runs the login program, which reads the password. If the username and password are correct, login runs the shell.

### 3. Logging and Auditing (syslog)

The kernel and many system programs produce error, warning, and other messages. It is often important that these messages can be viewed later, so they should be written to a file. The program doing this logging operation is known as **syslog**.

### 4. Periodic command execution (cron & at)

Both users and system administrators often need to run commands periodically. For example, the system administrator might want to run a command to clean the directories with temporary files from old files, to keep the disks from filling up, since not all programs clean up after themselves correctly.

o The **cron** service is set up to do this. Each user can have a *crontab* file, where the lists the commands wish to execute and the times they should be executed.

o The **at** service is similar to cron, but it is once only: the command is executed at the given time, but it is not repeated.

### 5. Graphical user interface

o UNIX and Linux don't incorporate the user interface into the kernel; instead, they let it be implemented by user level programs. This applies for both text mode and graphical

environments. This arrangement makes the system more flexible.

 The graphical environment primarily used with Linux is called the X Window System (X for short) that provides tools with which a GUI can be implemented. Some popular window managers are blackbox and windowmaker. There are also two popular desktop managers, KDE and Gnome.

## 6. Network logins (telnet, rlogin & ssh)

Network logins work a little differently than normal logins. For each person logging in via the network there is a separate virtual network connection. It is therefore not possible to run a separate getty for each virtual connection. There are several different ways to log in via a network, **telnet** and **ssh** being the major ones in TCP/IP networks.

Most of Linux system administrators consider telnet and rlogin to be insecure and prefer ssh, the ``secure shell'', which encrypts traffic going over the network, thereby making it far less likely that the malicious can ``sniff'' the connection and gain sensitive data like usernames and passwords.

## 7. Network File System (NFS & CIFS)

One of the more useful things that can be done with networking services is sharing files via a network file system. Depending on your network this could be done over the Network File System (NFS), or over the Common Internet File System (CIFS).

NFS is typically a 'UNIX' based service. In Linux, NFS is supported by the kernel. CIFS however is not. In Linux, CIFS is supported by **Samba**. With a network file system any file operations done by a program on one machine are sent over the network to another computer.

**UNIX Timeline -** Simplified history of Unix-like operating systems

## 5.7 SYSTEM ADMINISTRATOR

☐ A system administrator is a person who is responsible for the configuration and reliable operation of computer systems, especially multi-user computers, such as servers.

☐ The system administrator seeks to ensure that the uptime, performance, resources, and security of the computers without exceeding the budget.

☐ To meet these needs, a system administrator may acquire, install, or upgrade computer components and software, provide routine automation, maintain security policies AND troubleshoot.

## 5.7.1 Responsibilities of a System Administrator

A system administrator's responsibilities might include:

☐ Installing and configuring new hardware and software.

☐ Applying operating system updates, patches, and configuration changes.

☐ Analyzing system logs and identifying potential issues with computer systems.

☐ Introducing and integrating new technologies into existing data center environments and configuring, adding, and deleting file systems.

☐ Performing routine audit of systems and software.

☐ Adding, removing, or updating user account information, resetting passwords, etc.

☐ Responsibility for security and documenting the configuration of the system.

☐ Troubleshooting any reported problems.

☐ System performance tuning.

## 5.7.2 Various System Administrator Roles

In a larger company, these may all be separate positions within a computer support or Information Services (IS) department. In a smaller group they may be shared by a few sysadmins, or even a single person.

☐ A **database administrator** (DBA) maintains a database system, and is responsible for the integrity of the data and the efficiency and performance of the system.

☐ A **network administrator** maintains network infrastructure such as switches and routers, and diagnoses problems with these or with the behaviour of network-attached computers.

☐ A **security administrator** is a specialist in computer and network security, including the administration of security devices such as firewalls, as well as consulting on general security measures.

☐ A **web administrator** maintains web server services (such as Apache or IIS) that allow for internal or external access to web sites. Tasks include managing multiple sites, administering security, and configuring necessary components and software.

☐ A **computer operator** performs routine maintenance and upkeep, such as changing backup tapes or replacing failed drives in a redundant array of independent disks (RAID).

☐ A **postmaster** administers a mail server.

☐ A **Storage Administrator (SAN)** can create, provision, add or remove Storage to/from Computer systems. Storage can be attached locally to the system or from a storage area network (SAN) or network-attached storage (NAS).

## 5.7.3 Requirements for LINUX system administrator

1. While specific knowledge is a boon, system administrator should possess basic knowledge about all aspects of Linux. For example, a little knowledge about Solaris, BSD, nginx or various flavors of Linux.

2. Knowledge in at least one of the upper tier scripting language such as Python, Perl, Ruby or more.

3. To be a system administrator, he/she at least needs to have some hands-on experience of system management, system setup and managing Linux or Solaris based servers as well as configuring them.

4. Knowledge in shell programming such as Buorne or Korn and architecture.

5. Knowledge about storage technologies like FC, NFS or iSCSI is great, while knowledge regarding backup technologies is a must for a system administrator.

6. Knowledge in testing methodologies like Subversion or Git is great, while knowledge of version control is also an advantage.

7. Knowledge about basics of configuration management tools like Puppet and Chef.

8. Skills with system and application monitoring tools like SNMP or Nagios are also important, as they show your ability as an administrator in a team setting.

9. Knowledge about how to operate virtualized VMWare or Xen Server, Multifunction Server and Samba

10. An ITIL Foundation certification for Linux system administrator.

## 5.8 SETTING UP A LINUX MULTIFUNCTION SERVER

A Linux machine can be configured as a server either by compiling several well-defined scripts and off-line downloaded packages or through on-line installation method. Setting up a multifunction server, the system administrator should have knowledge about a series of shell commands. A Linux machine can be configured as any of following application servers such as,

• A Web Server (Apache 2.0.x)

• A Mail Server (Postfix)

• A DNS Server (BIND 9)

• An FTP Server (ProFTPD)

• Mail Delivery Agents (POP3/POP3s/IMAP/IMAPs)

• Webalizer for web site statistics

Files and directories shared by Linux system, as viewed from a Windows PC

## 5.8.1 Server Requirements

To set up a Linux Internet server, we will need a connection to the Internet and a static IP address. The system can also be setup with the address leased by ISP and configure it statically. Computer with at least a Pentium III CPU, a minimum of 256 MB of RAM, and a 10 GB hard drive is preferred. Obviously, a newer CPU and additional memory will provide better performance. This chapter is based on Debian's stable version. We strongly suggest using a CD with the Netinstall kernel. The Debian web site provides downloadable CD images.

### 5.8.2 Installing & Configuring Network Services

Administrator should log into the server from a remote console on desktop. It is recommended to do further administration from another system (even a laptop), because a secure server normally runs in what is called headless mode—that is, it has no monitor or keyboard.
Get used to administering the server like this. A SSH client on the remote machine is needed which virtually all Linux distributions have and which can be downloaded for other operating systems as well.

**Configuring the Network**

If DHCP is used during the Debian installation, Server with a static IP address should be configured as follows,

1. To change the settings to use a static IP address, you'll need to become root and edit the file /etc/network/interfaces to suit your needs. As an example, we'll use the IP address 70.153.258.42.

2. To add the IP address 70.153.258.42 to the interface eth0, we must change the file to look like this (you'll have to obtain some of the information from your ISP):

auto eth0
iface eth0 inet static
address 70.153.258.42
netmask 255.255.255.248
network 70.153.258.0
broadcast 70.153.258.47
gateway 70.153.258.46

3. After editing the /etc/network/interfaces file, restart the network by entering:
# /etc/init.d/networking restart

4. To edit /etc/resolv.conf and add nameservers to resolve Internet hostnames to their corresponding IP addresses. At this point, we will simply set up a minimal DNS server. Our **resolv.conf** looks as follows:

search server
nameserver 70.153.258.42
nameserver 70.253.158.45
nameserver 151.164.1.8

5. Now edit /etc/hosts and add your IP addresses:
127.0.0.1 localhost.localdomain localhost server1
70.153.258.42 server1.centralsoft.org server1

6. Now, to set the hostname, enter these commands:
# echo server1.centralsoft.org > /etc/hostname
# /bin/hostname -F /etc/hostname

7. verify that you configured your hostname correctly by running the *hostname* command:
~$ hostname -f
server1.centralsoft.org

### 5.9 Providing Domain Name Services (BIND - the ubiquitous DNS server)

☐ Debian provides a stable version of BIND in its repositories. BIND can be installed, setup

and secure it in a *chroot* environment, meaning it won't be able to see or access files outside its own directory tree. This is an important security technique.

☐ The term *chroot* refers to the trick of changing the root filesystem (the /directory) that a process sees, so that most of the system is effectively inaccessible to it.

☐ The BIND server also can be configured to run as a non-root user. That way, if someone gains access to BIND, he/she won't gain root privileges or be able to control other processes.

1. To install BIND on your Debian server, run this command:

# apt-get install bind9

Debian downloads and configures the file as an Internet service and the status can be seen on the console:

Setting up bind9 (9.2.4-1)
Adding group `bind' (104) - Done.
Adding system user `bind'
Adding new user `bind' (104) with group `bind'.
Not creating home directory.
Starting domain name service: named.

2. To put BIND in a secured environment, create a directory where the service can run unexposed to other processes. First stop the service by running the following command:

# /etc/init.d/bind9 stop

3. Edit the file /etc/default/bind9 so that the daemon will run as the unprivileged user bind, chrooted to /var/lib/named. Change the line:

OPTS="-u bind"

So that it reads:

OPTIONS="-u bind -t /var/lib/named"

4. To provide a complete environment for running BIND, create the necessary directories under /var/lib:

# mkdir -p /var/lib/named/etc
# mkdir /var/lib/named/dev
# mkdir -p /var/lib/named/var/cache/bind
# mkdir -p /var/lib/named/var/run/bind/run

Then move the config directory from /etc to /var/lib/named/etc:

# mv /etc/bind /var/lib/named/etc

Next, create a symbolic link to the new config directory from the old location, to avoid problems when BIND is upgraded in the future:

# ln -s /var/lib/named/etc/bind /etc/bind

Make null and random devices for use by BIND, and fix the permissions of the directories:

# mknod /var/lib/named/dev/null c 1 3
# mknod /var/lib/named/dev/random c 1 8

Then change permissions and ownership on the files:

# chmod 666 /var/lib/named/dev/null
/var/lib/named/dev/random
# chown -R bind:bind /var/lib/named/var/*
# chown -R bind:bind /var/lib/named/etc/bind

5. Finally, start BIND:

# /etc/init.d/bind9 start

6. To check whether named is functioning without any trouble.
Execute this command:

**server1:/home/admin# rndc status**

number of zones: 6
debug level: 0

xfers running: 0
xfers deferred: 0
soa queries in progress: 0
query logging is OFF
server is up and running
server1:/home/admin#

**Setting up Ubuntu shares in a Windows environment**

**Ubuntu's setup screen for file-sharing services**

## 5.10 Virtualization

☐ Virtualization refers to the act of creating a virtual (rather than actual) version of something, including a virtual computer hardware platform, operating system (OS), storage device, or computer network resources.

**Traditional Architecture vs. Virtual Architecture**

**Virtual Machine Server – A Layered Approach**

☐ *Hardware virtualization* or *platform virtualization* refers to the creation of a virtual machine that acts like a real computer with an operating system. Software executed on these virtual machines is separated from the underlying hardware resources.

☐ *Hardware virtualization* hides the physical characteristics of a computing platform from users, instead showing another abstract computing platform.

☐ For example, a computer that is running Microsoft Windows may host a virtual machine that looks like a computer with the Ubuntu Linux operating system; Ubuntu-based software can be run on the virtual machine.

**Benefits of Virtualization**

**Hardware Virtualization**

1. Instead of deploying several physical servers for each service, only one server can be used. Virtualization let multiple OSs and applications to run on a server at a time. Consolidate hardware to get vastly higher productivity from fewer servers.

2. If the preferred operating system is deployed as an image, so we needed to go through the installation process only once for the entire infrastructure.

3. **Improve business continuity:** Virtual operating system images allow us for instant recovery in case of a system failure. The crashed system can be restored back by coping the virtual image.

4. **Increased uptime:** Most server virtualization platforms offer a number of advanced features that just aren't found on physical servers which increases servers' uptime. Some of features are live migration, storage migration, fault tolerance, high availability, and distributed resource scheduling.

5. **Reduce capital and operating costs:** Server consolidation can be done by running multiple virtual machines (VM) on a single physical server. Fewer servers means lower capital and operating costs.

**Architecture - Virtualization**

The heart of virtualization is the "virtual machine" (VM), a tightly isolated software container with an operating system and application inside. Because each virtual machine is completely separate and independent, many of them can run simultaneously on a single computer. A thin layer of software called a hypervisor decouples the virtual machines from the host and dynamically allocates computing resources to each virtual machine as needed. This architecture redefines your computing equation and delivers:

☐ **Many applications on each server:** As each virtual machine encapsulates an entire machine, many applications and operating systems can run on a single host at the same time.

☐ **Maximum server utilization, minimum server count**: Every physical machine is used to

its full capacity, allowing you to significantly reduce costs by deploying fewer servers overall.

 **Faster, easier application and resource provisioning:** As self-contained software files, virtual machines can be manipulated with copy-and-paste ease. Virtual machines can even be transferred from one physical server to another while running, via a process known as live migration.

## 5.10.1 Setting up a VMware Workstation

### VMware Workstation

VMware Workstation is developed and sold by VMware, Inc., a division of EMC Corporation. VMware Workstation is a hypervisor that runs on x86 or x86-64 computers; it enables users to set up one or more virtual machines (VMs) on a single physical machine, and use them simultaneously along with the actual machine.

Each virtual machine can execute its own operating system, including versions of Microsoft Windows, Linux, BSD, and MS-DOS. VMware Workstation supports bridging existing host network adapters and share physical disk drives and USB devices with a virtual machine. In addition, it can simulate disk drives. It can mount an existing ISO image file into a virtual optical disc drive so that the virtual machine sees it as a real one. Likewise, virtual hard disk drives are made via .*vmdk* files.

VMware Workstation can save the state of a virtual machine (a "snapshot") at any instant. These snapshots can later be restored, effectively returning the virtual machine to the saved state.

### VMware Workstation

VMware Workstation includes the ability to designate multiple virtual machines as a team which can then be powered on, powered off, suspended or resumed as a single object, making it particularly useful for testing client-server environments.

### VMWare Player

The VMware Player, a virtualization package of basically similar, but reduced, functionality, is also available, and is free of charge for non-commercial use, or for distribution or other use by written agreement.

VMware Player is a virtualization software package supplied free of charge by VMware, Inc. VMware Player can run existing virtual appliances and create its own virtual machines. It uses the same virtualization core as VMware Workstation, a similar program with more features, but not free of charge. VMware Player is available for personal non-commercial use, or for distribution or other use by written agreement.

VMware claims the Player offers better graphics, faster performance, and tighter integration for running Windows XP under Windows Vista or Windows 7 than Microsoft's Windows XP Mode running on Windows Virtual PC, which is free of charge for all purposes.

### VMware Tools

VMware Tools is a package with drivers and other software that can be installed in guest operating systems to increase their performance. It has several components, including the following drivers for the emulated hardware:

 VESA-compliant graphics for the guest machine to access high screen resolutions
 Network drivers for the vmxnet2  Mouse integration, Drag-and-drop file support
 Clipboard sharing between host and guest
 Time synchronization capabilities (guest syncs with host machine's clock)
 Support for Unity, a feature that allows seamless integration of applications with the host desktop

### Installing and Configuring VMWare

1. Download VMware Server 2. VMware management console on a remote Ubuntu desktop behind a firewall at a remote location. Run the following command:

$gksu vmware-server-console

2. Install the VMware Server 2.0.2 rpm as shown below.

\# rpm -ivh VMware-server-2.0.2-203138.i386.rpm

Preparing...

1:VMware-server

\############################################# [100%]

The installation of VMware Server 2.0.2 for Linux completed successfully.

You can decide to remove this software from your system at any time by invoking the following command:

rpm -e VMware-server

Before running VMware Server for the first time, you need to configure it for your running kernel by invoking the following command:

/usr/bin/vmware-config.pl

3. Configure VMware Server 2 using *vmware-config.pl*. Execute the vmware-config.pl as shown below. Accept default values for everything. Partial output of the vmwareconfig. pl is shown below.

\# /usr/bin/vmware-config.pl

4. Go to VMware Infrastructure Webaccess. Go to **https://{host-os-ip}:8333/ui** to access the VMware Infrastructure web access console.

**VMware Web Access Login**

**Installing a VMware Guest OS**

1. **Start VMware Workstation**

*Windows host*: Double-click the VMware Workstation icon on your desktop or use the Start menu (Start > Programs > VMware > VMware Workstation).

*Linux host*: In a terminal window, enter the command

**vmware &**

2. **Start the New Virtual Machine Wizard**

When you start VMware Workstation, you can open an existing virtual machine or create a new one. Choose File > New > Virtual Machine to begin creating your virtual machine.

3. Select the method you want to use for configuring your virtual machine.

If you select *Typical*, the wizard prompts you to specify or accept defaults for the following choices:

☐ The guest operating system

☐ The virtual machine name and the location of the virtual machine's files

☐ The network connection type

☐ Whether to allocate all the space for a virtual disk at the time you create it

☐ Whether to split a virtual disk into 2GB files

If you select *Custom*, the wizard prompts you to specify or accept defaults for the following choices:

☐ Make a legacy virtual machine that is compatible with Workstation 4.x, GSX Server 3.x, ESX Server 2.x and VMware ACE 1.x.

☐ Use an IDE virtual disk for a guest operating system that would otherwise have a SCSI virtual disk created by default

☐ Use a physical disk rather than a virtual disk and Set memory options that are different from the defaults

4. Select a guest operating system and type a name and folder for the virtual machine.

**Linux hosts:** The default location for this Windows XP Professional virtual machine is <homedir>/vmware/winXPPro, where <homedir> is the home directory of the user who is

currently logged on.

5. Specify the number of processors for the virtual machine. The setting Two is supported only for host machines with at least two logical processors.

If you selected **Custom** as your configuration path, you may adjust the memory settings or accept the defaults, then click Next to continue.

6. Configure the networking capabilities of the virtual machine.

If you selected *Typical* as your configuration path, click Finish and the wizard sets up the files needed for your virtual machine.

If you selected *Custom* as your configuration path, continue with the steps below to configure a disk for your virtual machine.

7. Select whether to create an IDE or SCSI disk and specify the capacity of the virtual disk.

8. Click Finish. The wizard sets up the files needed for your virtual machine.

## 5.10.2 Setting up a XEN Workstation

### XEN Workstation

Xen is a hypervisor using a microkernel design, providing services that allow multiple computer operating systems to execute on the same computer hardware concurrently.

The University of Cambridge Computer Laboratory developed the first versions of Xen.

The Xen community develops and maintains Xen as free and open-source software, subject to the requirements of the GNU General Public License (GPL), version 2. Xen is currently available for the IA-32, x86-64 and ARM instruction sets.

XenServer runs directly on server hardware without requiring an underlying operating system, which results in an efficient and scalable system. XenServer works by abstracting elements from the physical machine (such as hard drives, resources and ports) and allocating them to the virtual machines running on it.

### XEN Environment

Responsibilities of the hypervisor include memory management and CPU scheduling of all virtual machines, and for launching the most privileged domain - the only virtual machine which by default has direct access to hardware. From the dom0 the hypervisor can be managed and unprivileged domains can be launched.

### Benefits of Using XenServer

### 1. Using XenServer reduces costs by:

• Consolidating multiple VMs onto physical servers

• Reducing the number of separate disk images that need to be managed

• Allowing for easy integration with existing networking and storage infrastructures

### 2. Using XenServer increases flexibility by:

• Allowing you to schedule zero downtime maintenance by using XenMotion to live migrate VMs between XenServer hosts

• Increasing availability of VMs by using High Availability to configure policies that restart VMs on another XenServer host if one fails

• Increasing portability of VM images, as one VM image will work on a range of deployment infrastructures

### Course material

### Administering XenServer

☐ There are two methods by which to administer XenServer: XenCenter and the XenServer Command-Line Interface (CLI).

☐ **XenCenter** is a graphical, Windows-based user interface. XenCenter allows you to manage XenServer hosts, pools and shared storage, and to deploy, manage and monitor VMs from your Windows desktop machine.

☐ The XenCenter on-line Help is a useful resource for getting started with XenCenter and for context-sensitive assistance.

## Installing and Configuring XenServer

1. Type the following command to get information about xen server package

# yum info xen

2. Run the system-config-securitylevel program or edit /etc/selinux/config to looks as follows:

SELINUX=Disabled

SELINUXTYPE=targeted

If you changed the SELINUX value from enforcing, you'll need to reboot Fedora before proceeding.

3. This command will install the Xen hypervisor, a Xen-modified Fedora kernel called *domain 0*, and various utilities:

# yum install kernel-xen0

4. To make the Xen kernel the default, change this line:

default=1

to

default=0

5. Now you can reboot. Xen should start automatically, but let's check:

# /usr/sbin/xm list

Name ID Mem(MiB) VCPUs State Time(s)

Domain-0 0 880 1 r----- 20.5

The output should show that Domain-0 is running. Domain 0 controls all the guest operating systems that run on the processor, similarly to how the kernel controls processes in an operating system.

## Installing a Xen Guest OS from the Command-line

1. **Preparing the System for virt-install**

Fedora Linux does not install VNC by default. To verify whether VNC is installed, run the following command from a Terminal Window:

rpm -q vnc

**Course material**

If rpm reports that VNC is not installed, it may be installed from root as follows:

yum install vnc

2. **Running virt-install to Build the Xen Guest System**

virt-install must be run as root and, once invoked, will ask a number of questions before creating the guest system. The question are as follows:

i. *What is the name of your virtual machine and install location?*

ii. *How much RAM should be allocated (in megabytes)?*

iii. *What would you like to use as the disk (path)?*

iv. *Would you like to enable graphics support? (yes or no)*

The following transcript shows a typical virt-install session:

# virt-install

3. Once the guest system has been created, the vncviewer screen will appear containing the operating system installer:

## Installing a Xen Guest OS (Fedora Core 5)

1. Fedora Core 5 has a Xen guest installation script that simplifies the process, although it installs only FC5 guests. The script expects to access the FC5 install tree via FTP, the Web, or NFS; for some reason, you can't specify a directory or file.

# mkdir /var/www/html/dvd

# mount -t iso9660 /dev/dvd /var/www/html/dvd

# apachectl start

Now we'll run the installation script and answer its questions:

# xenguest-install.py

2. Xen does not start the guest operating system automatically. You need to type this command on the host:

# xm create guest1

**Course material**

3. To prove that both servers are running, try these commands:

# xm list

# xentop

4. To start Xen domains automatically, use these commands:

# /sbin/chkconfig --level 345 xendomains on

# /sbin/service xendomains start

5. To Edit A Xen Guest Configuration File, Which Is A Text File (Actually, A Python Script) In The /Etc/Xen Directory.

# man xmdomain.cfg

And edit as follows,

# Automatically generated Xen config file

name = "guest1"

memory = "256"

disk = [ 'file:/xenguest,xvda,w' ]

vif = [ 'mac=00:16:3e:63:c7:76' ]

uuid = "bc2c1684-c057-99ea-962b-de44a038bbda"

bootloader="/usr/bin/pygrub"

on_reboot = 'restart'

on_crash = 'restart'

6. Once you have a guest configuration file, create the Xen guest with this command:

where

# xm create -c guest_name

**guest_name** can be a full pathname or a relative filename (in which case Xen places it in /etc/xen/guest_name).

Xen will create the guest domain and try to boot it from the given file or device.

The **-c** option attaches a console to the domain when it starts, so you can answer the installation questions that appear.