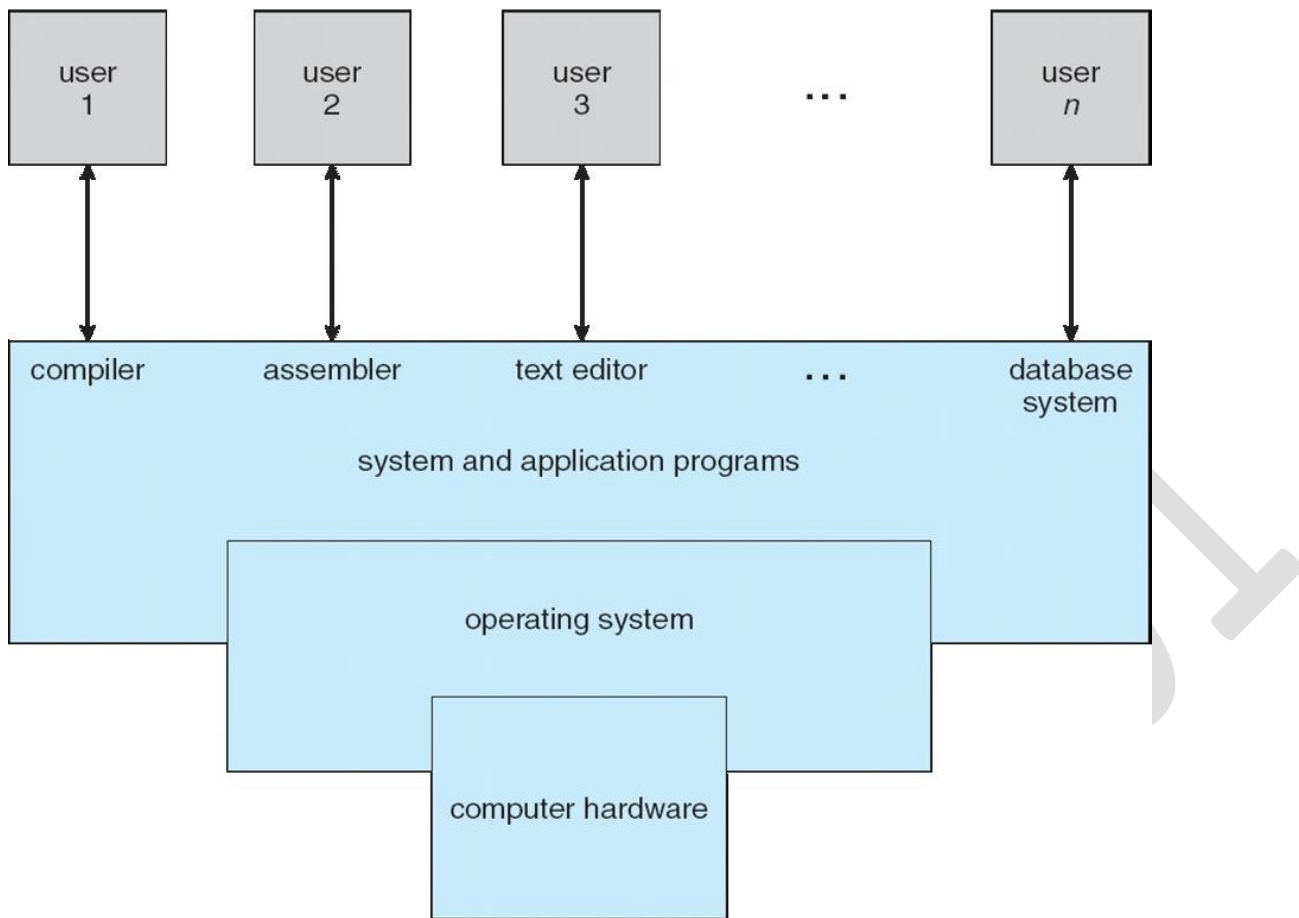# UNIT I

OS is 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

COMPUTER SYSTEM ARCHITECTURE

- Computer system can be divided into four components:

- Hardware – provides basic computing resources

  - CPU, memory, I/O devices

- Operating system

  - Controls and coordinates use of hardware among various applications and users

- Application programs – define the ways in which the system resources are used to solve the computing problems of the users

  - Word processors, compilers, web browsers, database systems, video games

- Users

  - People, machines, other computers

user 1    user 2    user 3    ...    user n

compiler    assembler    text editor    ...    database system

system and application programs

operating system

computer hardware

- Depends on the point of view

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

  o Don't care about **resource utilization**

- But shared computer such as **mainframe** or **minicomputer** must keep all users happy

- Users of dedicate systems such as **workstations** have dedicated resources but frequently use shared resources from **servers**

- Handheld computers are resource poor, optimized for usability and battery life

- Some computers have little or no user interface, such as embedded computers in devices and automobiles

- OS is a **resource allocator**

  o Manages all resources

  o Decides between conflicting requests for efficient and fair resource use

- OS is a **control program**

  o Controls execution of programs to prevent errors and improper use of the computer
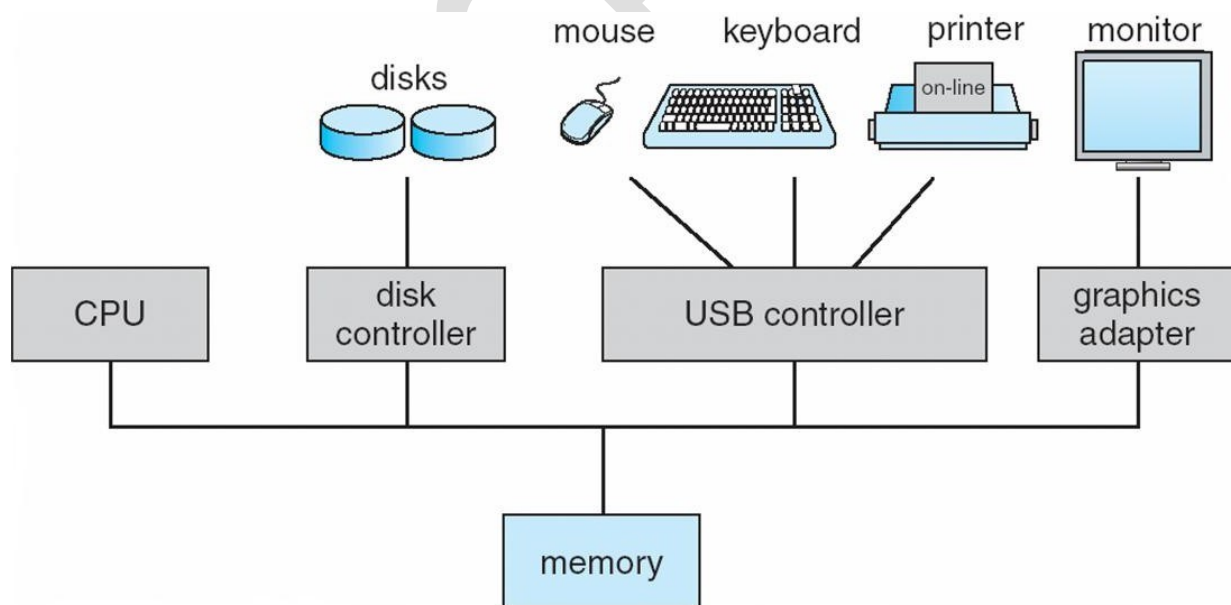
- No universally accepted definition

approximation

- o But varies wildly

- –The one program running at all times on the computer‖ is the **kernel**.

- Everything else is either

    - o a system program (ships with the operating system) , or

    - o an application program.

COMPUTER STARTUP

- **Bootstrap program** is loaded at power-up or reboot

    - o Typically stored in ROM or EPROM, generally known as **firmware**

    - o Initializes all aspects of system

    - o Loads operating system kernel and starts execution

COMPUTER-SYSTEM OPERATION

- One or more CPUs, device controllers connect through common bus providing access to shared memory



- I/O devices and the CPU can execute concurrently

- Each device controller is in charge of a particular device type

- Each device controller has a local buffer

- CPU moves data from/to main memory to/from local buffers

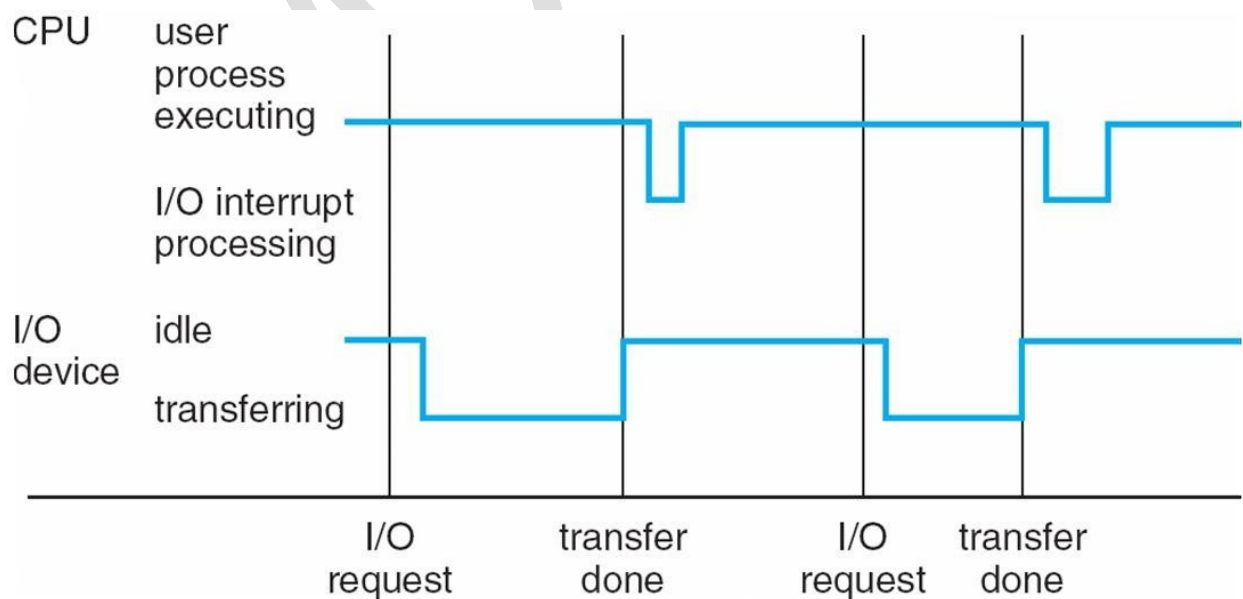- I/O is from the device to local buffer of controller

INTERRUPTS

- Device controller informs CPU that it has finished its operation by causing an interrupt

- Interrupt transfers control to the interrupt service routine generally, through the **interrupt vector**, which contains the addresses of all the service routines

- Interrupt architecture must save the address of the interrupted instruction

- A **trap** or **exception** is a software-generated interrupt caused either by an error or a user request

- An operating system is **interrupt driven**

INTERRUPT HANDLING

- The OS preserves the state of the CPU by storing registers and the program counter

- Determines which type of interrupt has occurred:

  o polling

    ▪ The interrupt controller polls (send a signal out to) each device to determine which one made the request

  o **vectored** interrupt system

- Separate segments of code determine what action should be taken for each type of interrupt

INTERRUPT TIMELINE

I/O STRUCTURE

- Synchronous (blocking) I/O

  o Waiting for I/O to complete

  o Easy to program, not always efficient

  o Wait instruction idles the CPU until the next interrupt

  o At most one I/O request is outstanding at a time

    ▪ no simultaneous I/O processing

- Asynchronous (nonblocking) I/O

  o After I/O starts, control returns to user program without waiting for I/O completion

  o Harder to program, more efficient

  o **System call** – request to the OS to allow user to wait for I/O completion (polling periodically to check busy/done)

  o **Device-status table** contains entry for each I/O device indicating its type, address, and state

STORAGE HIERARCHY

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

Computer storage, along with most computer throughput, is generally measured and manipulated in bytes and collections of bytes.

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

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

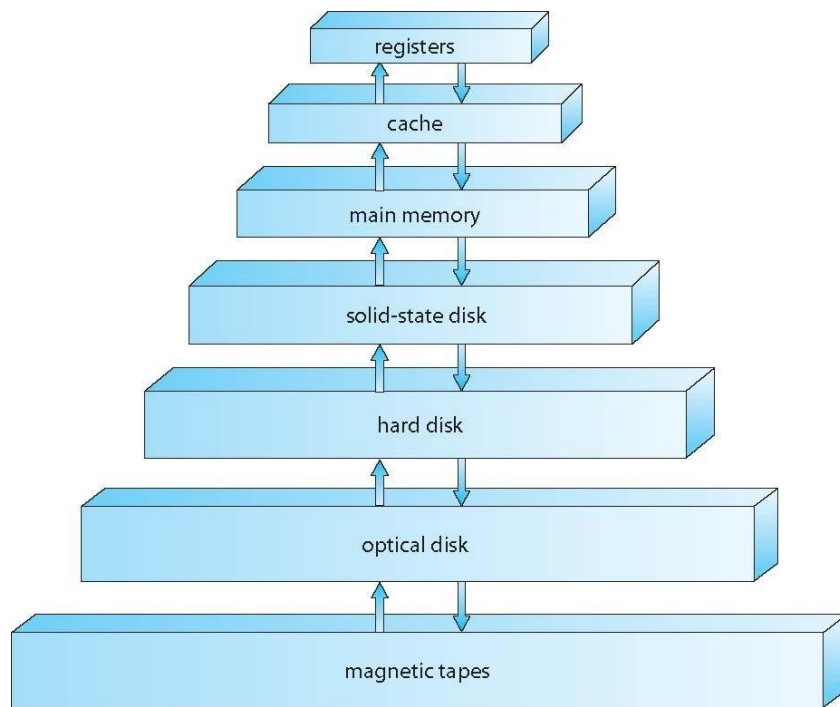a **gigabyte**, or **GB**, is $1,024^3$ bytes a

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

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

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).

STORAGE STRUCTURE

- Main memory – only large storage media that the CPU can access directly
    - Random access
    - Typically **volatile**
- Secondary storage – extension of main memory that provides large **nonvolatile** storage capacity
- Hard disks – rigid metal or glass platters covered with magnetic recording material
    - Disk surface is logically divided into **tracks**, which are subdivided into **sectors**
    - The **disk controller** determines the logical interaction between the device and the computer
- **Solid-state disks** – faster than hard disks, nonvolatile
    - Various technologies
    - Becoming more popular
- Storage systems organized in hierarchy
    - Speed
    - Cost (per byte of storage)
    - Volatility
- **Device Driver** for each device controller to manage I/O
    - Provides uniform interface between controller and kernel

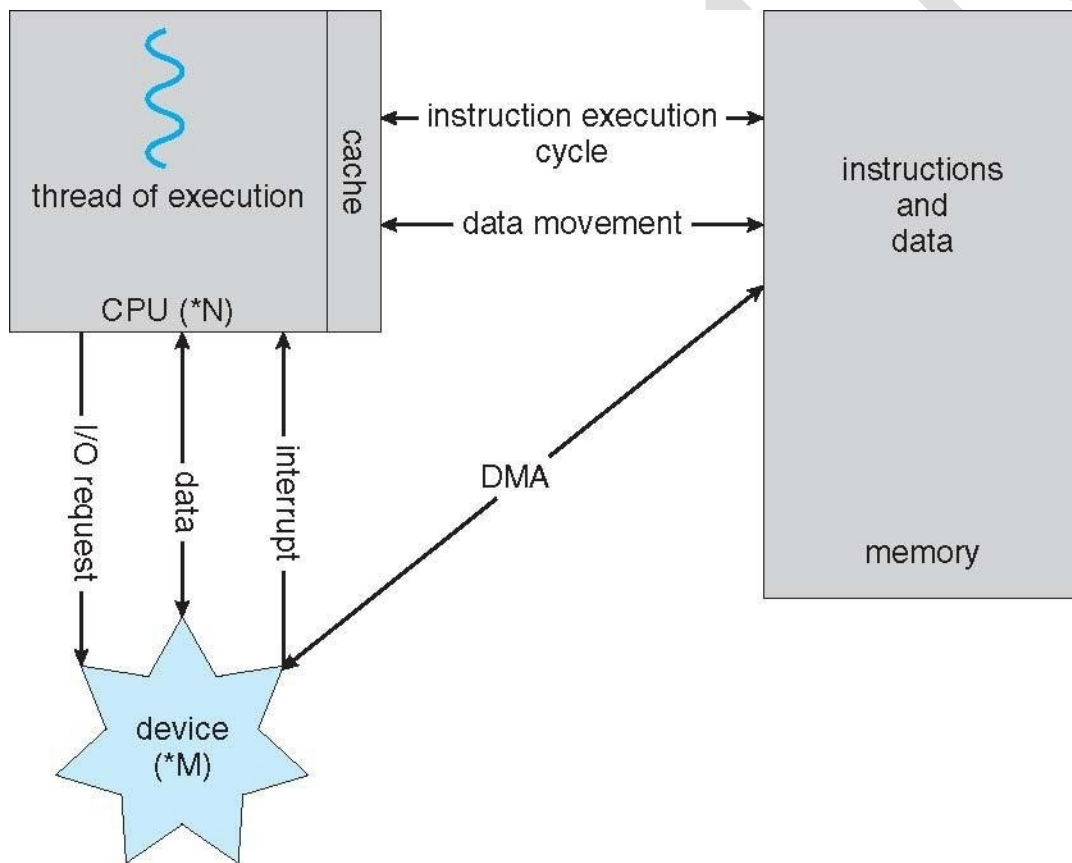| Level | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| Name | registers | cache | main memory | solid state disk | magnetic disk |
| Typical size | < 1 KB | < 16MB | < 64GB | < 1 TB | < 10 TB |
| Implementation technology | custom memory with multiple ports CMOS | on-chip or off-chip CMOS SRAM | CMOS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

CACHING

- Important principle

- Performed at many levels in a computer

    o in hardware,

    o operating system,

    o software

- Information in use copied from slower to faster storage temporarily

    o Efficiency

- Faster storage (cache) checked first to determine if information is there

- o   If it is, information used directly from the cache (fast)

- o   If not, data copied to cache and used there

- Cache smaller than storage being cached

  - o   Cache management important design problem

  - o   Cache size and replacement policy

DIRECT MEMORY ACCESS STRUCTURE

- Typically used for I/O devices that generate data in blocks, or generate data fast

- Device controller transfers blocks of data from buffer storage directly to main memory without CPU intervention

- Only one interrupt is generated per block, rather than the one interrupt per byte
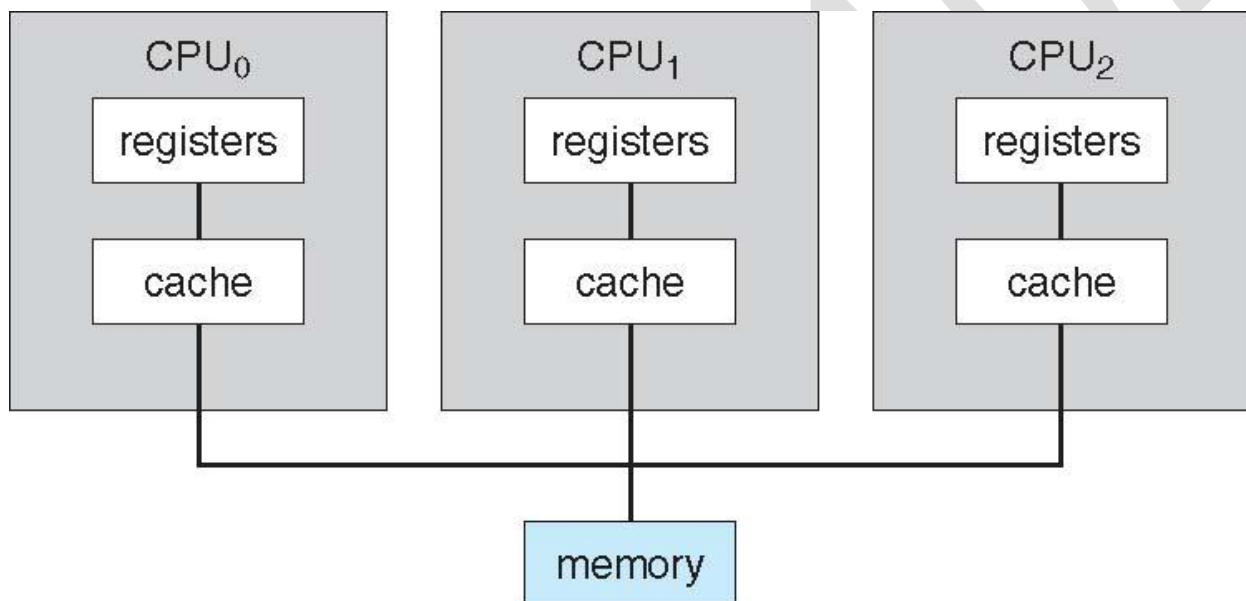
HOW A MODERN COMPUTER SYSTEM WORKS



TYPES OF SYSTEMS

- Most systems use a single general-purpose processor

  - o   Most systems have special-purpose processors as well

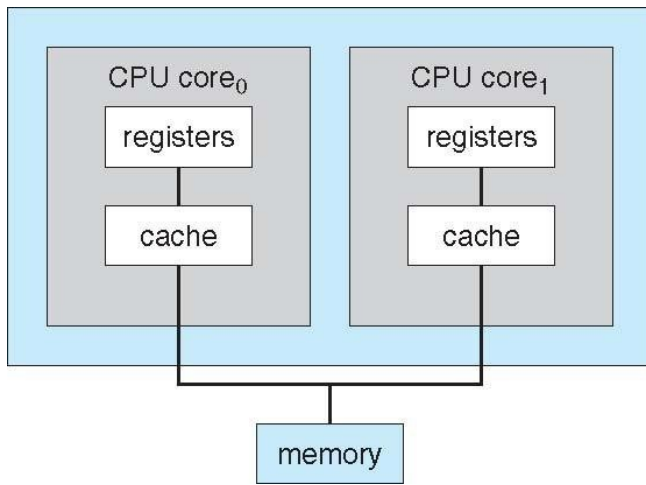- **Multiprocessors** systems growing in use and importance

o Advantages include:

1. Increased throughput

   **2. Economy of scale**

   3. **Increased reliability** – graceful degradation or fault tolerance

o Two types:

- **Asymmetric Multiprocessing** – each processor is assigned a specific task
- **Symmetric Multiprocessing** – each processor performs all tasks
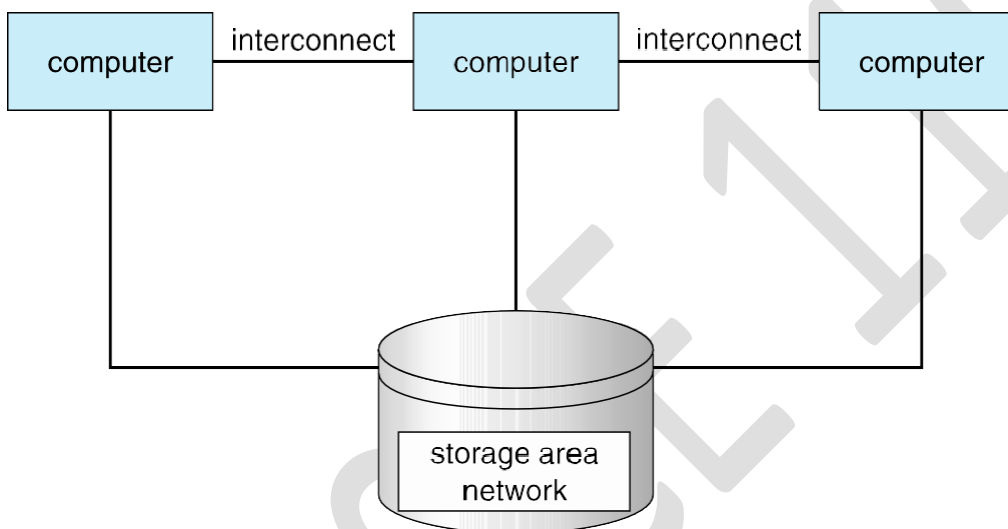
SYMMETRIC MULTIPROCESSING ARCHITECTURE



DUAL CORE DESIGN

- Multicore

  o Several cores on a single chip

  o On chip communication is faster than between-chip

  o Less power used

CLUSTERED SYSTEMS



- Like multiprocessor systems, but multiple systems working together

- Provides a **high-availability** service which survives failures

- **Asymmetric clustering** has one machine in hot-standby mode

- **Symmetric clustering** has multiple nodes running applications, monitoring each other

- Some clusters are for **high-performance computing (HPC)**

- Applications must be written to use **parallelization**

MEMORY LAYOUT FOR MULTIPROGRAMMED SYSTEMS



OPERATING SYSTEM OPERATIONS

- **Interrupt driven** (hardware and software)
    - Hardware interrupt by one of the devices

o Software interrupt (**exception** or **trap**):

  ▪ Software error (e.g., division by zero)

  ▪ Request for operating system service

  ▪ Other process problems include infinite loop, processes modifying each other or the operating system
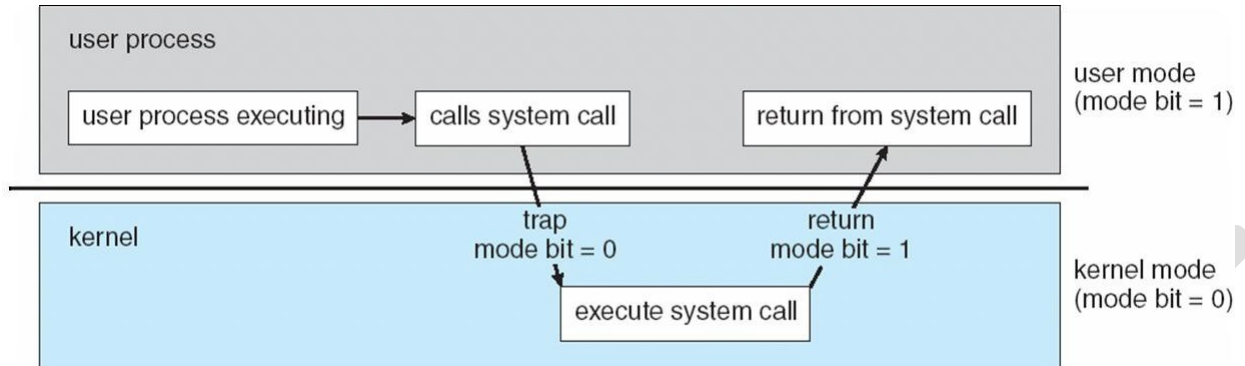


- **Dual-mode** operation allows OS to protect itself and other system components

- User mode and kernel mode

  - **Mode bit** provided by hardware

  - Provides ability to distinguish when system is running user code or kernel code

  - Some instructions designated as **privileged**, only executable in kernel mode

  - System call changes mode to kernel, return from call resets it to user

  - Timer to prevent infinite loop / process hogging resources

  - Timer is set to interrupt the computer after some time period

  - Keep a counter that is decremented by the physical clock

  - Operating system set the counter (privileged instruction)

  - When counter zero generate an interrupt

  - Set up before scheduling process to

  - regain control, or

  - terminate program that exceeds allotted time

PROCESS MANAGEMENT

- A process is a program in execution. It is a unit of work within the system. Program is a *passive entity*, process is an *active entity*.

- o Software interrupt (**exception** or **trap):**
- Process needs resources to accomplish its task

- o CPU, memory, I/O, files

- o Initialization data

- Typically system has many processes, some user, some operating system running concurrently on one or more CPUs

  - o Concurrency by multiplexing the CPUs among the processes / threads

ACTIVITIES

- Creating and deleting both user and system processes

- Suspending and resuming processes

- Providing mechanisms for process synchronization

- Providing mechanisms for process communication

- Providing mechanisms for deadlock handling

MEMORY MANAGEMENT

- To execute a program all (or part) of the instructions must be in memory

- All (or part) of the data that is needed by the program must be in memory.

- Memory management determines what is in memory and when

  - o Optimizing CPU utilization and computer response to users

- Memory management activities

  - o Keeping track of which parts of memory are currently being used and by whom

  - o Deciding which processes (or parts thereof) and data to move into and out of memory

  - o Allocating and deallocating memory space as needed

STORAGE MANAGEMENT

- OS provides uniform, logical view of information storage

  - o Different devices, same view

  - o Abstracts physical properties to logical storage unit - **file**

  - o Each medium is controlled by device (i.e., disk drive, tape drive)

    - ▪ Varying properties include access speed, capacity, data-transfer rate, access method (sequential or random)

    - ▪ File-System management

- o Files usually organized into directories

- o Access control on most systems to determine who can access what

- o OS activities include

  - Creating and deleting files and directories

  - Primitives to manipulate files and directories

  - Mapping files onto secondary storage

  - Backup files onto stable (non-volatile) storage media

## MASS STORAGE MANAGEMENT

- Usually disks used to store

- data that does not fit in main memory, or

- data that must be kept for a –long‖ period of time

- Proper management is of central importance

- Entire speed of computer operation hinges on disk subsystem and its algorithms

- Disk is slow, its I/O is often a bottleneck

- OS activities

- Free-space management

- Storage allocation

- Disk scheduling

- Multitasking environments must be careful to use most recent value, no matter where it is stored in the storage hierarchy

- Multiprocessor environment must provide **cache coherency** in hardware such that all CPUs have the most recent value in their cache

- Distributed environment situation even more complex

- Several copies of a datum can exist

## I/O SUBSYSTEM

- One purpose of OS is to hide peculiarities of hardware devices from the user

- I/O subsystem responsible for

  - o Memory management of I/O including buffering (storing data temporarily while it is being transferred),

- o   General device-driver interface

- o   Drivers for specific hardware devices

- **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

  - o   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

  - o   Access control for users and groups

COMPUTING ENVIRONMENTS

TRADITIONAL

- Stand-alone general purpose machines

- But blurred as most systems interconnect with others (i.e., the Internet)

- **Portals** provide web access to internal systems

- **Network computers** (**thin clients**) are like Web terminals

- Mobile computers interconnect via **wireless networks**

- Networking becoming ubiquitous – even home systems use **firewalls** to protect home computers from Internet attacks

MOBILE

- Handheld smartphones, tablets, etc

- What is the functional difference between them and a ‑traditional‖ laptop?

- Extra feature – more OS features (GPS, gyroscope)

- Allows new types of apps like *augmented reality*

- Use IEEE 802.11 wireless, or cellular data networks for connectivity

- Leaders are **Apple iOS** and **Google Android**


- Distributed computing

- Collection of separate, possibly heterogeneous, systems networked together

- **Network** is a communications path, **TCP/IP** most common

- Local Area Network (LAN)

- **Wide Area Network (WAN)**

- Metropolitan Area Network (MAN)

- **Personal Area Network (PAN)**

- **Network Operating System** provides features between systems across network

- Communication scheme allows systems to exchange messages

- Illusion of a single system

- Client-Server Computing

- Dumb terminals supplanted by smart PCs

- Many systems now **servers**, responding to requests generated by **clients**

- **Compute-server system** provides an interface to client to request services (i.e., database)

- **File-server system** provides interface for clients to store and retrieve files

  P2P does not distinguish clients and servers
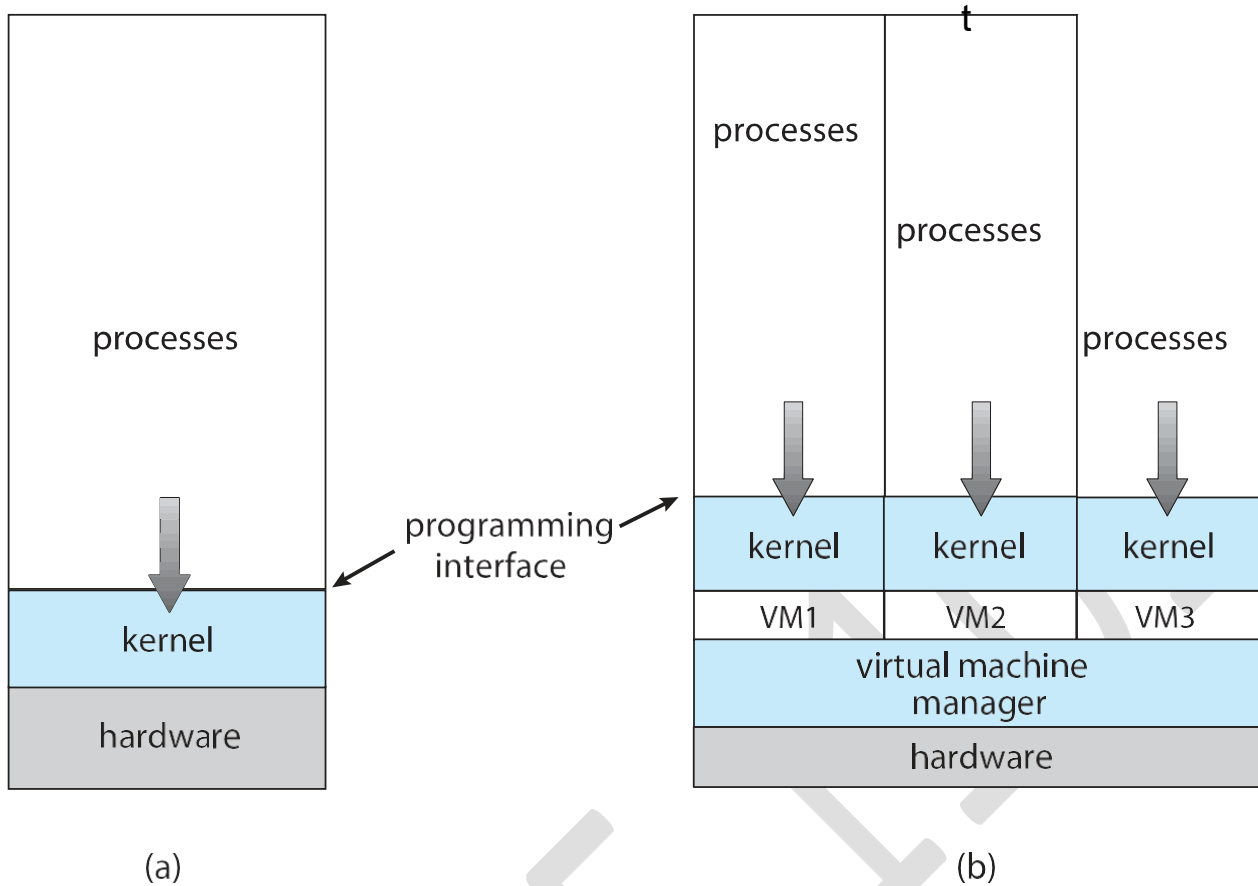
  - o Instead all nodes are considered peers

  - o May each act as client, server or both

  - o Node must join P2P network

    - Registers its service with central lookup service on network, or

    - Broadcast request for service and respond to requests for service via *discovery protocol*

  - o Examples include Napster and Gnutella, **Voice over IP (VoIP)** such as Skype

Virtualization

- **Host** OS, natively compiled for CPU

- **VMM** - virtual machine manager

- Creates and runs virtual machines

- VMM runs **guest** OSes, also natively compiled for CPU

- Applications run within these guest OSes

- Example: Parallels for OS X running Win and/or Linux and their apps

- Some VMM'es run within a host OS

- But, some act as a specialized OS

- Example. VMware ESX: installed on hardware, runs when hardware boots, provides services to apps, runs guest OSes

- Vast and growing industry

- Use cases

- Developing apps for multiple different OSes on 1 PC

- Very important for **cloud computing**

- Executing and managing **compute environments** in data centers

(a)                                    (b)

- Operating systems made available in source-code format rather than just binary closed-source

- Counter to the **copy protection** and **Digital Rights Management (DRM)** movement

- Started by **Free Software Foundation (FSF)**, which has ―copyleft‖ **GNU Public License (GPL)**

- Examples include **GNU/Linux** and **BSD UNIX** (including core of **Mac OS X**), and many more

- Can use VMM like VMware Player (Free on Windows), Virtualbox (open source and free on many platforms - http://www.virtualbox.com)

- Use to run guest operating systems for exploration

OPERATING SYSTEM SERVICES

- Operating System Services

- User Operating System Interface

- System Calls

- Types of System Calls

- System Programs

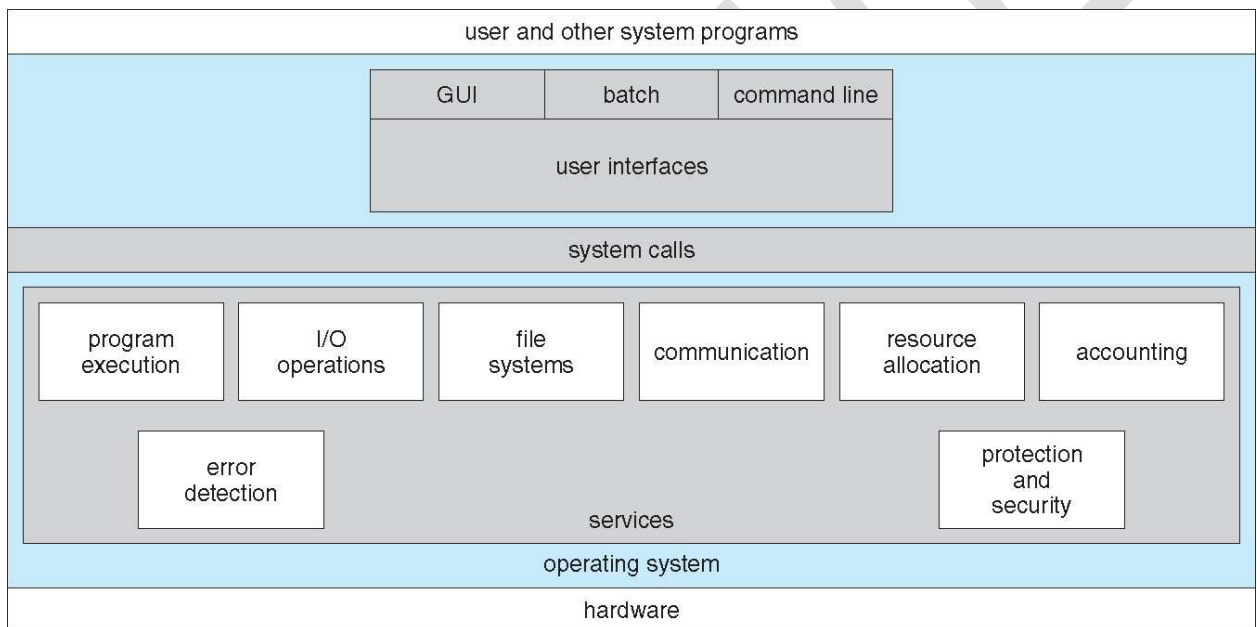- Operating System Design and Implementation

- Operating System Structure

- Operating System Debugging

- Operating System Generation

- System Boot

- Operating systems provide an environment for execution of programs and services (helpful functions) to programs and users

- User services:

    - **User interface**

        - No UI, Command-Line (CLI), Graphics User Interface (GUI), Batch

    - **Program execution** - Loading a program into memory and running it, end execution, either normally or abnormally (indicating error)

    - **I/O operations** - A running program may require I/O, which may involve a file or an I/O device

    - User services (Cont.):

    - **File-system manipulation** - Programs need to read and write files and directories, create and delete them, search them, list file Information, permission management.

    - **Communications** – Processes may exchange information, on the same computer or between computers over a network

        - Communications may be via shared memory or through message passing (packets moved by the OS)

    - **Error detection** – OS needs to be constantly aware of possible errors

        - May occur in the CPU and memory hardware, in I/O devices, in user program

        - For each type of error, OS should take the appropriate action to ensure correct and consistent computing

        - Debugging facilities can greatly enhance the user's and programmer's abilities to efficiently use the system
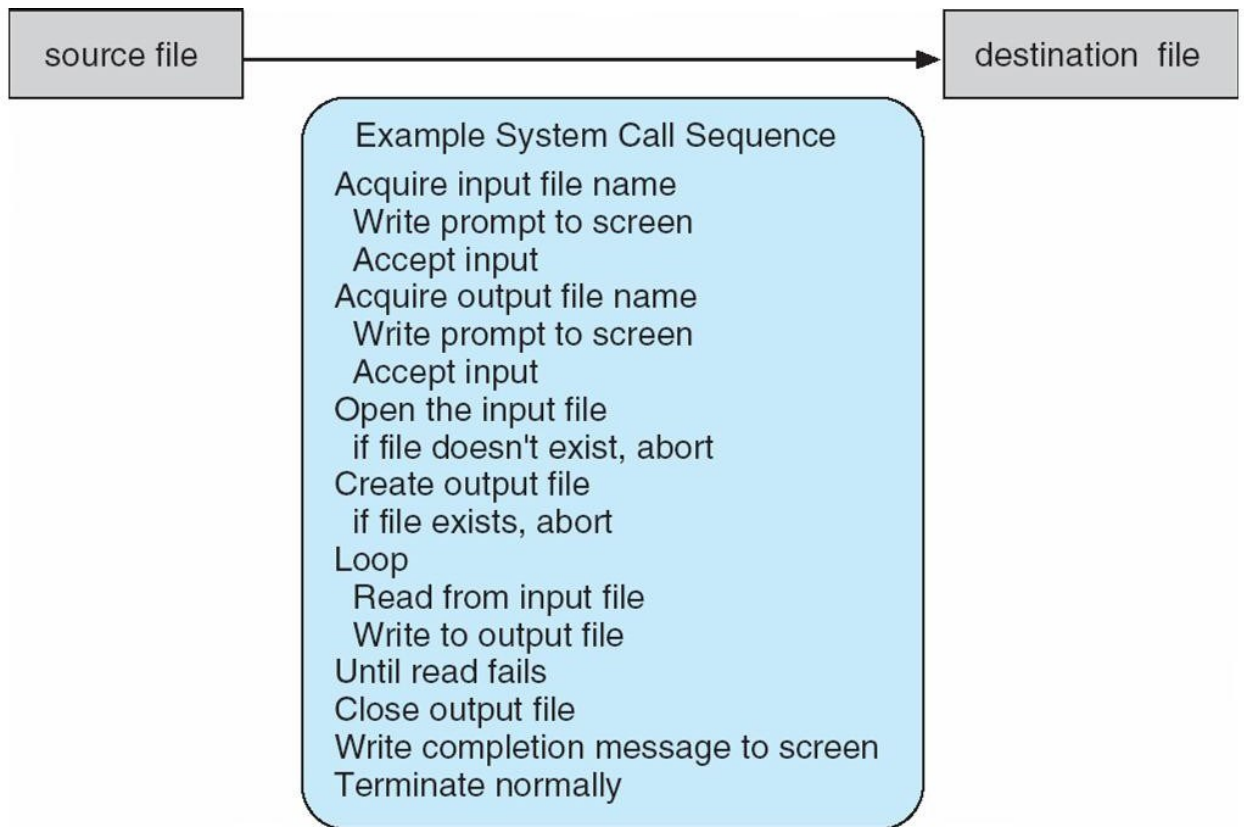
System services:

    - For ensuring the efficient operation of the system itself via resource sharing

    - **Resource allocation -** When multiple users or multiple jobs running concurrently, resources must be allocated to each of them

- devices.

  o **Accounting -** To keep track of which users use how much and what kinds of computer resources

  o **Protection and security -** The owners of information stored in a multiuser or networked computer system may want to control use of that information, concurrent processes should not interfere with each other

    ▪ **Protection** involves ensuring that all access to system resources is controlled

    ▪ **Security** of the system from outsiders requires user authentication, extends to defending external I/O devices from invalid access attempts

- 

| user and other system programs |
|---|

| GUI | batch | command line |
|---|---|---|
| user interfaces | | |

| system calls |
|---|

| program execution | I/O operations | file systems | communication | resource allocation | accounting |
|---|---|---|---|---|---|

| error detection | | | services | | protection and security |

| operating system |
|---|

| hardware |
|---|

- **CLI** or **command interpreter** allows direct command entry

  o Sometimes implemented in kernel, sometimes by systems program

  o Sometimes multiple flavors implemented – **shells**

  o Primarily fetches a command from user and executes it

  o Sometimes commands built-in, sometimes just names of programs

    ▪ If the latter, adding new features doesn't require shell modification

- Systems calls: 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 Programming Interface (API)** rather than direct system call use

- Three most common APIs are

- o   Win32 API for Windows,

- o   POSIX API for POSIX-based systems

    - ▪   including virtually all versions of UNIX, Linux, and Mac OS X,

- o   Java API for the Java virtual machine (JVM)

| source file | → | destination  file |
|---|---|---|

```
    Example System Call Sequence
Acquire input file name
  Write prompt to screen
  Accept input
Acquire output file name
  Write prompt to screen
  Accept input
Open the input file
  if file doesn't exist, abort
Create output file
  if file exists, abort
Loop
  Read from input file
  Write to output file
Until read fails
Close output file
Write completion message to screen
Terminate normally
```
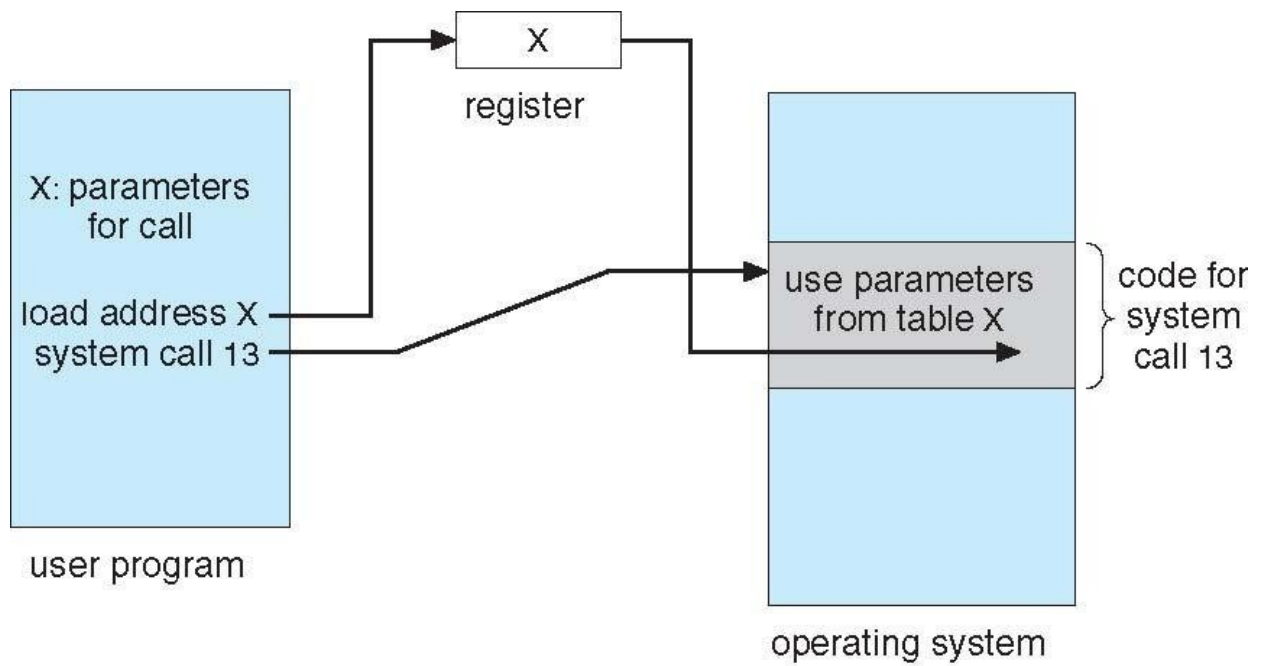
- •   

- •   Typically, a number associated with each system call

    - o   **System-call interface** maintains a table indexed according to these numbers

- •   The system call interface invokes the intended system call in OS kernel and returns status of the system call and any return values

- •   The caller need know nothing about how the system call is implemented

    - o   Just needs to obey API and understand what OS will do as a result call

    - o   Most details of OS interface hidden from programmer by API

        - ▪   Managed by run-time support library (set of functions built into libraries included with compiler)

- 

- Three general methods used to pass parameters to the OS in system calls

  - o Simplest: in registers

    - ▪ In some cases, may be more parameters than registers

  - o Parameters stored in a block, or table, in memory, and address of block passed as a parameter in a register

    - ▪ This approach taken by Linux and Solaris

  - o Parameters placed, or **pushed**, onto the **stack** by the program and **popped** off the stack by the operating system

  - o Block and stack methods do not limit the number or length of parameters being passed

TYPES OF SYSTEM CALLS

- Process control

    - create process, terminate process

    - end, abort

    - load, execute

    - get process attributes, set process attributes

    - wait for time

    - wait event, signal event

    - allocate and free memory

    - Dump memory if error

    - **Debugger** for determining **bugs, single step** execution

    - **Locks** for managing access to shared data between processes

- File management

    - create file, delete file

    - open, close file

    - read, write, reposition

    - get and set file attributes

- Device management
    - request device, release device

- o read, write, reposition

- o get device attributes, set device attributes

- o logically attach or detach devices

- Information maintenance

  - o get time or date, set time or date

  - o get system data, set system data

  - o get and set process, file, or device attributes

- Communications

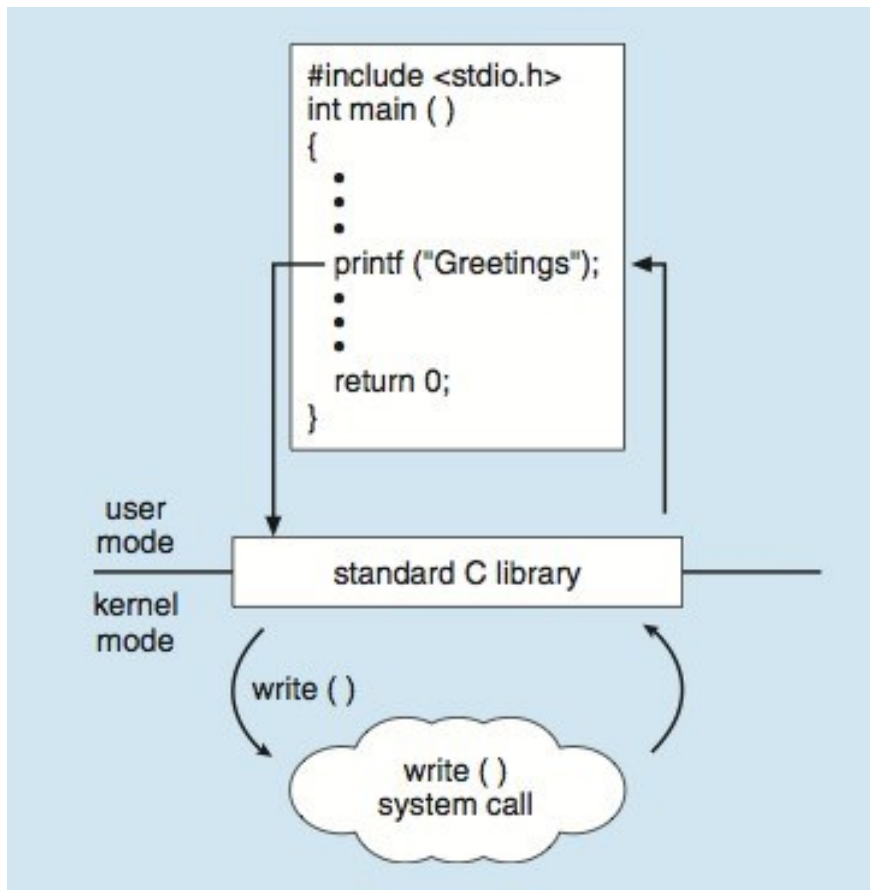  - o create, delete communication connection

- send, receive messages if **message passing model** to **host name** or **process name**
  - From **client** to **server**
- **Shared-memory model** create and gain access to memory regions
- transfer status information
- attach and detach remote devices

- Protection
  - Control access to resources
  - Get and set permissions
  - Allow and deny user access

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

•

```
#include <stdio.h>
int main ( )
{
    .
    .
    .
    printf ("Greetings");
    .
    .
    .
    return 0;
}
```

user
mode
kernel
mode

standard C library

write ( )

write ( )
system call

SYSTEM PROGRAMS

- System programs provide a convenient environment for program development and execution.

- Most users' view of the operation system is defined by system programs, not the actual system calls

- They can be divided into:

  - File manipulation

    - rm, ls, cp, mv, etc in Unix

  - Status information sometimes stored in a File modification

  - Programming language support

  - Program loading and execution

  - Communications

  - Background services

  - Application programs

- Provide a convenient environment for program development and execution

  - Some of them are simply user interfaces to system calls; others are considerably more complex

  - **File management** - Create, delete, copy, rename, print, dump, list, and generally manipulate files and directories

- Status information

  - Some ask the system for info - date, time, amount of available memory, disk space, number of users

  - Others provide detailed performance, logging, and debugging information

  - Typically, these programs format and print the output to the terminal or other output devices

  - Some systems implement a **registry** - used to store and retrieve configuration information

- File modification

  - Text editors to create and modify files

text

- **Programming-language support** - Compilers, assemblers, debuggers and interpreters sometimes provided

- **Program loading and execution**- Absolute loaders, relocatable loaders, linkage editors, and overlay-loaders, debugging systems for higher-level and machine language

- **Communications** - Provide the mechanism for creating virtual connections among processes, users, and computer systems

  - Allow users to send messages to one another's screens, browse web pages, send electronic-mail messages, log in remotely, transfer files from one machine to another

- Background Services

  - Launch at boot time

    - Some for system startup, then terminate

    - Some from system boot to shutdown

  - Provide facilities like disk checking, process scheduling, error logging, printing

  - Run in user context not kernel context

  - Known as **services**, **subsystems**, **daemons**

- Application programs

  - Don't pertain to system

  - Run by users

  - Not typically considered part of OS

  - Launched by command line, mouse click, finger poke

- Design and Implementation of OS not ‒solvable‖, but some approaches have proven successful

- Internal structure of different Operating Systems can vary widely

- Start the design by defining goals and specifications

- Affected by choice of hardware, type of system

- **User** goals and **System** goals

  - User goals – operating system should be convenient to use, easy to learn, reliable, safe, and fast

maintain, as well as flexible, reliable, error-free, and efficient

- Important principle to separate

- **Policy:** *What* will be done?
  **Mechanism:** *How* to do it?

- Mechanisms determine how to do something, policies decide what will be done

- The separation of policy from mechanism is a very important principle, it allows maximum flexibility if policy decisions are to be changed later (example – timer)

- Specifying and designing an OS is highly creative task of **software engineering**

- Much variation

  o Early OSes in assembly language

  o Then system programming languages like Algol, PL/1

  o Now C, C++

- Actually usually a mix of languages

  o Lowest levels in assembly

  o Main body in C

  o Systems programs in C, C++, scripting languages like PERL, Python, shell scripts

- More high-level language easier to **port** to other hardware

  o But slower

- General-purpose OS is a **very large program**

  o How to implement and structure it?

  o Can apply many ideas from software engineering

    ▪ Software engineering - a separate area in CS

    ▪ Studies design, development, and maintenance of software

- A common approach is to partition OS into modules/components

  o Each modules is responsible for one (or several) aspect of the desired functionality

  o Each module has carefully defined interfaces

  o Advantages:

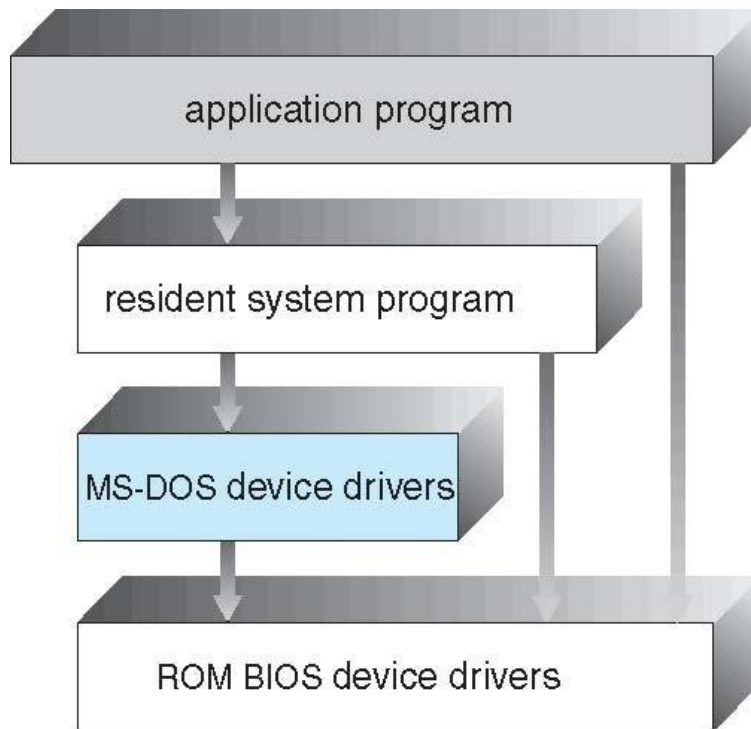    ▪ Traditional advantages of modular programming:

etc

- Modules can be developed independently from each other

- o Disadvantages:
  - Efficiency can decrease (vs monolithic approach)
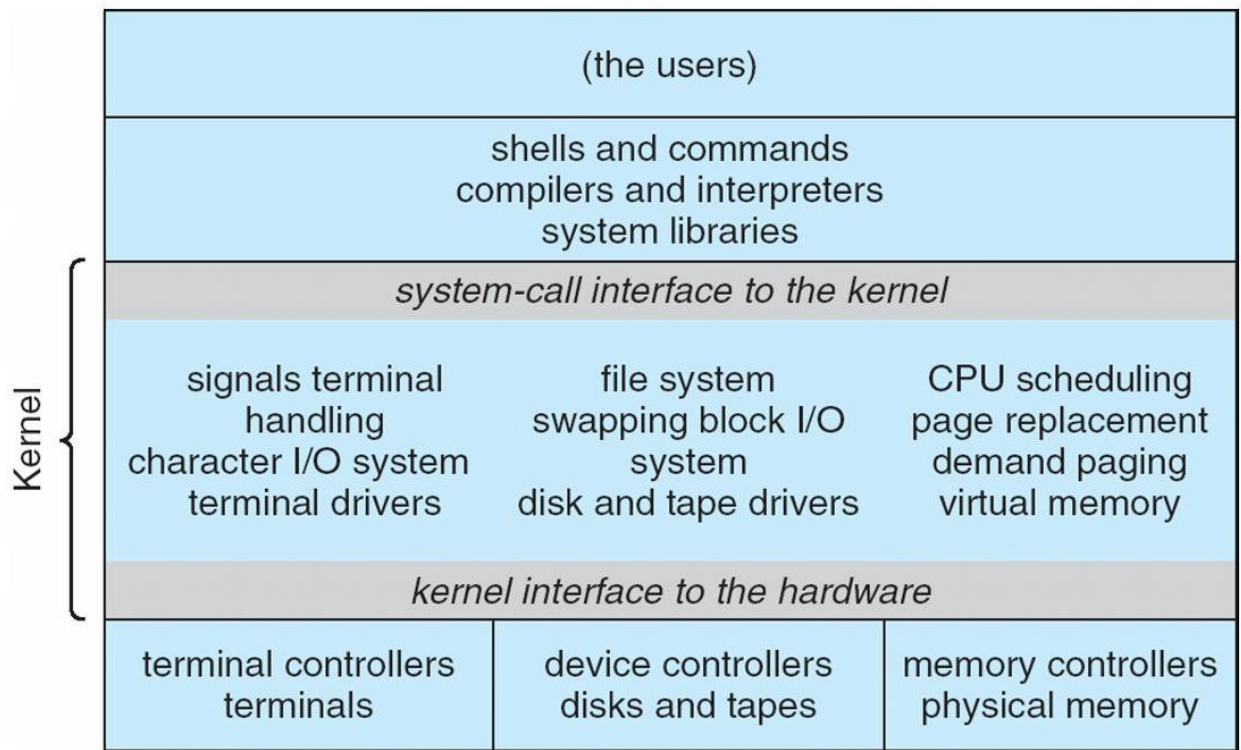
## OPERATING SYSTEM STRUCTURES

- In general, various ways are used to structure OSes

- Many OS'es don't have well-defined structures

  - o Not one pure model: Hybrid systems

  - o Combine multiple approaches to address performance, security, usability needs

- Simple structure – MS-DOS

- More complex structure - UNIX

- Layered OSes

- Microkernel OSes

## SIMPLE STRUCTURE

- MS-DOS was created to provide the most functionality in the least space

- Not divided into modules

- MS-DOS has some structure

  - o But its interfaces and levels of functionality are not well separated

- No dual mode existed for Intel 8088

  - o MS-DOS was developed for 8088

  - o Direct access to hardware is allowed
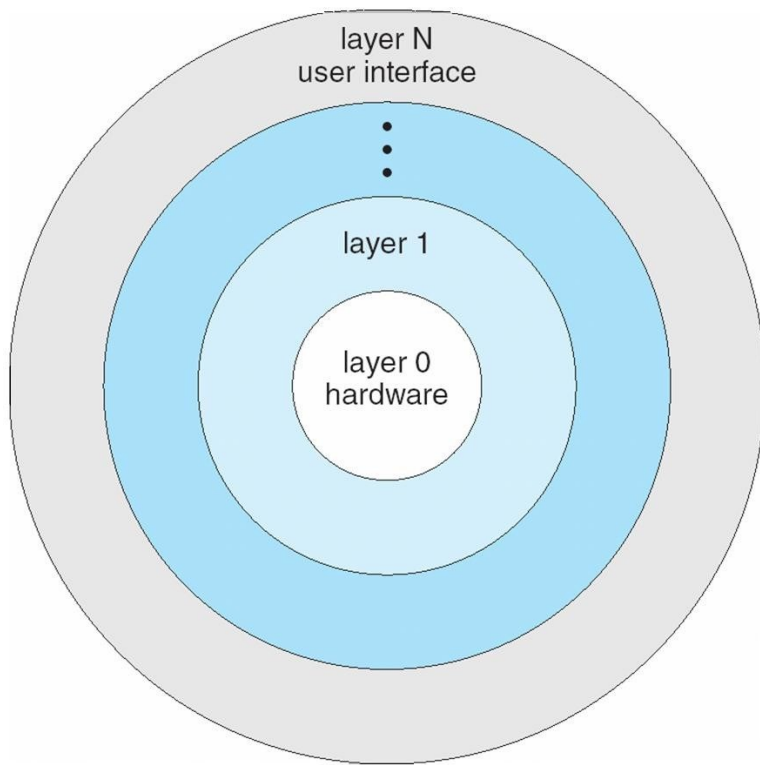
    - System crashes possible

- 

- Traditional UNIX has limited structuring

- UNIX consists of 2 separable parts:

  o Systems programs

  o Kernel

- UNIX Kernel

  o Consists of everything that is

    ▪ below the system-call interface and

    ▪ above the physical hardware

  o Kernel provides

    ▪ File system, CPU scheduling, memory management, and other operating-system functions

    ▪ This is a lot of functionality for just 1 layer

    ▪ Rather monolithic

      - But fast -- due to lack of overhead in communication inside kernel

|                                                                   |
|:-----------------------------------------------------------------:|
| (the users)                                                       |
| shells and commands<br>compilers and interpreters<br>system libraries |
| *system-call interface to the kernel*                             |
| signals terminal handling / character I/O system / terminal drivers  —  file system / swapping block I/O system / disk and tape drivers  —  CPU scheduling / page replacement / demand paging / virtual memory |
| *kernel interface to the hardware*                                |
| terminal controllers / terminals  —  device controllers / disks and tapes  —  memory controllers / physical memory |

(Kernel brace spans the system-call interface through kernel interface to the hardware rows.)

- 

  LAYERED APPROACH

- One way to make OS modular – layered approach

- The OS is divided into a number of layers (levels)

- Each layer is built on top of lower layers

- The bottom layer (layer 0), is the hardware

- The highest (layer N) is the user interface

- Layers are selected such that each uses functions (operations) and services of only lower-level layers

- Advantages:

  o simplicity of construction and debugging

- Disadvantages:

  o can be hard to decide how to split functionality into layers

  o less efficient due to high overhead

- 

OPERATING SYSTEM DEBUGGING

- **Debugging** is finding and fixing errors, or **bugs**

- OS generate **log files** containing error information

- Failure of an application can generate **core dump** file capturing memory of the process

- Operating system failure can generate **crash dump** file containing kernel memory

- Beyond crashes, performance tuning can optimize system performance

    o Sometimes using *trace listings* of activities, recorded for analysis

    o **Profiling** is periodic sampling of instruction pointer to look for statistical trends

- Kernighan's Law: ‖Debugging is twice as hard as writing the code in the first place. Therefore, if you write the code as cleverly as possible, you are, by definition, not smart enough to debug it.‖

PERFORMANCE TUNNING

- Improve performance by removing bottlenecks

- OS must provide means of computing and displaying measures of system behavior

- For example, ‖top‖ program or Windows Task Manager

OS SYSGEN

- Operating systems are designed to run on any of a class of machines

- The process of configuration is known as system generation **SYSGEN**

    o How to format partitions

    o Which hardware is present

    o Etc

- Used to build system-specific compiled kernel or system-tuned

- Can general more efficient code than one

general kernelSYSTEM BOOT

- How OS is loaded?

- When power is initialized on system, execution starts at a predefined memory location

    o Firmware ROM is used to hold initial bootstrap program (=bootstrap loader)

- Bootstrap loader

    o small piece of code

    o locates the kernel, loads it into memory, and starts it

- Sometimes 2-step process is used instead

    o Simple bootstrap loader in ROM loads a more complex boot program from (afixed location on) disk
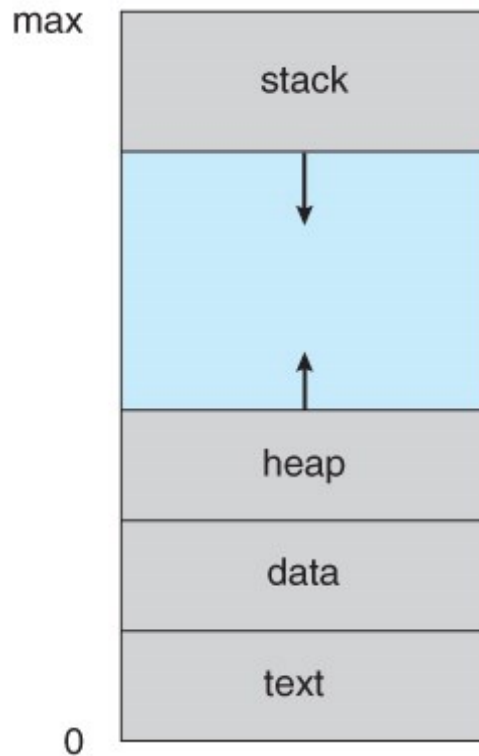
    o This more complex loader loads the kernel

# UNIT II

## 3.1 Process Concept

- A process is an instance of a program in execution.
- Batch systems work in terms of "jobs". Many modern process concepts are still expressed in terms of jobs, ( e.g. job scheduling ), and the two terms are often used interchangeably.
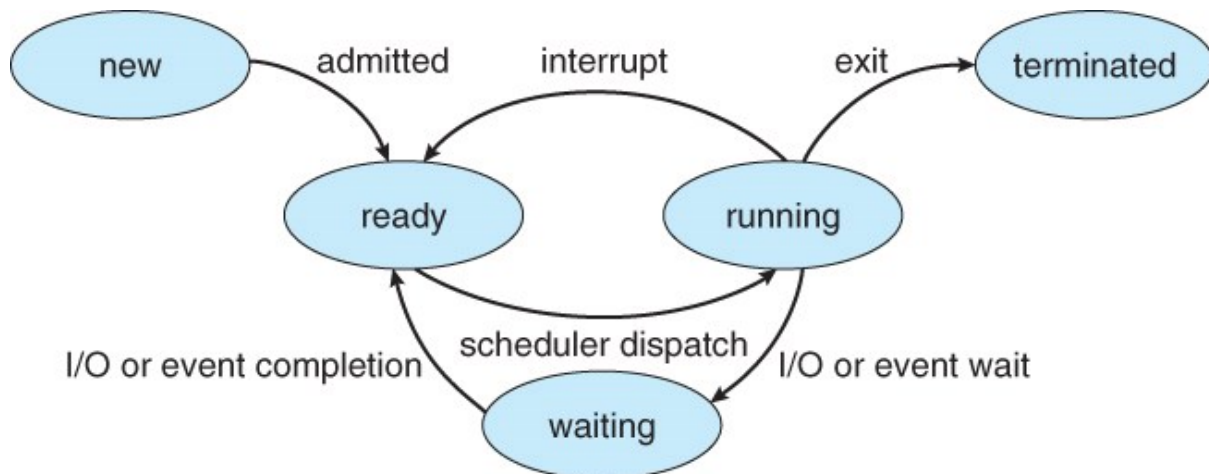
### 3.1.1 The Process

- Process memory is divided into four sections as shown in Figure 3.1 below:
  - The text section comprises the compiled program code, read in from non-volatile storage when the program is launched.
  - The data section stores global and static variables, allocated and initialized prior to executing main.
  - The heap is used for dynamic memory allocation, and is managed via calls to new, delete, malloc, free, etc.
  - The stack is used for local variables. Space on the stack is reserved for local variables when they are declared ( at function entrance or elsewhere, depending on the language ), and the space is freed up when the variables go out of scope. Note that the stack is also used for function return values, and the exact mechanisms of stack management may be language specific.
  - Note that the stack and the heap start at opposite ends of the process's free space and grow towards each other. If they should ever meet, then either a stack overflow error will occur, or else a call to new or malloc will fail due to insufficient memory available.
- When processes are swapped out of memory and later restored, additional information must also be stored and restored. Key among them are the program counter and the value of all program registers.

**Figure 3.1 - A process in memory**

### 3.1.2 Process State

- Processes may be in one of 5 states, as shown in Figure 3.2 below.
    - **New** - The process is in the stage of being created.
    - **Ready** - The process has all the resources available that it needs to run, but the CPU is not currently working on this process's instructions.
    - **Running** - The CPU is working on this process's instructions.
    - **Waiting** - The process cannot run at the moment, because it is waiting for some resource to become available or for some event to occur. For example the process may be waiting for keyboard input, disk access request, inter-process messages, a timer to go off, or a child process to finish.
    - **Terminated -** The process has completed.
- The load average reported by the "w" command indicate the average number of processes in the "Ready" state over the last 1, 5, and 15 minutes, i.e. processes who have everything they need to run but cannot because the CPU is busy doing something else.
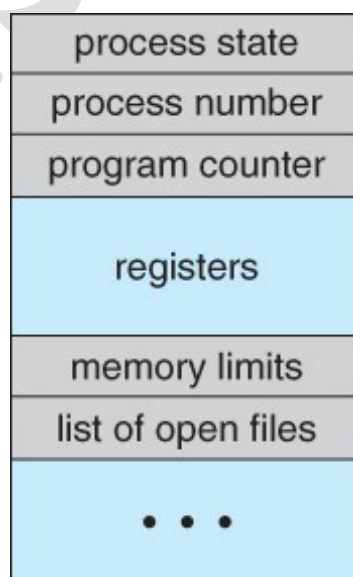- Some systems may have other states besides the ones listed here.

**Figure 3.2 - Diagram of process state**
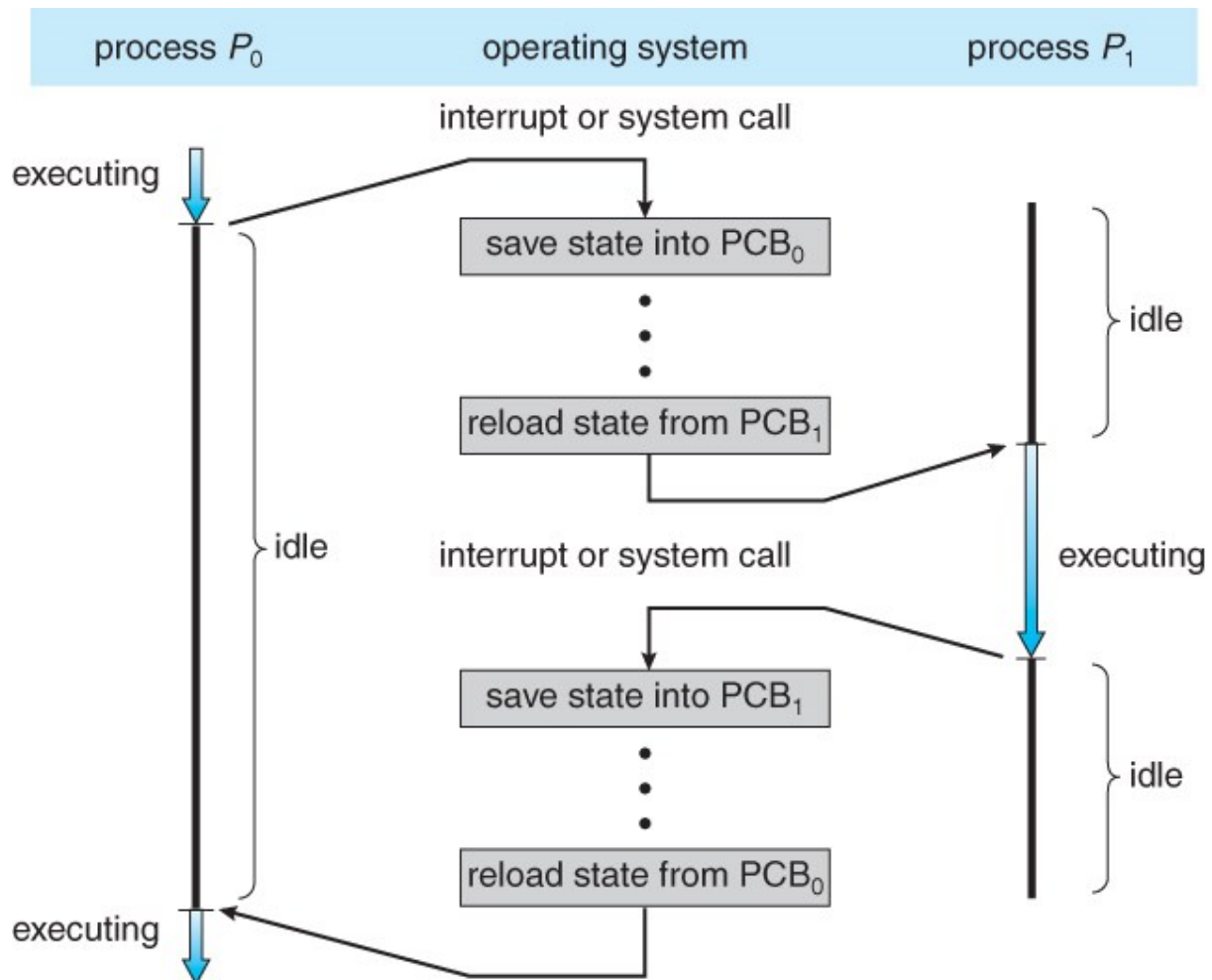
### 3.1.3 Process Control Block

For each process there is a Process Control Block, PCB, which stores the following ( types of ) process-specific information, as illustrated in Figure 3.1. ( Specific details may vary from system to system. )

- **Process State** - Running, waiting, etc., as discussed above.
- **Process ID**, and parent process ID.
- **CPU registers and Program Counter** - These need to be saved and restored when swapping processes in and out of the CPU.
- **CPU-Scheduling information** - Such as priority information and pointers to scheduling queues.
- **Memory-Management information** - E.g. page tables or segment tables.
- **Accounting information** - user and kernel CPU time consumed, account numbers, limits, etc.
- **I/O Status information** - Devices allocated, open file tables, etc.

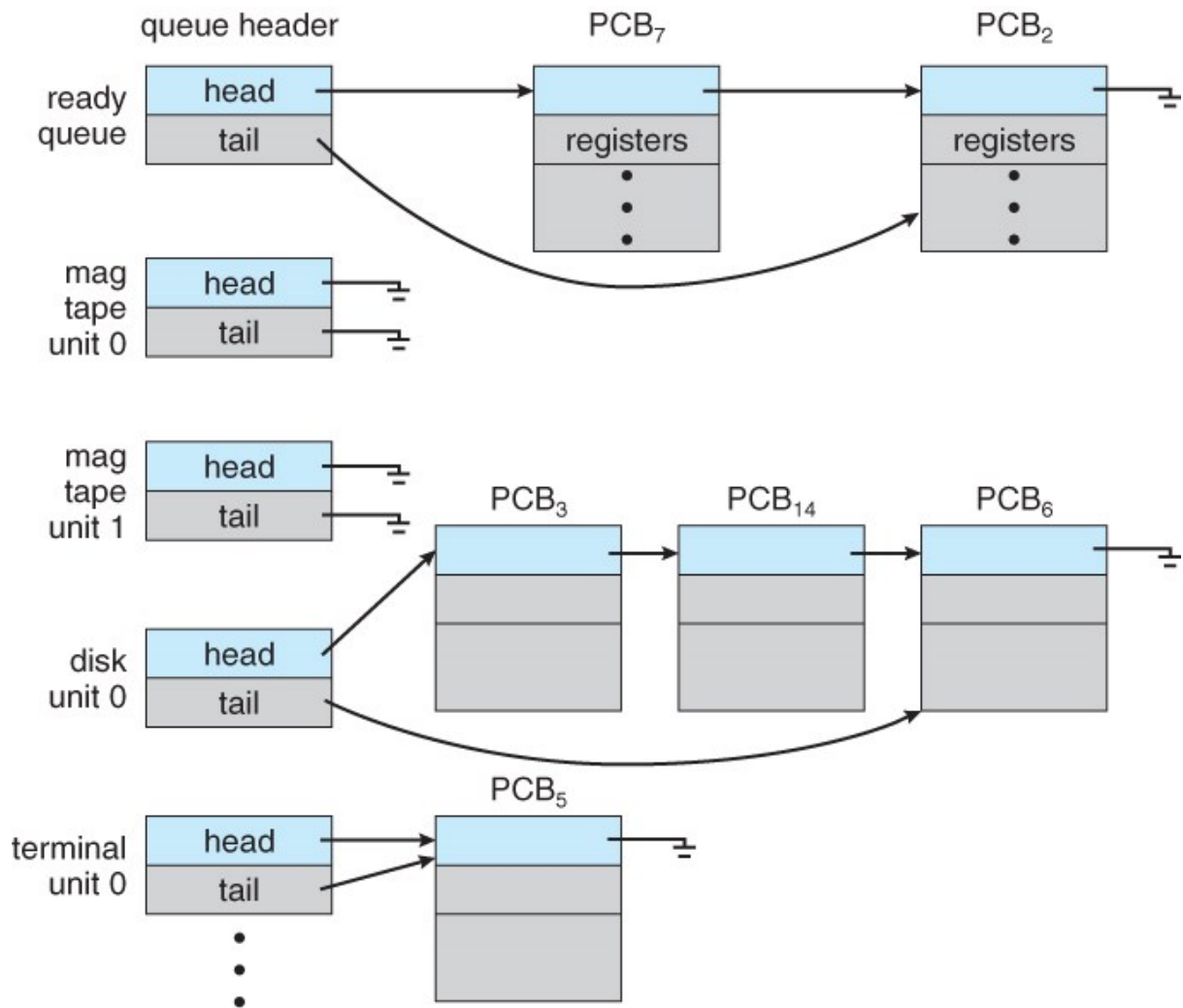

**Figure 3.3 - Process control block ( PCB )**

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

## 3.2 Process Scheduling

- The two main objectives of the process scheduling system are to keep the CPU busy at all times and to deliver "acceptable" response times for all programs, particularly for interactive ones.
- The process scheduler must meet these objectives by implementing suitable policies for swapping processes in and out of the CPU.
- ( Note that these objectives can be conflicting. In particular, every time the system steps in to swap processes it takes up time on the CPU to do so, which is thereby "lost" from doing any useful productive work. )
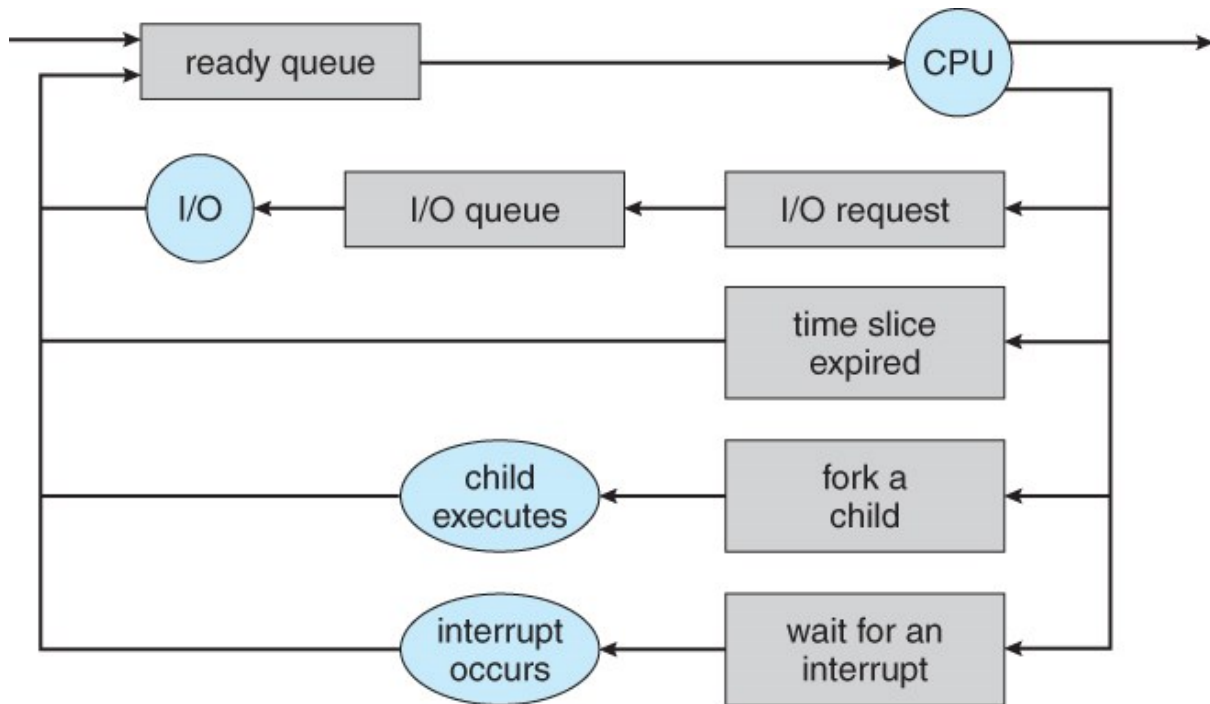
### 3.2.1 Scheduling Queues

- All processes are stored in the **job queue.**
- Processes in the Ready state are placed in the **ready queue.**
- Processes waiting for a device to become available or to deliver data are placed in **device queues**. There is generally a separate device queue for each device.
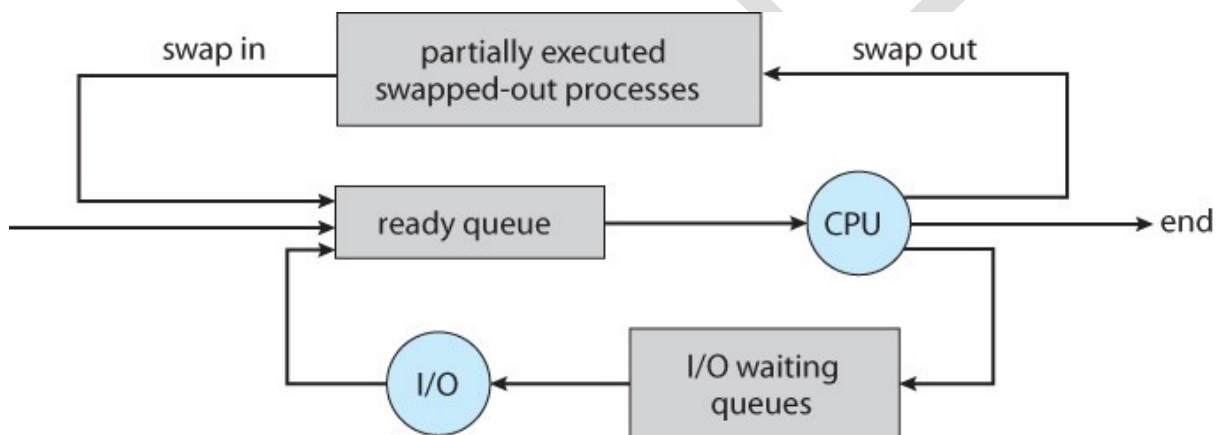- Other queues may also be created and used as needed.

**Figure 3.5 - The ready queue and various I/O device queues**

### 3.2.2 Schedulers

- A **long-term scheduler** is typical of a batch system or a very heavily loaded system. It runs infrequently, ( such as when one process ends selecting one more to be loaded in from disk in its place ), and can afford to take the time to implement intelligent and advanced scheduling algorithms.
- The **short-term scheduler**, or CPU Scheduler, runs very frequently, on the order of 100 milliseconds, and must very quickly swap one process out of the CPU and swap in another one.
- Some systems also employ a **medium-term scheduler**. When system loads get high, this scheduler will swap one or more processes out of the ready queue system for a few seconds, in order to allow smaller faster jobs to finish up quickly and clear the system. See the differences in Figures 3.7 and 3.8 below.
- An efficient scheduling system will select a good **process mix** of **CPU-bound** processes and **I/O bound** processes.

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



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

**3.2.3 Context Switch**

- Whenever an interrupt arrives, the CPU must do a **state-save** of the currently running process, then switch into kernel mode to handle the interrupt, and then do a **state-restore** of the interrupted process.
- Similarly, a **context switch** occurs when the time slice for one process has expired and a new process is to be loaded from the ready queue. This will be instigated by a timer interrupt, which will then cause the current process's state to be saved and the new process's state to be restored.
- Saving and restoring states involves saving and restoring all of the registers and program counter(s), as well as the process control blocks described above.
- Context switching happens VERY VERY frequently, and the overhead of doing the switching is just lost CPU time, so context switches ( state saves & restores ) need to be as fast as possible. Some hardware has special provisions for

speeding this up, such as a single machine instruction for saving or restoring all registers at once.

- Some Sun hardware actually has multiple sets of registers, so the context switching can be speeded up by merely switching which set of registers are currently in use. Obviously there is a limit as to how many processes can be switched between in this manner, making it attractive to implement the medium-term scheduler to swap some processes out as shown in Figure 3.8 above.

## MULTITASKING IN MOBILE SYSTEMS

Because of the constraints imposed on mobile devices, early versions of iOS did not provide user-application multitasking; only one application runs in the foreground and all other user applications are suspended. Operating-system tasks were multitasked because they were written by Apple and well behaved. However, beginning with iOS 4, Apple now provides a limited form of multitasking for user applications, thus allowing a single foreground application to run concurrently with multiple background applications. (On a mobile device, the foreground application is the application currently open and appearing on the display. The background application remains in memory, but does not occupy the display screen.) The iOS 4 programming API provides support for multitasking, thus allowing a process to run in the background without being suspended. However, it is limited and only available for a limited number of application types, including applications

- running a single, finite-length task (such as completing a download of content from a network);

- receiving notifications of an event occurring (such as a new email message);

- with long-running background tasks (such as an audio player.)

Apple probably limits multitasking due to battery life and memory use concerns. The CPU certainly has the features to support multitasking, but Apple chooses to not take advantage of some of them in order to better manage resource use.
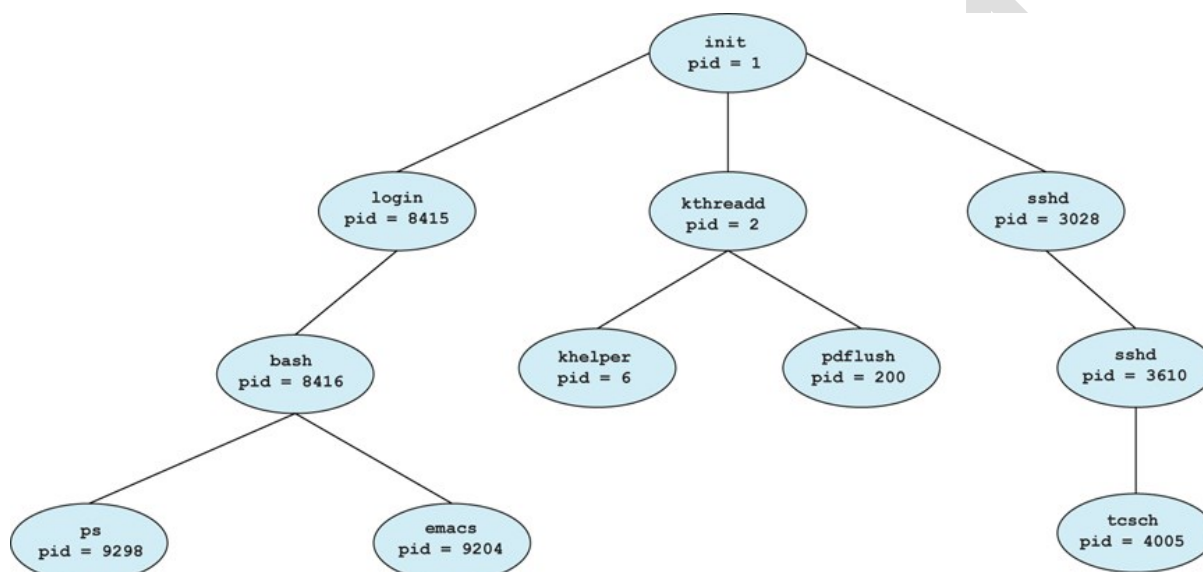
Android does not place such constraints on the types of applications that can run in the background. If an application requires processing while in the background, the application must use a service, a separate application component that runs on behalf of the background process. Consider a streaming audio application: if the application moves to the background, the service continues to send audio files to the audio device driver on behalf of the background application. In fact, the service will continue to run even if the background application is suspended. Services do not have a user interface and have a small memory footprint, thus providing an efficient technique for multitasking in a mobile environment.

## 3.3 Operations on Processes

### 3.3.1 Process Creation

- Processes may create other processes through appropriate system calls, such as **fork** or **spawn**. The process which does the creating is termed the **parent** of the other process, which is termed its **child**.

- Each process is given an integer identifier, termed its **process identifier**, or PID. The parent PID ( PPID ) is also stored for each process.
- On typical UNIX systems the process scheduler is termed **sched**, and is given PID 0. The first thing it does at system startup time is to launch **init**, which gives that process PID 1. Init then launches all system daemons and user logins, and becomes the ultimate parent of all other processes. Figure 3.9 shows a typical process tree for a Linux system, and other systems will have similar though not identical trees:



**Figure 3.8 - A tree of processes on a typical Linux system**

- Depending on system implementation, a child process may receive some amount of shared resources with its parent. Child processes may or may not be limited to a subset of the resources originally allocated to the parent, preventing runaway children from consuming all of a certain system resource.
- There are two options for the parent process after creating the child:
    1. Wait for the child process to terminate before proceeding. The parent makes a wait( ) system call, for either a specific child or for any child, which causes the parent process to block until the wait( ) returns. UNIX shells normally wait for their children to complete before issuing a new prompt.
    2. Run concurrently with the child, continuing to process without waiting. This is the operation seen when a UNIX shell runs a process as a background task. It is also possible for the parent to run for a while, and then wait for the child later, which might occur in a sort of a parallel processing operation. ( E.g. the parent may fork off a number of children without waiting for any of them, then do a little work of its own, and then wait for the children. )
- Two possibilities for the address space of the child relative to the parent:
    1. The child may be an exact duplicate of the parent, sharing the same program and data segments in memory. Each will have their own PCB,

including program counter, registers, and PID. This is the behavior of the **fork** system call in UNIX.

2. The child process may have a new program loaded into its address space, with all new code and data segments. This is the behavior of the **spawn** system calls in Windows. UNIX systems implement this as a second step, using the **exec** system call.

- Figures 3.10 and 3.11 below shows the fork and exec process on a UNIX system. Note that the **fork** system call returns the PID of the processes child to each process - It returns a zero to the child process and a non-zero child PID to the parent, so the return value indicates which process is which. Process IDs can be looked up any time for the current process or its direct parent using the getpid( ) and getppid( ) system calls respectively.

```c
#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");
        exit(-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");
        exit(0);
    }
}
```
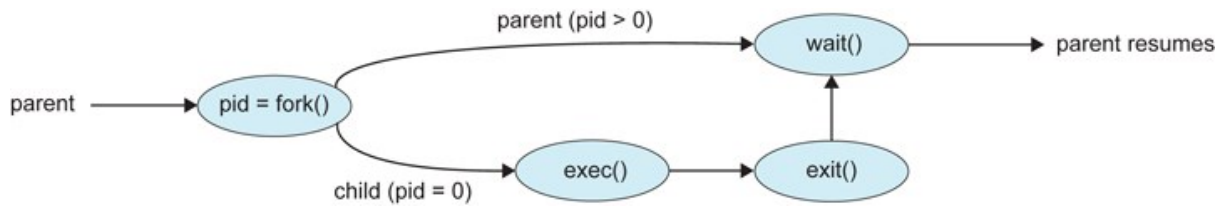
**Figure 3.10** C program forking a separate process.

**Figure 3.9 Creating a separate process using the UNIX fork( ) system call.**

**Figure 3.10 - Process creation using the fork( ) system call**

- Related man pages:
  - fork( 2 )
  - exec( 3 )
  - wait( 2 )

- Figure 3.12 shows the more complicated process for Windows, which must provide all of the parameter information for the new process as part of the forking process.

```
#include <stdio.h>
#include <windows.h>

int main(VOID)
{
STARTUPINFO si;
PROCESS_INFORMATION pi;

    // allocate memory
    ZeroMemory(&si, sizeof(si));
    si.cb = sizeof(si);
    ZeroMemory(&pi, sizeof(pi));

    // create child process
    if (!CreateProcess(NULL, // use command line
        "C:\\WINDOWS\\system32\\mspaint.exe", // command line
        NULL, // don't inherit process handle
        NULL, // don't inherit thread handle
        FALSE, // disable handle inheritance
        0, // no creation flags
        NULL, // use parent's environment block
        NULL, // use parent's existing directory
        &si,
        &pi))
    {
        fprintf(stderr, "Create Process Failed");
        return -1;
    }
    // parent will wait for the child to complete
    WaitForSingleObject(pi.hProcess, INFINITE);
    printf("Child Complete");

    // close handles
    CloseHandle(pi.hProcess);
    CloseHandle(pi.hThread);
}
```

Figure 3.12 Creating a separate process using the Win32 API.

**Figure 3.11**

### 3.3.2 Process Termination

- Processes may request their own termination by making the **exit( )** system call, typically returning an int. This int is passed along to the parent if it is doing

a **wait( )**, and is typically zero on successful completion and some non-zero code in the event of problems.

- o child code:

```
o               int exitCode;
                exit( exitCode );  // return exitCode; has the same
        effect when executed from main( )
```

- o parent code:

```
o               pid_t pid;
o               int status
o               pid = wait( &status );
o               // pid indicates which child exited. exitCode in low-
        order bits of status
                // macros can test the high-order bits of status for why
        it stopped
```

- Processes may also be terminated by the system for a variety of reasons, including:
  - o The inability of the system to deliver necessary system resources.
  - o In response to a KILL command, or other un handled process interrupt.
  - o A parent may kill its children if the task assigned to them is no longer needed.
  - o If the parent exits, the system may or may not allow the child to continue without a parent. ( On UNIX systems, orphaned processes are generally inherited by init, which then proceeds to kill them. The UNIX *nohup* command allows a child to continue executing after its parent has exited. )
- When a process terminates, all of its system resources are freed up, open files flushed and closed, etc. The process termination status and execution times are returned to the parent if the parent is waiting for the child to terminate, or eventually returned to init if the process becomes an orphan. ( Processes which are trying to terminate but which cannot because their parent is not waiting for them are termed **zombies**. These are eventually inherited by init as orphans and killed off. Note that modern UNIX shells do not produce as many orphans and zombies as older systems used to. )
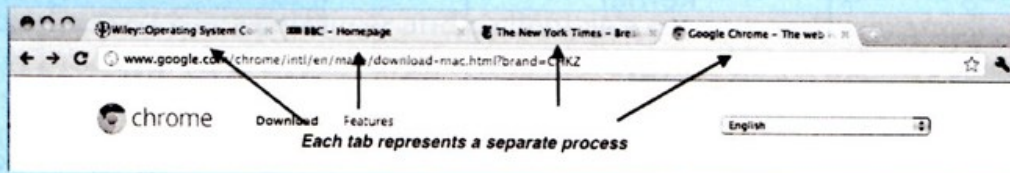
## 3.4 Interprocess Communication

- **Independent Processes** operating concurrently on a systems are those that can neither affect other processes or be affected by other processes.
- **Cooperating Processes** are those that can affect or be affected by other processes. There are several reasons why cooperating processes are allowed:
  - o Information Sharing - There may be several processes which need access to the same file for example. ( e.g. pipelines. )
  - o Computation speedup - Often a solution to a problem can be solved faster if the problem can be broken down into sub-tasks to be solved simultaneously ( particularly when multiple processors are involved. )

- Modularity - The most efficient architecture may be to break a system down into cooperating modules. ( E.g. databases with a client-server architecture. )
- Convenience - Even a single user may be multi-tasking, such as editing, compiling, printing, and running the same code in different windows.

## MULTIPROCESS ARCHITECTURE—CHROME BROWSER

Many websites contain active content such as JavaScript, Flash, and HTML5 to provide a rich and dynamic web-browsing experience. Unfortunately, these web applications may also contain software bugs, which can result in sluggish response times and can even cause the web browser to crash. This isn't a big problem in a web browser that displays content from only one website. But most contemporary web browsers provide tabbed browsing, which allows a single instance of a web browser application to open several websites at the same time, with each site in a separate tab. To switch between the different sites , a user need only click on the appropriate tab. This arrangement is illustrated below:
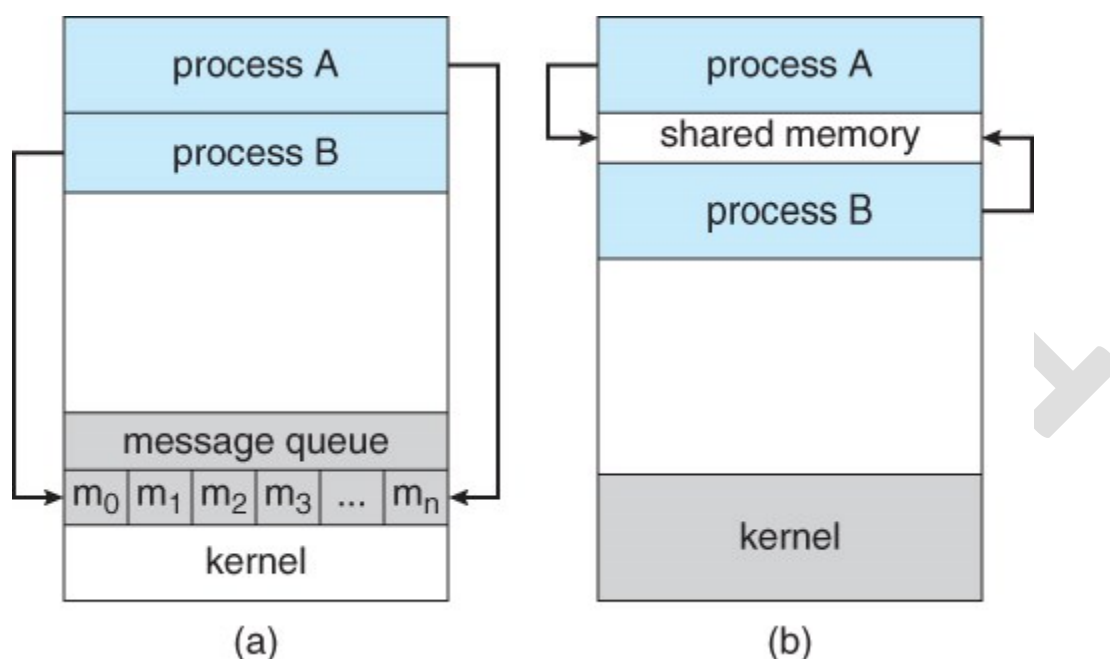


*Each tab represents a separate process*

A problem with this approach is that if a web application in any tab crashes, the entire process—including all other tabs displaying additional websites —crashes as well.

Google's Chrome web browser was designed to address this issue by using a multiprocess architecture. Chrome identifies three different types of processes: browser, renderers, and plug-ins.

- The **browser** process is responsible for managing the user interface as well as disk and network I/O. A new browser process is created when Chrome is started. Only one browser process is created.

- **Renderer** processes contain logic for rendering web pages. Thus, they contain the logic for handling HTML, Javascript, images, and so forth. As a general rule, a new renderer process is created for each website opened in a new tab, and so several renderer processes may be active at the same time.

- A **plug-in** process is created for each type of plug-in (such as Flash or QuickTime) in use. Plug-in processes contain the code for the plug-in as well as additional code that enables the plug-in to communicate with associated renderer processes and the browser process.

The advantage of the multiprocess approach is that websites run in isolation from one another. If one website crashes, only its renderer process is affected; all other processes remain unharmed. Furthermore, renderer processes run in a **sandbox**, which means that access to disk and network I/O is restricted, minimizing the effects of any security exploits.

- Cooperating processes require some type of inter-process communication, which is most commonly one of two types: Shared Memory systems or Message Passing systems. Figure 3.13 illustrates the difference between the two systems:



**Figure 3.12 - Communications models: (a) Message passing. (b) Shared memory.**

- Shared Memory is faster once it is set up, because no system calls are required and access occurs at normal memory speeds. However it is more complicated to set up, and doesn't work as well across multiple computers. Shared memory is generally preferable when large amounts of information must be shared quickly on the same computer.
- Message Passing requires system calls for every message transfer, and is therefore slower, but it is simpler to set up and works well across multiple computers. Message passing is generally preferable when the amount and/or frequency of data transfers is small, or when multiple computers are involved.

**3.4.1 Shared-Memory Systems**

- In general the memory to be shared in a shared-memory system is initially within the address space of a particular process, which needs to make system calls in order to make that memory publicly available to one or more other processes.
- Other processes which wish to use the shared memory must then make their own system calls to attach the shared memory area onto their address space.
- Generally a few messages must be passed back and forth between the cooperating processes first in order to set up and coordinate the shared memory access.

**Producer-Consumer Example Using Shared Memory**

- This is a classic example, in which one process is producing data and another process is consuming the data. ( In this example in the order in which it is produced, although that could vary. )
- The data is passed via an intermediary buffer, which may be either unbounded or bounded. With a bounded buffer the producer may have to wait until there is space available in the buffer, but with an unbounded buffer the producer will never need to wait. The consumer may need to wait in either case until there is data available.
- This example uses shared memory and a circular queue. Note in the code below that only the producer changes "in", and only the consumer changes "out", and that they can never be accessing the same array location at the same time.
- First the following data is set up in the shared memory area:

```
#define BUFFER_SIZE 10

typedef struct {
     . . .
} item;

item buffer[ BUFFER_SIZE ];
int in = 0;
int out = 0;
```

- Then the producer process. Note that the buffer is full when "in" is one less than "out" in a circular sense:

```
// Code from Figure 3.13

item nextProduced;

while( true ) {

/* Produce an item and store it in nextProduced */
nextProduced = makeNewItem( . . . );

/* Wait for space to become available */
while( ( ( in + 1 ) % BUFFER_SIZE ) == out )
    ; /* Do nothing */

/* And then store the item and repeat the loop. */
buffer[ in ] = nextProduced;
in = ( in + 1 ) % BUFFER_SIZE;

}
```

- Then the consumer process. Note that the buffer is empty when "in" is equal to "out":

```
// Code from Figure 3.14

item nextConsumed;

while( true ) {

/* Wait for an item to become available */
while( in == out )
      ; /* Do nothing */

/* Get the next available item */
nextConsumed = buffer[ out ];
out = ( out + 1 ) % BUFFER_SIZE;

/* Consume the item in nextConsumed
      ( Do something with it ) */

}
```

### 3.4.2 Message-Passing Systems

- Message passing systems must support at a minimum system calls for "send message" and "receive message".
- A communication link must be established between the cooperating processes before messages can be sent.
- There are three key issues to be resolved in message passing systems as further explored in the next three subsections:
    - Direct or indirect communication ( naming )
    - Synchronous or asynchronous communication
    - Automatic or explicit buffering.

#### 3.4.2.1 Naming

- With **direct communication** the sender must know the name of the receiver to which it wishes to send a message.
    - There is a one-to-one link between every sender-receiver pair.
    - For **symmetric** communication, the receiver must also know the specific name of the sender from which it wishes to receive messages.
      For **asymmetric** communications, this is not necessary.
- **Indirect communication** uses shared mailboxes, or ports.
    - Multiple processes can share the same mailbox or boxes.
    - Only one process can read any given message in a mailbox. Initially the process that creates the mailbox is the owner, and is the only one allowed to read mail in the mailbox, although this privilege may be transferred.

55

- ( Of course the process that reads the message can immediately turn around and place an identical message back in the box for someone else to read, but that may put it at the back end of a queue of messages. )
  - The OS must provide system calls to create and delete mailboxes, and to send and receive messages to/from mailboxes.

### 3.4.2.2 Synchronization

- Either the sending or receiving of messages ( or neither or both ) may be either **blocking** or **non-blocking**.

### 3.4.2.3 Buffering

- Messages are passed via queues, which may have one of three capacity configurations:
  1. **Zero capacity** - Messages cannot be stored in the queue, so senders must block until receivers accept the messages.
  2. **Bounded capacity**- There is a certain pre-determined finite capacity in the queue. Senders must block if the queue is full, until space becomes available in the queue, but may be either blocking or non-blocking otherwise.
  3. **Unbounded capacity** - The queue has a theoretical infinite capacity, so senders are never forced to block.

```
message next_produced;

while (true) {
    /* produce an item in next_produced */

    send(next_produced);
}
```

**Figure 3.15** The producer process using message passing.

```
message next_consumed;

while (true) {
    receive(next_consumed);

    /* consume the item in next_consumed */
}
```

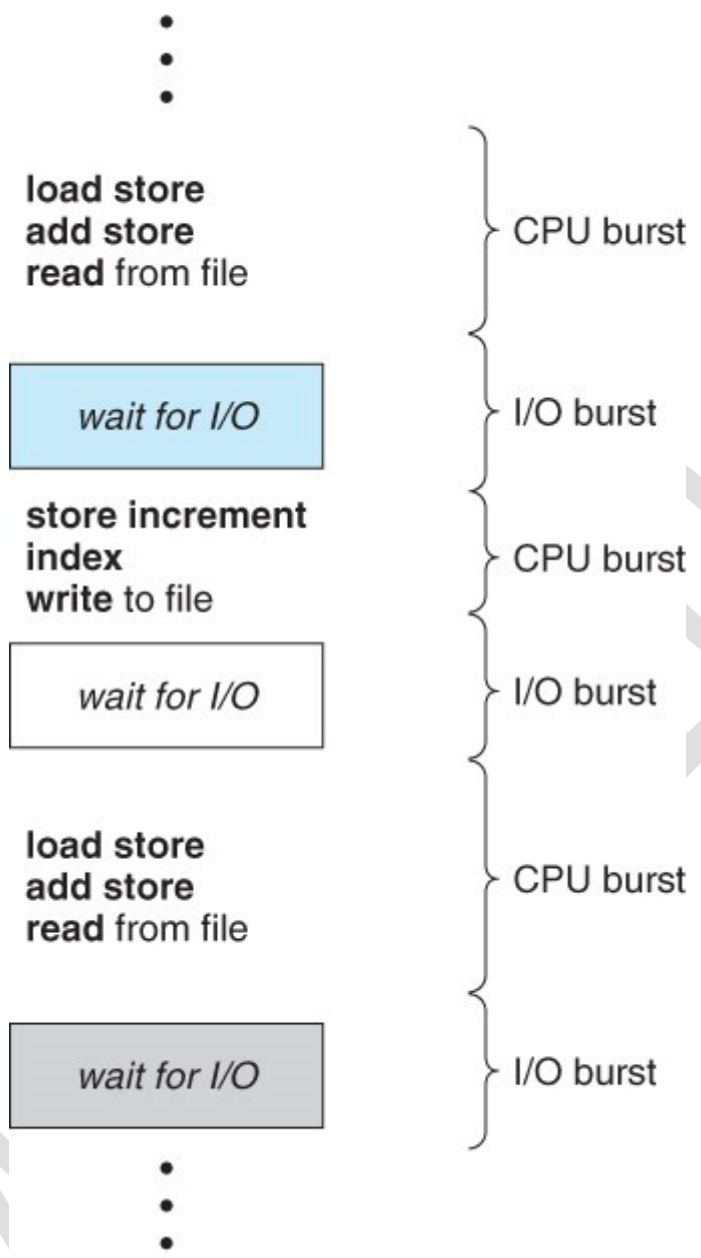**Figure 3.16** The consumer process using message passing.

## CPU Scheduling

## 6.1 Basic Concepts

- Almost all programs have some alternating cycle of CPU number crunching and waiting for I/O of some kind. ( Even a simple fetch from memory takes a long time relative to CPU speeds. )
- In a simple system running a single process, the time spent waiting for I/O is wasted, and those CPU cycles are lost forever.
- A scheduling system allows one process to use the CPU while another is waiting for I/O, thereby making full use of otherwise lost CPU cycles.
- The challenge is to make the overall system as "efficient" and "fair" as possible, subject to varying and often dynamic conditions, and where "efficient" and "fair" are somewhat subjective terms, often subject to shifting priority policies.
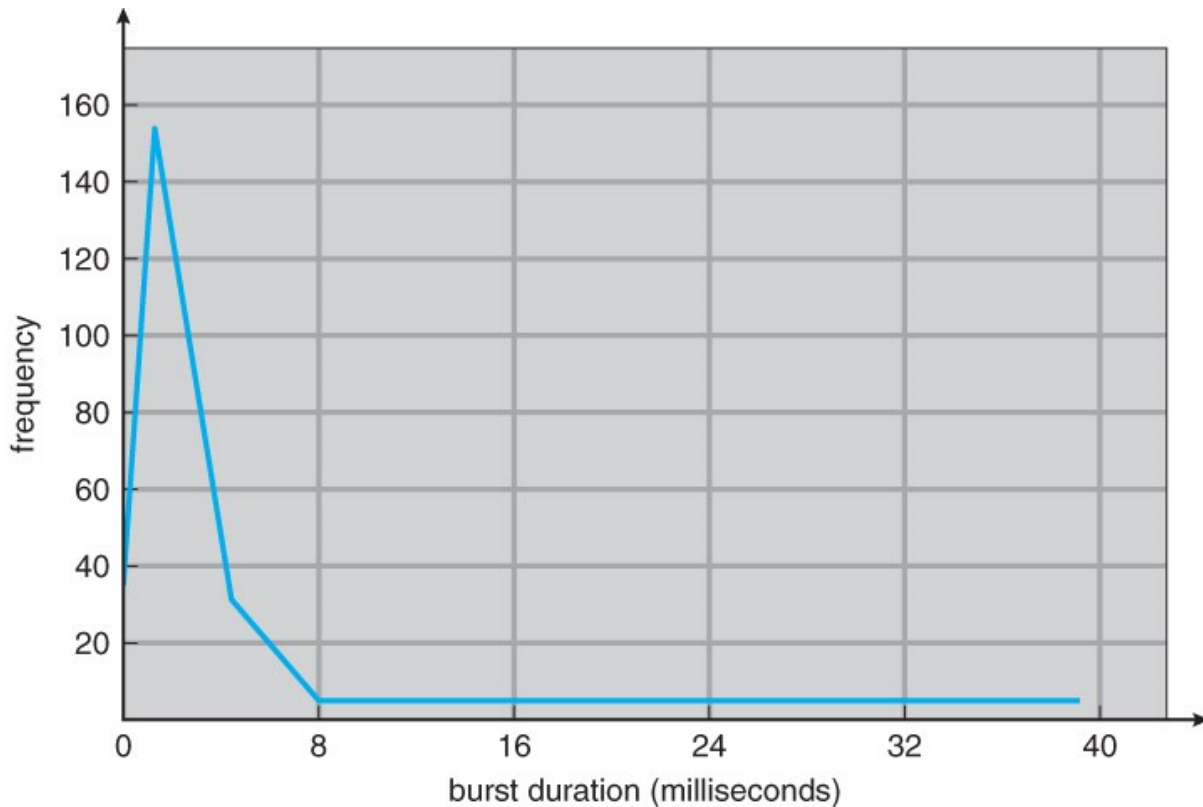
### 6.1.1 CPU-I/O Burst Cycle

- Almost all processes alternate between two states in a continuing *cycle*, as shown in Figure 6.1 below *:*
  - A CPU burst of performing calculations, and
  - An I/O burst, waiting for data transfer in or out of the system.

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

- CPU bursts vary from process to process, and from program to program, but an extensive study shows frequency patterns similar to that shown in Figure 6.2:

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

### 6.1.2 CPU Scheduler

- Whenever the CPU becomes idle, it is the job of the CPU Scheduler ( a.k.a. the short-term scheduler ) to select another process from the ready queue to run next.
- The storage structure for the ready queue and the algorithm used to select the next process are not necessarily a FIFO queue. There are several alternatives to choose from, as well as numerous adjustable parameters for each algorithm, which is the basic subject of this entire chapter.

### 6.1.3. Preemptive Scheduling

- CPU scheduling decisions take place under one of four conditions:
    1. When a process switches from the running state to the waiting state, such as for an I/O request or invocation of the wait( ) system call.
    2. When a process switches from the running state to the ready state, for example in response to an interrupt.
    3. When a process switches from the waiting state to the ready state, say at completion of I/O or a return from wait( ).
    4. When a process terminates.
- For conditions 1 and 4 there is no choice - A new process must be selected.
- For conditions 2 and 3 there is a choice - To either continue running the current process, or select a different one.
- If scheduling takes place only under conditions 1 and 4, the system is said to be *non-preemptive*, or *cooperative*. Under these conditions, once a

process starts running it keeps running, until it either voluntarily blocks or until it finishes. Otherwise the system is said to be ***preemptive.***

- Windows used non-preemptive scheduling up to Windows 3.x, and started using pre-emptive scheduling with Win95. Macs used non-preemptive prior to OSX, and pre-emptive since then. Note that pre-emptive scheduling is only possible on hardware that supports a timer interrupt.
- Note that pre-emptive scheduling can cause problems when two processes share data, because one process may get interrupted in the middle of updating shared data structures. Chapter 5 examined this issue in greater detail.
- Preemption can also be a problem if the kernel is busy implementing a system call ( e.g. updating critical kernel data structures ) when the preemption occurs. Most modern UNIXes deal with this problem by making the process wait until the system call has either completed or blocked before allowing the preemption Unfortunately this solution is problematic for real-time systems, as real-time response can no longer be guaranteed.
- Some critical sections of code protect themselves from con currency problems by disabling interrupts before entering the critical section and re-enabling interrupts on exiting the section. Needless to say, this should only be done in rare situations, and only on very short pieces of code that will finish quickly, ( usually just a few machine instructions. )

### 6.1.4 Dispatcher

- The **dispatcher** is the module that gives control of the CPU to the process selected by the scheduler. This function involves:
  - Switching context.
  - Switching to user mode.
  - Jumping to the proper location in the newly loaded program.
- The dispatcher needs to be as fast as possible, as it is run on every context switch. The time consumed by the dispatcher is known as **dispatch latency.**

## 6.2 Scheduling Criteria

- There are several different criteria to consider when trying to select the "best" scheduling algorithm for a particular situation and environment, including:
  - **CPU utilization** - Ideally the CPU would be busy 100% of the time, so as to waste 0 CPU cycles. On a real system CPU usage should range from 40% ( lightly loaded ) to 90% ( heavily loaded. )

- o **Throughput** - Number of processes completed per unit time. May range from 10 / second to 1 / hour depending on the specific processes.
- o **Turnaround time** - Time required for a particular process to complete, from submission time to completion. ( Wall clock time. )
- o **Waiting time** - How much time processes spend in the ready queue waiting their turn to get on the CPU.
  - ▪ ( **Load average** - The average number of processes sitting in the ready queue waiting their turn to get into the CPU. Reported in 1-minute, 5-minute, and 15-minute averages by "uptime" and "who". )
- o **Response time** - The time taken in an interactive program from the issuance of a command to the *commence* of a response to that command.
- In general one wants to optimize the average value of a criteria ( Maximize CPU utilization and throughput, and minimize all the others. ) However some times one wants to do something different, such as to minimize the maximum response time.
- Sometimes it is most desirable to minimize the *variance* of a criteria than the actual value. I.e. users are more accepting of a consistent predictable system than an inconsistent one, even if it is a little bit slower.
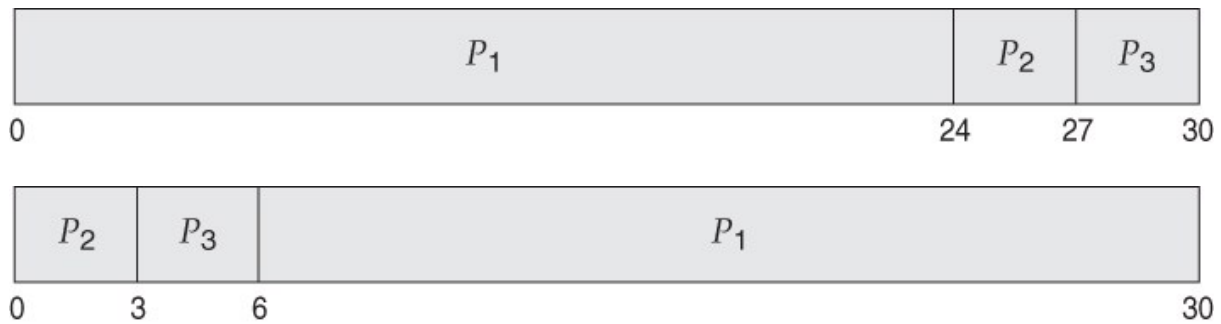
## 6.3 Scheduling Algorithms

The following subsections will explain several common scheduling strategies, looking at only a single CPU burst each for a small number of processes. Obviously real systems have to deal with a lot more simultaneous processes executing their CPU-I/O burst cycles.

### 6.3.1 First-Come First-Serve Scheduling, FCFS

- FCFS is very simple - Just a FIFO queue, like customers waiting in line at the bank or the post office or at a copying machine.
- Unfortunately, however, FCFS can yield some very long average wait times, particularly if the first process to get there takes a long time. For example, consider the following three processes:

| Process | Burst Time |
|---------|------------|
| P1 | 24 |
| P2 | 3 |
| P3 | 3 |

- In the first Gantt chart below, process P1 arrives first. The average waiting time for the three processes is ( 0 + 24 + 27 ) / 3 = 17.0 ms.
- In the second Gantt chart below, the same three processes have an average wait time of ( 0 + 3 + 6 ) / 3 = 3.0 ms. The total run time for the three bursts is the same, but in the second case two of the three finish much quicker, and the other process is only delayed by a short amount.
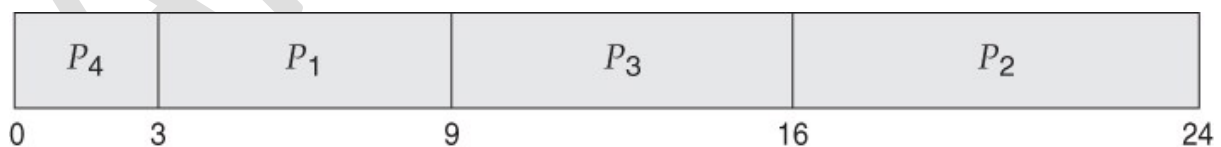
| | | | | | $P_1$ | | | | | | $P_2$ | $P_3$ |

0                                                      24     27     30

| $P_2$ | $P_3$ | $P_1$ |

0     3     6                                                    30

- FCFS can also block the system in a busy dynamic system in another way, known as the **convoy effect**. When one CPU intensive process blocks the CPU, a number of I/O intensive processes can get backed up behind it, leaving the I/O devices idle. When the CPU hog finally relinquishes the CPU, then the I/O processes pass through the CPU quickly, leaving the CPU idle while everyone queues up for I/O, and then the cycle repeats itself when the CPU intensive process gets back to the ready queue.

**6.3.2 Shortest-Job-First Scheduling, SJF**

- The idea behind the SJF algorithm is to pick the quickest fastest little job that needs to be done, get it out of the way first, and then pick the next smallest fastest job to do next.
- ( Technically this algorithm picks a process based on the next shortest **CPU burst**, not the overall process time. )
- For example, the Gantt chart below is based upon the following CPU burst times, ( and the assumption that all jobs arrive at the same time. )

| Process | Burst Time |
|---------|------------|
| P1 | 6 |
| P2 | 8 |
| P3 | 7 |
| P4 | 3 |

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

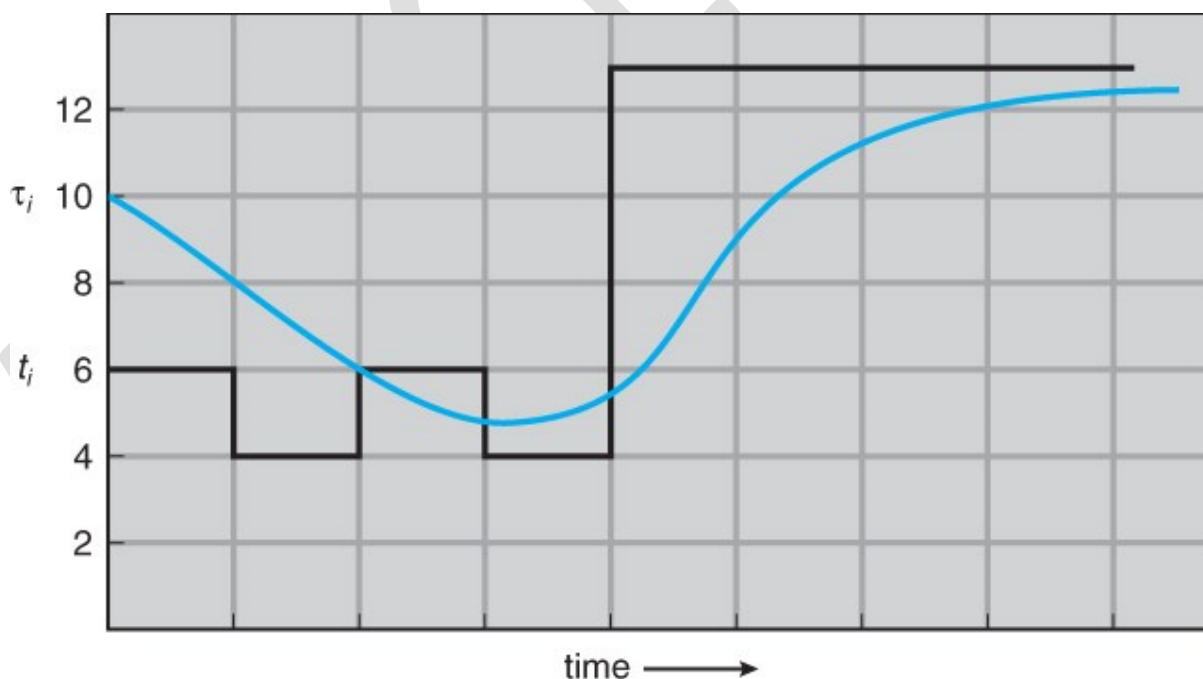0     3                   9                   16                 24

- In the case above the average wait time is ( 0 + 3 + 9 + 16 ) / 4 = 7.0 ms, ( as opposed to 10.25 ms for FCFS for the same processes. )
- SJF can be proven to be the fastest scheduling algorithm, but it suffers from one important problem: How do you know how long the next CPU burst is going to be?
  - o For long-term batch jobs this can be done based upon the limits that users set for their jobs when they submit them, which encourages

them to set low limits, but risks their having to re-submit the job if they set the limit too low. However that does not work for short-term CPU scheduling on an interactive system.

o Another option would be to statistically measure the run time characteristics of jobs, particularly if the same tasks are run repeatedly and predictably. But once again that really isn't a viable option for short term CPU scheduling in the real world.

o A more practical approach is to *predict* the length of the next burst, based on some historical measurement of recent burst times for this process. One simple, fast, and relatively accurate method is the *exponential average*, which can be defined as follows. ( The book uses tau and t for their variables, but those are hard to distinguish from one another and don't work well in HTML. )

estimate[ i + 1 ] = alpha * burst[ i ] + ( 1.0 - alpha ) * estimate[ i ]

o In this scheme the previous estimate contains the history of all previous times, and alpha serves as a weighting factor for the relative importance of recent data versus past history. If alpha is 1.0, then past history is ignored, and we assume the next burst will be the same length as the last burst. If alpha is 0.0, then all measured burst times are ignored, and we just assume a constant burst time. Most commonly alpha is set at 0.5, as illustrated in Figure 5.3:
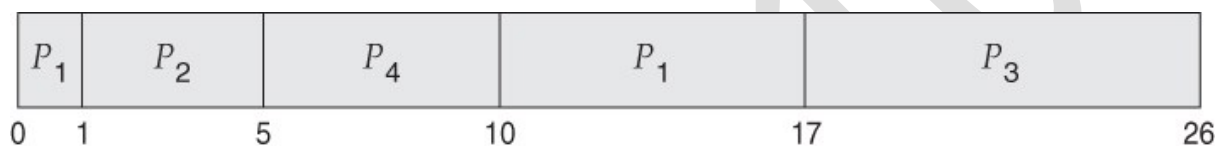


| CPU burst ($t_i$) | | 6 | 4 | 6 | 4 | 13 | 13 | 13 | ... |
|---|---|---|---|---|---|---|---|---|---|
| "guess" ($\tau_i$) | 10 | 8 | 6 | 6 | 5 | 9 | 11 | 12 | ... |

**Figure 6.3 - Prediction of the length of the next CPU burst.**

- SJF can be either preemptive or non-preemptive. Preemption occurs when a new process arrives in the ready queue that has a predicted burst time shorter than the time remaining in the process whose burst is currently on the CPU. Preemptive SJF is sometimes referred to as *shortest remaining time first scheduling.*
- For example, the following Gantt chart is based upon the following data:

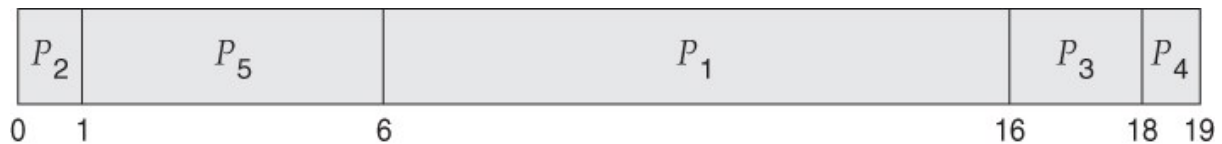| Process | Arrival Time | Burst Time |
|---------|--------------|------------|
| P1 | 0 | 8 |
| P2 | 1 | 4 |
| P3 | 2 | 9 |
| p4 | 3 | 5 |



- The average wait time in this case is $(( 5 - 3 ) + ( 10 - 1 ) + ( 17 - 2 )) / 4 = 26 / 4 = 6.5$ ms. ( As opposed to 7.75 ms for non-preemptive SJF or 8.75 for FCFS. )

### 6.3.3 Priority Scheduling

- Priority scheduling is a more general case of SJF, in which each job is assigned a priority and the job with the highest priority gets scheduled first. ( SJF uses the inverse of the next expected burst time as its priority - The smaller the expected burst, the higher the priority. )
- Note that in practice, priorities are implemented using integers within a fixed range, but there is no agreed-upon convention as to whether "high" priorities use large numbers or small numbers. This book uses low number for high priorities, with 0 being the highest possible priority.
- For example, the following Gantt chart is based upon these process burst times and priorities, and yields an average waiting time of 8.2 ms:

| Process | Burst Time | Priority |
|---------|------------|----------|
| P1 | 10 | 3 |
| P2 | 1 | 1 |
| P3 | 2 | 4 |
| P4 | 1 | 5 |
| P5 | 5 | 2 |

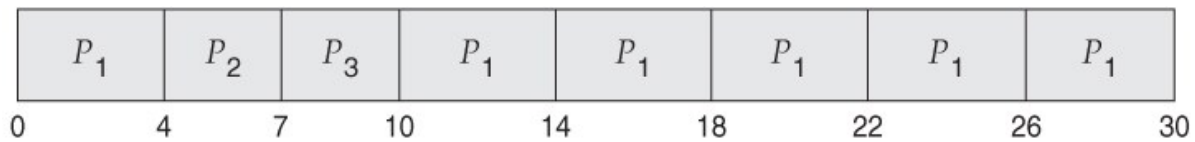| $P_2$ | $P_5$ | $P_1$ | $P_3$ | $P_4$ |
|---|---|---|---|---|
| 0   1 | 6 | 16 | 18 | 19 |

- Priorities can be assigned either internally or externally. Internal priorities are assigned by the OS using criteria such as average burst time, ratio of CPU to I/O activity, system resource use, and other factors available to the kernel. External priorities are assigned by users, based on the importance of the job, fees paid, politics, etc.
- Priority scheduling can be either preemptive or non-preemptive.
- Priority scheduling can suffer from a major problem known as *indefinite blocking*, or *starvation*, in which a low-priority task can wait forever because there are always some other jobs around that have higher priority.
    - o If this problem is allowed to occur, then processes will either run eventually when the system load lightens ( at say 2:00 a.m. ), or will eventually get lost when the system is shut down or crashes. ( There are rumors of jobs that have been stuck for years. )
    - o One common solution to this problem is *aging*, in which priorities of jobs increase the longer they wait. Under this scheme a low-priority job will eventually get its priority raised high enough that it gets run.
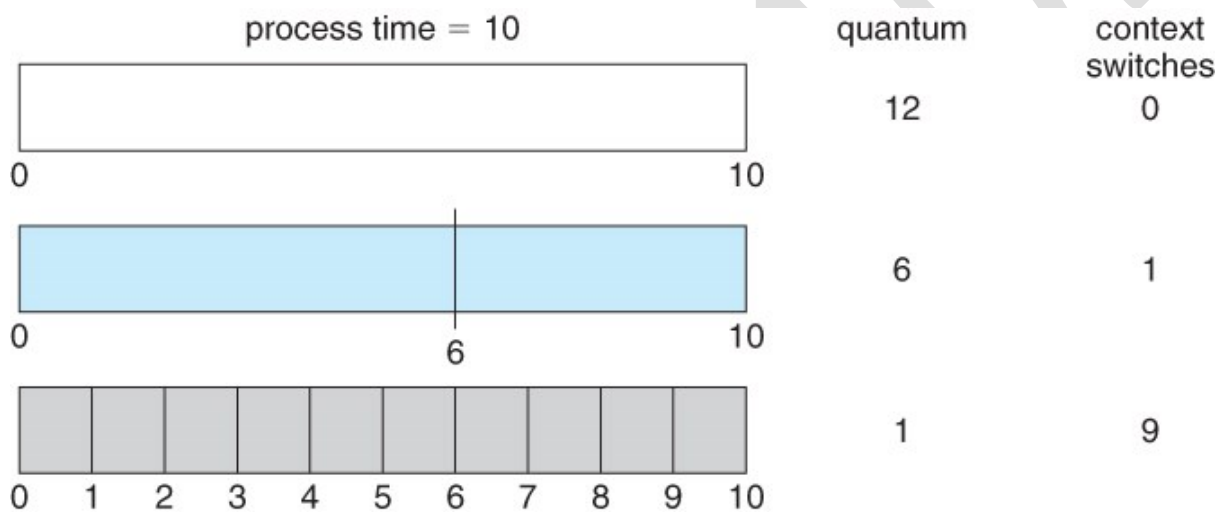
**6.3.4 Round Robin Scheduling**

- Round robin scheduling is similar to FCFS scheduling, except that CPU bursts are assigned with limits called *time quantum*.
- When a process is given the CPU, a timer is set for whatever value has been set for a time quantum.
    - o If the process finishes its burst before the time quantum timer expires, then it is swapped out of the CPU just like the normal FCFS algorithm.
    - o If the timer goes off first, then the process is swapped out of the CPU and moved to the back end of the ready queue.
- The ready queue is maintained as a circular queue, so when all processes have had a turn, then the scheduler gives the first process another turn, and so on.
- RR scheduling can give the effect of all processors sharing the CPU equally, although the average wait time can be longer than with other scheduling algorithms. In the following example the average wait time is 5.66 ms.

| Process | Burst Time |
|---|---|
| P1 | 24 |
| P2 | 3 |
| P3 | 3 |

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

```
0     4     7    10    14    18    22    26    30
```
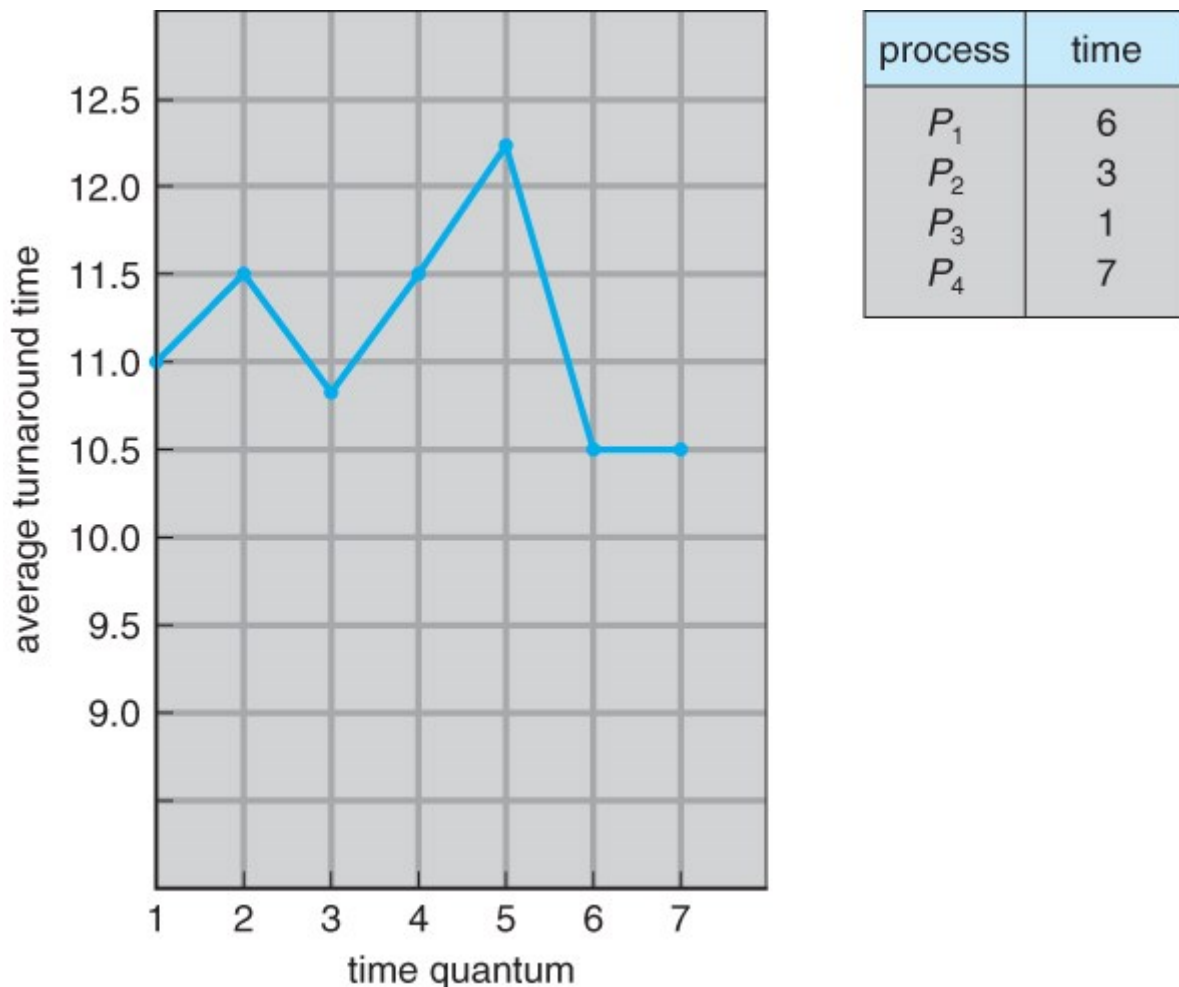
- The performance of RR is sensitive to the time quantum selected. If the quantum is large enough, then RR reduces to the FCFS algorithm; If it is very small, then each process gets 1/nth of the processor time and share the CPU equally.
- **BUT**, a real system invokes overhead for every context switch, and the smaller the time quantum the more context switches there are. ( See Figure 6.4 below. ) Most modern systems use time quantum between 10 and 100 milliseconds, and context switch times on the order of 10 microseconds, so the overhead is small relative to the time quantum.



| process time = 10 | quantum | context switches |
|---|---|---|
| 0 ... 10 | 12 | 0 |
| 0 ... 6 ... 10 | 6 | 1 |
| 0 1 2 3 4 5 6 7 8 9 10 | 1 | 9 |

**Figure 6.4 - The way in which a smaller time quantum increases context switches.**

- Turn around time also varies with quantum time, in a non-apparent manner. Consider, for example the processes shown in Figure 6.5:

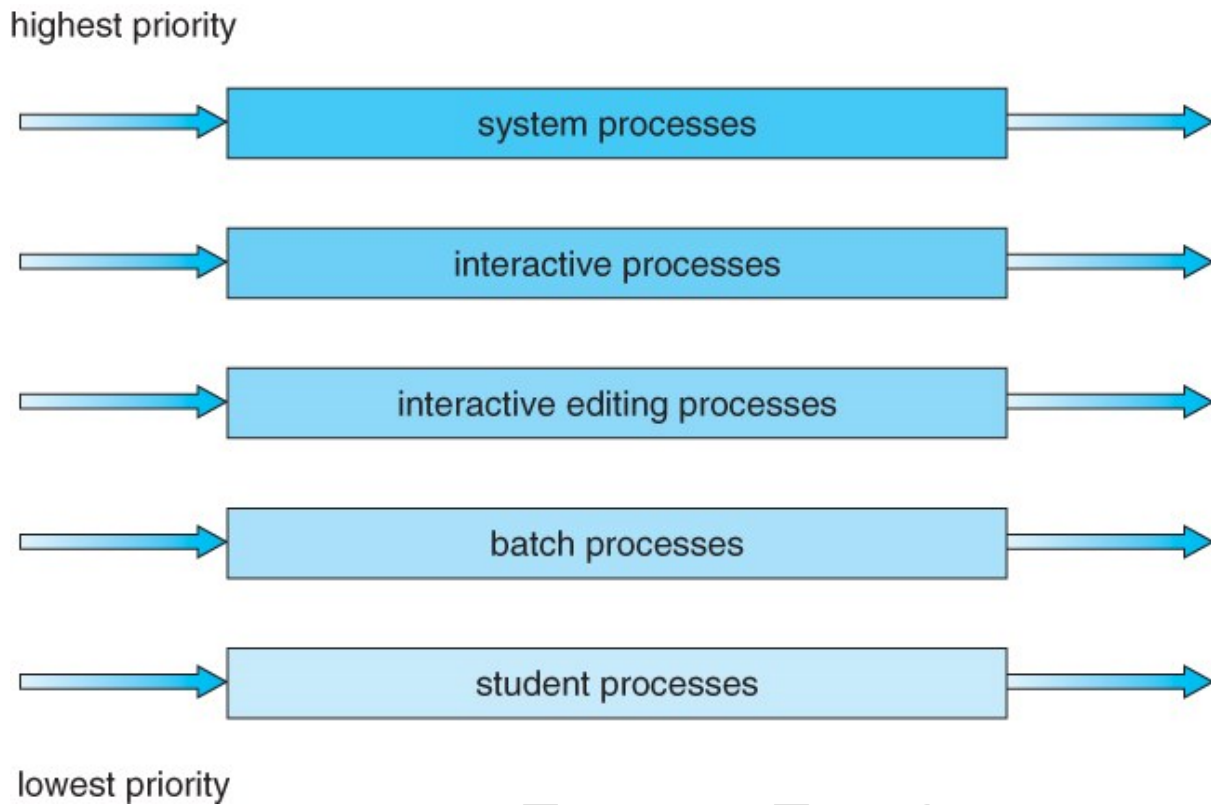| process | time |
|:---:|:---:|
| $P_1$ | 6 |
| $P_2$ | 3 |
| $P_3$ | 1 |
| $P_4$ | 7 |

**Figure 6.5 - The way in which turnaround time varies with the time quantum.**

- In general, turnaround time is minimized if most processes finish their next cpu burst within one time quantum. For example, with three processes of 10 ms bursts each, the average turnaround time for 1 ms quantum is 29, and for 10 ms quantum it reduces to 20. However, if it is made too large, then RR just degenerates to FCFS. A rule of thumb is that 80% of CPU bursts should be smaller than the time quantum.

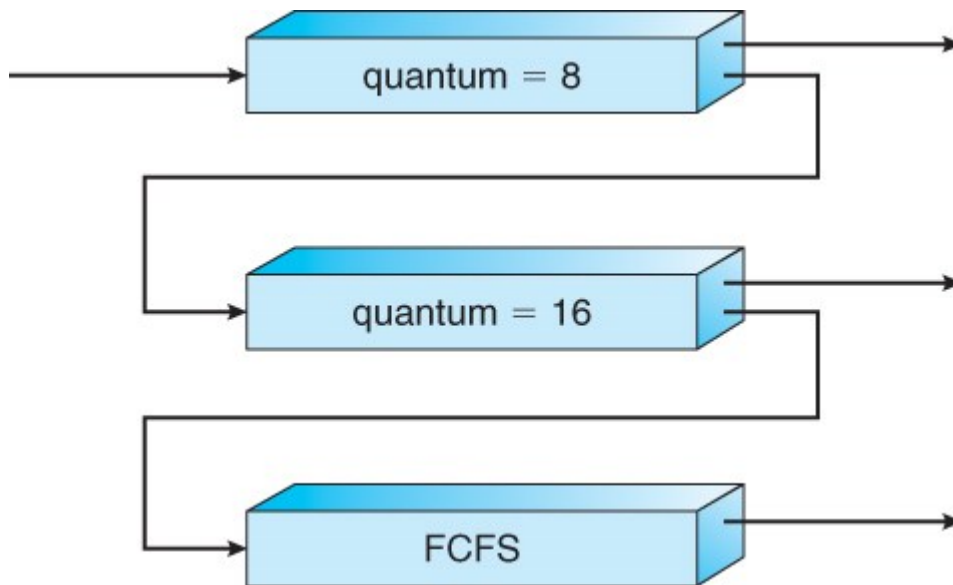### 6.3.5 Multilevel Queue Scheduling

- When processes can be readily categorized, then multiple separate queues can be established, each implementing whatever scheduling algorithm is most appropriate for that type of job, and/or with different parametric adjustments.
- Scheduling must also be done between queues, that is scheduling one queue to get time relative to other queues. Two common options are strict priority ( no job in a lower priority queue runs until all higher priority queues are empty ) and round-robin ( each queue gets a time slice in turn, possibly of different sizes. )
- Note that under this algorithm jobs cannot switch from queue to queue - Once they are assigned a queue, that is their queue until they finish.

**Figure 6.6 - Multilevel queue scheduling**

### 6.3.6 Multilevel Feedback-Queue Scheduling

- Multilevel feedback queue scheduling is similar to the ordinary multilevel queue scheduling described above, except jobs may be moved from one queue to another for a variety of reasons:
    - If the characteristics of a job change between CPU-intensive and I/O intensive, then it may be appropriate to switch a job from one queue to another.
    - Aging can also be incorporated, so that a job that has waited for a long time can get bumped up into a higher priority queue for a while.
- Multilevel feedback queue scheduling is the most flexible, because it can be tuned for any situation. But it is also the most complex to implement because of all the adjustable parameters. Some of the parameters which define one of these systems include:
    - The number of queues.
    - The scheduling algorithm for each queue.
    - The methods used to upgrade or demote processes from one queue to another. ( Which may be different. )
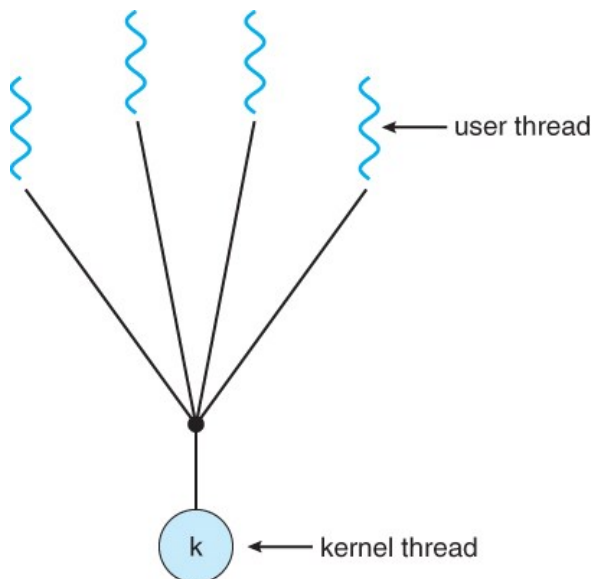    - The method used to determine which queue a process enters initially.

**Figure 6.7 - Multilevel feedback queues.**

## 4.3 Multithreading Models

- There are two types of threads to be managed in a modern system: User threads and kernel threads.
- User threads are supported above the kernel, without kernel support. These are the threads that application programmers would put into their programs.
- Kernel threads are supported within the kernel of the OS itself. All modern OSes support kernel level threads, allowing the kernel to perform multiple simultaneous tasks and/or to service multiple kernel system calls simultaneously.
- In a specific implementation, the user threads must be mapped to kernel threads, using one of the following strategies.
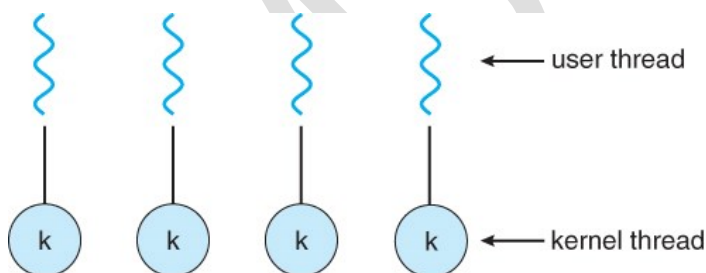
### 4.3.1 Many-To-One Model

- In the many-to-one model, many user-level threads are all mapped onto a single kernel thread.
- Thread management is handled by the thread library in user space, which is very efficient.
- However, if a blocking system call is made, then the entire process blocks, even if the other user threads would otherwise be able to continue.
- Because a single kernel thread can operate only on a single CPU, the many-to-one model does not allow individual processes to be split across multiple CPUs.
- Green threads for Solaris and GNU Portable Threads implement the many-to-one model in the past, but few systems continue to do so today.

**Figure 4.5 - Many-to-one model**

### 4.3.2 One-To-One Model

- The one-to-one model creates a separate kernel thread to handle each user thread.
- One-to-one model overcomes the problems listed above involving blocking system calls and the splitting of processes across multiple CPUs.
- However the overhead of managing the one-to-one model is more significant, involving more overhead and slowing down the system.
- Most implementations of this model place a limit on how many threads can be created.
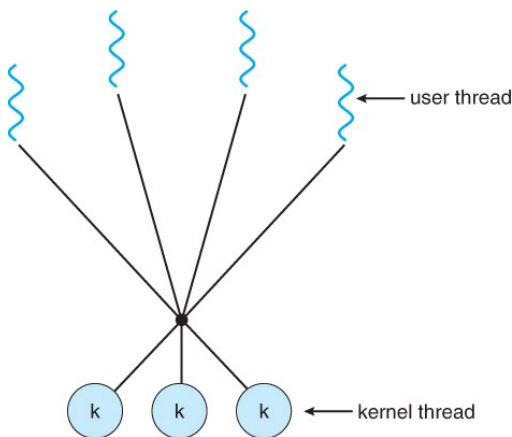- Linux and Windows from 95 to XP implement the one-to-one model for threads.



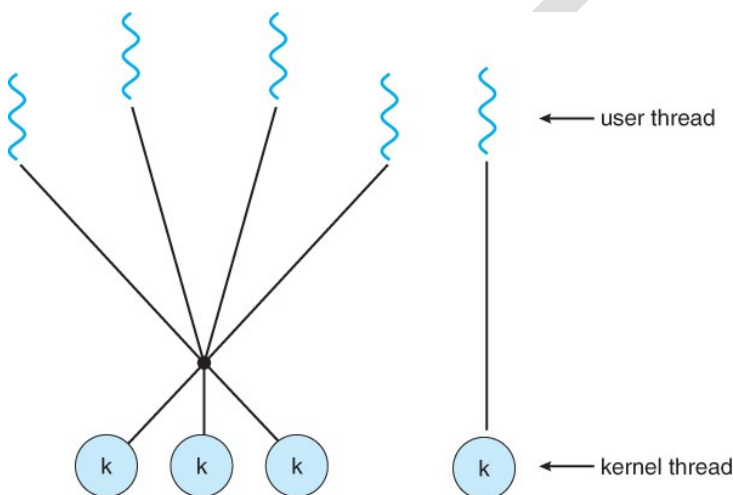**Figure 4.6 - One-to-one model**

### 4.3.3 Many-To-Many Model

- The many-to-many model multiplexes any number of user threads onto an equal or smaller number of kernel threads, combining the best features of the one-to-one and many-to-one models.
- Users have no restrictions on the number of threads created.
- Blocking kernel system calls do not block the entire process.
- Processes can be split across multiple processors.

- Individual processes may be allocated variable numbers of kernel threads, depending on the number of CPUs present and other factors.



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

- One popular variation of the many-to-many model is the two-tier model, which allows either many-to-many or one-to-one operation.
- IRIX, HP-UX, and Tru64 UNIX use the two-tier model, as did Solaris prior to Solaris 9.



**Figure 4.8 - Two-level model**

## 4.6 Threading Issues

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

- Q: If one thread forks, is the entire process copied, or is the new process single-threaded?
- A: System dependant.
- A: If the new process execs right away, there is no need to copy all the other threads. If it doesn't, then the entire process should be copied.

- A: Many versions of UNIX provide multiple versions of the fork call for this purpose.

### 4.6.2 Signal Handling

- Q: When a multi-threaded process receives a signal, to what thread should that signal be delivered?
- A: There are four major 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 to certain threads in the process.
  4. Assign a specific thread to receive all signals in a process.
- The best choice may depend on which specific signal is involved.
- UNIX allows individual threads to indicate which signals they are accepting and which they are ignoring. However the signal can only be delivered to one thread, which is generally the first thread that is accepting that particular signal.
- UNIX provides two separate system calls, **kill( pid, signal )** and **pthread_kill( tid, signal )**, for delivering signals to processes or specific threads respectively.
- Windows does not support signals, but they can be emulated using Asynchronous Procedure Calls ( APCs ). APCs are delivered to specific threads, not processes.

### 4.6.3 Thread Cancellation

- 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.
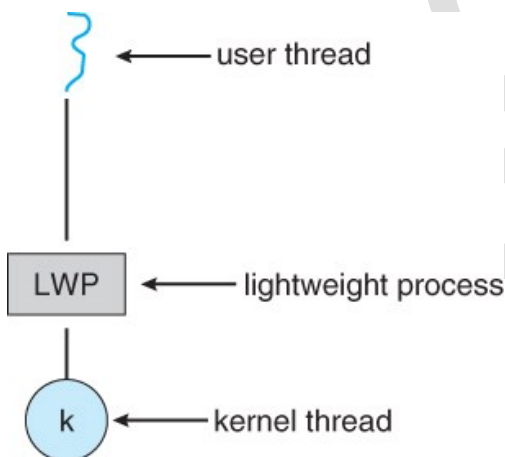
### 4.6.4 Thread-Local Storage ( was 4.4.5 Thread-Specific Data )

- 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.

### 4.6.5 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.



**Figure 4.13 - Lightweight process ( LWP )**

## 6.2 The Critical-Section Problem

- The producer-consumer problem described above is a specific example of a more general situation known as the *critical section* problem. The general idea is that in a number of cooperating processes, each has a critical section of code, with the following conditions and terminologies:
    - Only one process in the group can be allowed to execute in their critical section at any one time. If one process is already executing their critical section and another process wishes to do so, then the second process must

73

be made to wait until the first process has completed their critical section work.

- o The code preceding the critical section, and which controls access to the critical section, is termed the entry section. It acts like a carefully controlled locking door.
- o The code following the critical section is termed the exit section. It generally releases the lock on someone else's door, or at least lets the world know that they are no longer in their critical section.
- o The rest of the code not included in either the critical section or the entry or exit sections is termed the remainder section.

```
do {

    entry section

        critical section

    exit section

        remainder section

} while (TRUE);
```

**Figure 5.1 - General structure of a typical process Pi**

- A solution to the critical section problem must satisfy the following three conditions:
    1. **Mutual Exclusion** - Only one process at a time can be executing in their critical section.
    2. **Progress** - If no process is currently executing in their critical section, and one or more processes want to execute their critical section, then only the processes not in their remainder sections can participate in the decision, and the decision cannot be postponed indefinitely. ( I.e. processes cannot be blocked forever waiting to get into their critical sections. )
    3. **Bounded Waiting** - There exists a limit as to how many other processes can get into their critical sections after a process requests entry into their critical section and before that request is granted. ( I.e. a process requesting entry into their critical section will get a turn eventually, and there is a limit as to how many other processes get to go first. )
- We assume that all processes proceed at a non-zero speed, but no assumptions can be made regarding the *relative* speed of one process versus another.
- Kernel processes can also be subject to race conditions, which can be especially problematic when updating commonly shared kernel data structures such as open file tables or virtual memory management. Accordingly kernels can take on one of two forms:
    1. Non-preemptive kernels do not allow processes to be interrupted while in kernel mode. This eliminates the possibility of kernel-mode race conditions,

but requires kernel mode operations to complete very quickly, and can be problematic for real-time systems, because timing cannot be guaranteed.

2. Preemptive kernels allow for real-time operations, but must be carefully written to avoid race conditions. This can be especially tricky on SMP systems, in which multiple kernel processes may be running simultaneously on different processors.

- Non-preemptive kernels include Windows XP, 2000, traditional UNIX, and Linux prior to 2.6; Preemptive kernels include Linux 2.6 and later, and some commercial UNIXes such as Solaris and IRIX.

## 6.4 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.

- 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.

- 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:

```
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().

## Figures 5.3 and 5.4 illustrate "test_and_set( )" function

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

```
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.

- 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

76

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 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:

```
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);
```

**Figure 5.7 Bounded-waiting mutual exclusion with TestAndSet( ).**

- 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.

## 6.6 Semaphores

- A more robust alternative to simple mutexes is to use *semaphores*, which are integer variables for which only two ( atomic ) operations are defined, the wait and signal operations, as shown in the following figure.
- Note that not only must the variable-changing steps ( S-- and S++ ) be indivisible, it is also necessary that for the wait operation when the test proves false that there be no interruptions before S gets decremented. It IS okay, however, for the busy loop to be interrupted when the test is true, which prevents the system from hanging forever.

```
Wait:
  wait(S) {
      while S <= 0
          ; // no-op
      S--;
  }
```

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

### 6.6.1 Semaphore Usage

- In practice, semaphores can take on one of two forms:
    - **Binary semaphores** can take on one of two values, 0 or 1. They can be used to solve the critical section problem as described above, and can be used as mutexes on systems that do not provide a separate mutex mechanism.. The use of mutexes for this purpose is shown in Figure 6.9 ( from the 8th edition ) below.

```
do {
   waiting(mutex);

      // critical section

   signal(mutex);

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

**Mutual-exclusion implementation with semaphores. ( From 8th edition. )**

    - **Counting semaphores** can take on any integer value, and are usually used to count the number remaining of some limited resource. The counter is initialized to the number of such resources available in the system, and whenever the counting semaphore is greater than zero, then a process can enter a critical section and use one of the resources.

When the counter gets to zero ( or negative in some implementations ), then the process blocks until another process frees up a resource and increments the counting semaphore with a signal call. ( The binary semaphore can be seen as just a special case where the number of resources initially available is just one. )

o Semaphores can also be used to synchronize certain operations between processes. For example, suppose it is important that process P1 execute statement S1 before process P2 executes statement S2.

- First we create a semaphore named synch that is shared by the two processes, and initialize it to zero.
- Then in process P1 we insert the code:

  S1;
  signal( synch );

- and in process P2 we insert the code:

  wait(                    synch                    );
  S2;

- Because synch was initialized to 0, process P2 will block on the wait until after P1 executes the call to signal.

**6.6.2 Semaphore Implementation**

- The big problem with semaphores as described above is the busy loop in the wait call, which consumes CPU cycles without doing any useful work. This type of lock is known as a *spinlock*, because the lock just sits there and spins while it waits. While this is generally a bad thing, it does have the advantage of not invoking context switches, and so it is sometimes used in multi-processing systems when the wait time is expected to be short - One thread spins on one processor while another completes their critical section on another processor.
- An alternative approach is to block a process when it is forced to wait for an available semaphore, and swap it out of the CPU. In this implementation each semaphore needs to maintain a list of processes that are blocked waiting for it, so that one of the processes can be woken up and swapped back in when the semaphore becomes available. ( Whether it gets swapped back into the CPU immediately or whether it needs to hang out in the ready queue for a while is a scheduling problem. )
- The new definition of a semaphore and the corresponding wait and signal operations are shown as follows:

**Semaphore Structure:**

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

**Wait Operation:**

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

**Signal Operation:**

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

- Note that in this implementation the value of the semaphore can actually become negative, in which case its magnitude is the number of processes waiting for that semaphore. This is a result of decrementing the counter before checking its value.
- Key to the success of semaphores is that the wait and signal operations be atomic, that is no other process can execute a wait or signal on the same semaphore at the same time. ( Other processes could be allowed to do other things, including working with other semaphores, they just can't have access to **this** semaphore. ) On single processors this can be implemented by disabling interrupts during the execution of wait and signal; Multiprocessor systems have to use more complex methods, including the use of spinlocking.

### 6.6.3 Deadlocks and Starvation

- One important problem that can arise when using semaphores to block processes waiting for a limited resource is the problem of *deadlocks*, which occur when multiple processes are blocked, each waiting for a resource that can only be freed by one of the other ( blocked ) processes, as illustrated in the following example. ( Deadlocks are covered more completely in chapter 7. )

```
        P₀                  P₁

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

- Another problem to consider is that of *starvation*, in which one or more processes gets blocked forever, and never get a chance to take their turn in the critical section. For example, in the semaphores above, we did not specify the algorithms for adding processes to the waiting queue in the semaphore in the wait( ) call, or selecting one to be removed from the queue in the signal( ) call. If the method chosen is a FIFO queue, then every process will eventually get their turn, but if a LIFO queue is implemented instead, then the first process to start waiting could starve.

### 6.6.4 Priority Inversion

- A challenging scheduling problem arises when a high-priority process gets blocked waiting for a resource that is currently held by a low-priority process.
- If the low-priority process gets pre-empted by one or more medium-priority processes, then the high-priority process is essentially made to wait for the medium priority processes to finish before the low-priority process can release the needed resource, causing a *priority inversion.* If there are enough medium-priority processes, then the high-priority process may be forced to wait for a very long time.
- One solution is a *priority-inheritance protocol,* in which a low-priority process holding a resource for which a high-priority process is waiting will temporarily inherit the high priority from the waiting process. This prevents the medium-priority processes from preempting the low-priority process until it releases the resource, blocking the priority inversion problem.
- The book has an interesting discussion of how a priority inversion almost doomed the Mars Pathfinder mission, and how the problem was solved when the priority inversion was stopped. Full details are available online at http://research.microsoft.com/en-us/um/people/mbj/mars_pathfinder/authoritative_account.htm

### 6.7 Classic Problems of Synchronization

The following classic problems are used to test virtually every new proposed synchronization algorithm.

### 6.7.1 The Bounded-Buffer Problem

- This is a generalization of the producer-consumer problem wherein access is controlled to a shared group of buffers of a limited size.
- In this solution, the two counting semaphores "full" and "empty" keep track of the current number of full and empty buffers respectively ( and initialized to 0 and N respectively. ) The binary semaphore mutex controls access to the critical section. The producer and consumer processes are nearly identical - One can think of the producer as producing full buffers, and the consumer producing empty buffers.

```
do {
    . . .
   // produce an item in nextp
    . . .
   wait(empty);
   wait(mutex);
    . . .
   // add nextp to buffer
    . . .
   signal(mutex);
   signal(full);
}while (TRUE);
```

**Figure 5.9**  The structure of the producer process.

```
do {
   wait(full);
   wait(mutex);
    . . .
   // remove an item from buffer to nextc
    . . .
   signal(mutex);
   signal(empty);
    . . .
   // consume the item in nextc
    . . .
}while (TRUE);
```

**Figure 5.10**  The structure of the consumer process.

## Figures 5.9 and 5.10 use variables next_produced and next_consumed

### 6.7.2 The Readers-Writers Problem

- In the readers-writers problem there are some processes ( termed readers ) who only read the shared data, and never change it, and there are other processes ( termed writers ) who may change the data in addition to or instead of reading it. There is no limit to how many readers can access the

data simultaneously, but when a writer accesses the data, it needs exclusive access.

- There are several variations to the readers-writers problem, most centered around relative priorities of readers versus writers.
    - The *first* readers-writers problem gives priority to readers. In this problem, if a reader wants access to the data, and there is not already a writer accessing it, then access is granted to the reader. A solution to this problem can lead to starvation of the writers, as there could always be more readers coming along to access the data. ( A steady stream of readers will jump ahead of waiting writers as long as there is currently already another reader accessing the data, because the writer is forced to wait until the data is idle, which may never happen if there are enough readers. )
    - The *second* readers-writers problem gives priority to the writers. In this problem, when a writer wants access to the data it jumps to the head of the queue - All waiting readers are blocked, and the writer gets access to the data as soon as it becomes available. In this solution the readers may be starved by a steady stream of writers.
- The following code is an example of the first readers-writers problem, and involves an important counter and two binary semaphores:
    - readcount is used by the reader processes, to count the number of readers currently accessing the data.
    - mutex is a semaphore used only by the readers for controlled access to readcount.
    - rw_mutex is a semaphore used to block and release the writers. The first reader to access the data will set this lock and the last reader to exit will release it; The remaining readers do not touch rw_mutex. ( Eighth edition called this variable wrt. )
    - Note that the first reader to come along will block on rw_mutex if there is currently a writer accessing the data, and that all following readers will only block on mutex for their turn to increment readcount.

```
do {
   wait(rw_mutex);
      . . .
   /* writing is performed */
      . . .
   signal(rw_mutex);
} while (true);
```

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

```
do {
   wait(mutex);
   read_count++;
   if (read_count == 1)
      wait(rw_mutex);
   signal(mutex);
      . . .
   /* reading is performed */
      . . .
   wait(mutex);
   read_count--;
   if (read_count == 0)
      signal(rw_mutex);
   signal(mutex);
} while (true);
```
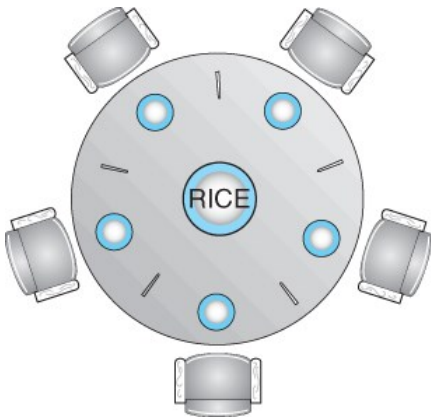
**Figure 5.12**   The structure of a reader process.

- Some hardware implementations provide specific reader-writer locks, which are accessed using an argument specifying whether access is requested for reading or writing. The use of reader-writer locks is beneficial for situation in which: (1) processes can be easily identified as either readers or writers, and (2) there are significantly more readers than writers, making the additional overhead of the reader-writer lock pay off in terms of increased concurrency of the readers.

### 6.7.3 The Dining-Philosophers Problem

- The dining philosophers problem is a classic synchronization problem involving the allocation of limited resources amongst a group of processes in a deadlock-free and starvation-free manner:
    - Consider five philosophers sitting around a table, in which there are five chopsticks evenly distributed and an endless bowl of rice in the center, as shown in the diagram below. ( There is exactly one chopstick between each pair of dining philosophers. )
    - These philosophers spend their lives alternating between two activities: eating and thinking.
    - When it is time for a philosopher to eat, it must first acquire two chopsticks - one from their left and one from their right.

- When a philosopher thinks, it puts down both chopsticks in their original locations.



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

- One possible solution, as shown in the following code section, is to use a set of five semaphores ( chopsticks[ 5 ] ), and to have each hungry philosopher first wait on their left chopstick ( chopsticks[ i ] ), and then wait on their right chopstick ( chopsticks[ ( i + 1 ) % 5 ] )
- But suppose that all five philosophers get hungry at the same time, and each starts by picking up their left chopstick. They then look for their right chopstick, but because it is unavailable, they wait for it, forever, and eventually all the philosophers starve due to the resulting deadlock.

```
do {
   wait(chopstick[i]);
   wait(chopstick[(i+1) % 5]);
      . . .
   // eat
      . . .
   signal(chopstick[i]);
   signal(chopstick[(i+1) % 5]);
      . . .
   // think
      . . .
}while (TRUE);
```

**Figure 5.14 - The structure of philosopher i.**

- Some potential solutions to the problem include:
  - Only allow four philosophers to dine at the same time. ( Limited simultaneous processes. )
  - Allow philosophers to pick up chopsticks only when both are available, in a critical section. ( All or nothing allocation of critical resources. )
  - Use an asymmetric solution, in which odd philosophers pick up their left chopstick first and even philosophers pick up their right chopstick

first. ( Will this solution always work? What if there are an even number of philosophers? )
- Note carefully that a deadlock-free solution to the dining philosophers problem does not necessarily guarantee a starvation-free one. ( While some or even most of the philosophers may be able to get on with their normal lives of eating and thinking, there may be one unlucky soul who never seems to be able to get both chopsticks at the same time. :-(

## 5.8 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**.

### 5.8.1 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.

```
monitor monitor name
{
   // shared variable declarations

   procedure P1 ( . . . ) {
      . . .
   }

   procedure P2 ( . . . ) {
      . . .
   }

        .
        .
        .
   procedure Pn ( . . . ) {
      . . .
   }

   initialization code ( . . . ) {
      . . .
   }
}
```
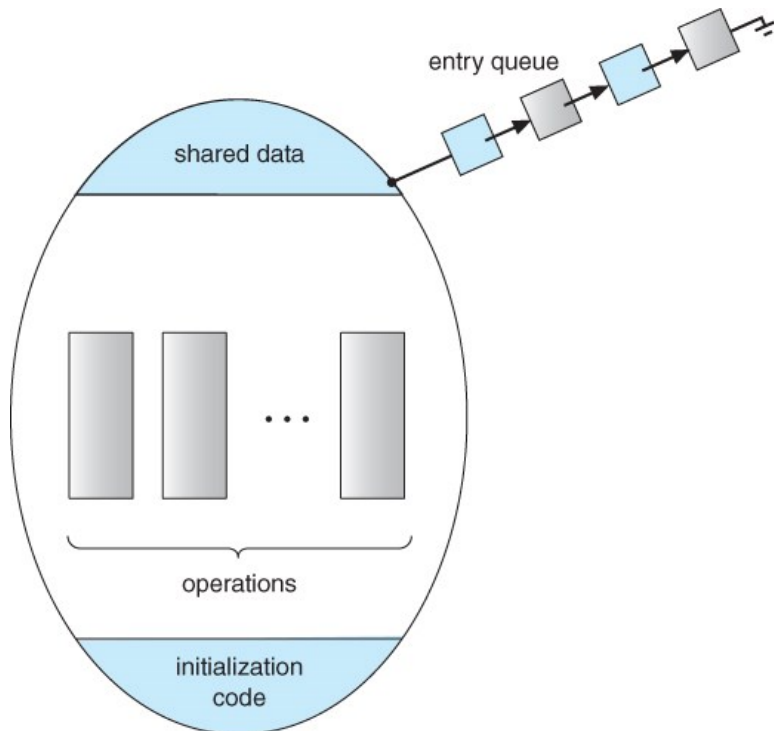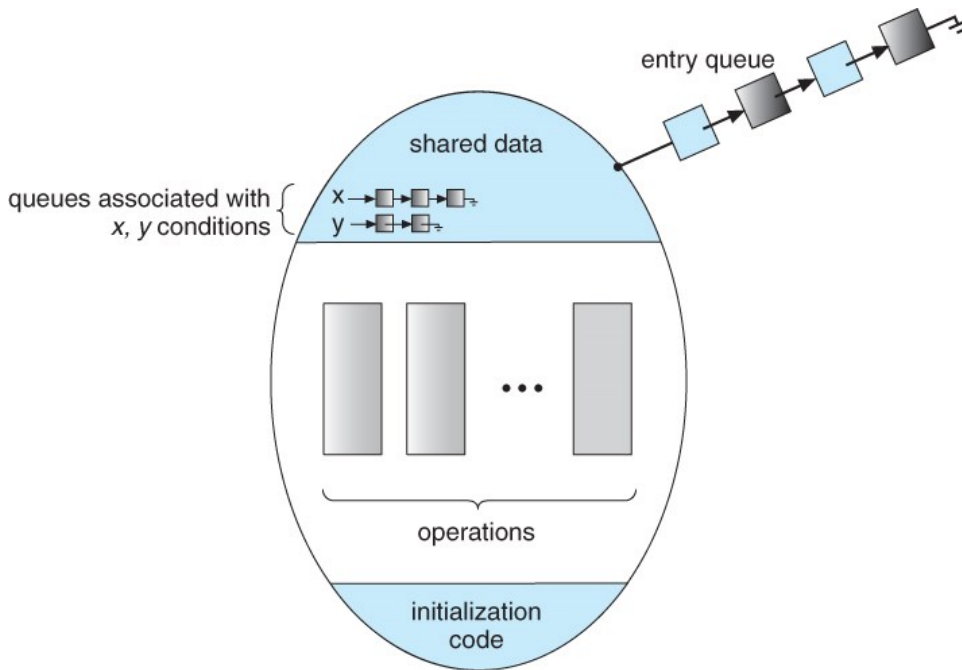
**Figure 5.15 - Syntax of a monitor.**

- Figure 5.16 shows a schematic of a monitor, with an entry queue of processes waiting their turn to execute monitor operations ( methods. )

**Figure 5.16 - 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 until some other process calls signal, and adds the blocked process onto a list associated with that condition.
    - 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. )
- Figure 6.18 below illustrates a monitor that includes condition variables within its data space. Note that the condition variables, along with the list of processes currently waiting for the conditions, are in the data space of the monitor - The processes on these lists are not "in" the monitor, in the sense that they are not executing any code in the monitor.

**Figure 5.17 - Monitor with condition variables**

- 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.

**5.8.2 Dining-Philosophers Solution Using Monitors**

- This solution to the dining philosophers uses monitors, and the restriction that a philosopher may only pick up chopsticks when both are available. There are also two key data structures in use in this solution:
    1. **enum { THINKING, HUNGRY,EATING } state[ 5 ];** A philosopher may only set their state to eating when neither of their adjacent neighbors is eating. ( state[ ( i + 1 ) % 5 ] != EATING && state[ ( i + 4 ) % 5 ] != EATING ).

2. **condition self[ 5 ];** This condition is used to delay a hungry philosopher who is unable to acquire chopsticks.
- In the following solution philosophers share a monitor, DiningPhilosophers, and eat using the following sequence of operations:
  1. DiningPhilosophers.pickup( ) - Acquires chopsticks, which may block the process.
  2. eat
  3. DiningPhilosophers.putdown( ) - Releases the chopsticks.

```
monitor DiningPhilosophers
{
   enum {THINKING, HUNGRY, EATING} state[5];
   condition self[5];

   void pickup(int i) {
      state[i] = HUNGRY;
      test(i);
      if (state[i] != EATING)
         self[i].wait();
   }

   void putdown(int i) {
      state[i] = THINKING;
      test((i + 4) % 5);
      test((i + 1) % 5);
   }

   void test(int i) {
      if ((state[(i + 4) % 5] != EATING) &&
        (state[i] == HUNGRY) &&
        (state[(i + 1) % 5] != EATING)) {
         state[i] = EATING;
         self[i].signal();
      }
   }

   initialization_code() {
      for (int i = 0; i < 5; i++)
         state[i] = THINKING;
   }
}
```

**Figure 5.18**   A monitor solution to the dining-philosopher problem.

### 5.8.3 Implementing a Monitor Using Semaphores

- One possible implementation of a monitor uses a semaphore "mutex" to control mutual exclusionary access to the monitor, and a counting semaphore "next" on which processes can suspend themselves after they are already "inside" the monitor ( in conjunction with condition variables, see below. ) The integer next_count keeps track of how many processes are waiting in

the next queue. Externally accessible monitor processes are then implemented as:

```
wait(mutex);
    ...
   body of F
    ...
if (next_count > 0)
   signal(next);
else
   signal(mutex);
```

- Condition variables can be implemented using semaphores as well. For a condition x, a semaphore "x_sem" and an integer "x_count" are introduced, both initialized to zero. The wait and signal methods are then implemented as follows. ( This approach to the condition implements the signal-and-wait option described above for ensuring that only one process at a time is active inside the monitor. )

**Wait:**

```
   x_count++;
   if (next_count > 0)
      signal(next);
   else
      signal(mutex);
   wait(x_sem);
   x_count--;
```

**Signal:**

```
   if (x_count > 0) {
      next_count++;
      signal(x_sem);
      wait(next);
      next_count--;
   }
```

**5.8.4 Resuming Processes Within a Monitor**

- When there are multiple processes waiting on the same condition within a monitor, how does one decide which one to wake up in response to a signal on that condition? One obvious approach is FCFS, and this may be suitable in many cases.
- Another alternative is to assign ( integer ) priorities, and to wake up the process with the smallest ( best ) priority.
- Figure 5.19 illustrates the use of such a condition within a monitor used for resource allocation. Processes wishing to access this resource must specify

the time they expect to use it using the acquire( time ) method, and must call the release( ) method when they are done with the resource.

```
monitor ResourceAllocator
{
   boolean busy;
   condition x;

   void acquire(int time) {
      if (busy)
         x.wait(time);
      busy = TRUE;
   }

   void release() {
      busy = FALSE;
      x.signal();
   }

   initialization_code() {
      busy = FALSE;
   }
}
```

**Figure 5.19 - A monitor to allocate a single resource.**

- Unfortunately the use of monitors to restrict access to resources still only works if programmers make the requisite acquire and release calls properly. One option would be to place the resource allocation code into the monitor, thereby eliminating the option for programmers to bypass or ignore the monitor, but then that would substitute the monitor's scheduling algorithms for whatever other scheduling algorithms may have been chosen for that particular resource. Chapter 14 on Protection presents more advanced methods for enforcing "nice" cooperation among processes contending for shared resources.
- Concurrent Pascal, Mesa, C#, and Java all implement monitors as described here. Erlang provides concurrency support using a similar mechanism.

## 7.3 Methods for Handling Deadlocks

- Generally speaking there are three ways of handling deadlocks:
  1. Deadlock prevention or avoidance - Do not allow the system to get into a deadlocked state.
  2. Deadlock detection and recovery - Abort a process or preempt some resources when deadlocks are detected.
  3. Ignore the problem all together - If deadlocks only occur once a year or so, it may be better to simply let them happen and reboot as necessary than to incur the constant overhead and system performance penalties associated

with deadlock prevention or detection. This is the approach that both Windows and UNIX take.

- In order to avoid deadlocks, the system must have additional information about all processes. In particular, the system must know what resources a process will or may request in the future. ( Ranging from a simple worst-case maximum to a complete resource request and release plan for each process, depending on the particular algorithm. )
- Deadlock detection is fairly straightforward, but deadlock recovery requires either aborting processes or preempting resources, neither of which is an attractive alternative.
- If deadlocks are neither prevented nor detected, then when a deadlock occurs the system will gradually slow down, as more and more processes become stuck waiting for resources currently held by the deadlock and by other waiting processes. Unfortunately this slowdown can be indistinguishable from a general system slowdown when a real-time process has heavy computing needs.

## 7.4 Deadlock Prevention

- Deadlocks can be prevented by preventing at least one of the four required conditions:

### 7.4.1 Mutual Exclusion

- Shared resources such as read-only files do not lead to deadlocks.
- Unfortunately some resources, such as printers and tape drives, require exclusive access by a single process.

### 7.4.2 Hold and Wait

- To prevent this condition processes must be prevented from holding one or more resources while simultaneously waiting for one or more others. There are several possibilities for this:
  - Require that all processes request all resources at one time. This can be wasteful of system resources if a process needs one resource early in its execution and doesn't need some other resource until much later.
  - Require that processes holding resources must release them before requesting new resources, and then re-acquire the released resources along with the new ones in a single new request. This can be a problem if a process has partially completed an operation using a resource and then fails to get it re-allocated after releasing it.
  - Either of the methods described above can lead to starvation if a process requires one or more popular resources.

### 7.4.3 No Preemption

- Preemption of process resource allocations can prevent this condition of deadlocks, when it is possible.

- One approach is that if a process is forced to wait when requesting a new resource, then all other resources previously held by this process are implicitly released, ( preempted ), forcing this process to re-acquire the old resources along with the new resources in a single request, similar to the previous discussion.
- Another approach is that when a resource is requested and not available, then the system looks to see what other processes currently have those resources *and* are themselves blocked waiting for some other resource. If such a process is found, then some of their resources may get preempted and added to the list of resources for which the process is waiting.
- Either of these approaches may be applicable for resources whose states are easily saved and restored, such as registers and memory, but are generally not applicable to other devices such as printers and tape drives.

### 7.4.4 Circular Wait

- One way to avoid circular wait is to number all resources, and to require that processes request resources only in strictly increasing ( or decreasing ) order.
- In other words, in order to request resource $R_j$, a process must first release all $R_i$ such that $i >= j$.
- One big challenge in this scheme is determining the relative ordering of the different resources
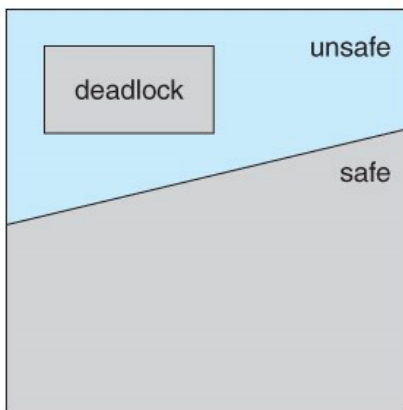
### 7.5 Deadlock Avoidance

- The general idea behind deadlock avoidance is to prevent deadlocks from ever happening, by preventing at least one of the aforementioned conditions.
- This requires more information about each process, AND tends to lead to low device utilization. ( I.e. it is a conservative approach. )
- In some algorithms the scheduler only needs to know the *maximum* number of each resource that a process might potentially use. In more complex algorithms the scheduler can also take advantage of the *schedule* of exactly what resources may be needed in what order.
- When a scheduler sees that starting a process or granting resource requests may lead to future deadlocks, then that process is just not started or the request is not granted.
- A resource allocation *state* is defined by the number of available and allocated resources, and the maximum requirements of all processes in the system.

### 7.5.1 Safe State

- A state is *safe* if the system can allocate all resources requested by all processes ( up to their stated maximums ) without entering a deadlock state.
- More formally, a state is safe if there exists a *safe sequence* of processes { P0, P1, P2, ..., PN } such that all of the resource requests for Pi can be

93

granted using the resources currently allocated to Pi and all processes Pj where j < i. ( I.e. if all the processes prior to Pi finish and free up their resources, then Pi will be able to finish also, using the resources that they have freed up. )

- If a safe sequence does not exist, then the system is in an unsafe state, which **MAY** lead to deadlock. ( All safe states are deadlock free, but not all unsafe states lead to deadlocks. )



**Figure 7.6 - Safe, unsafe, and deadlocked state spaces.**

- For example, consider a system with 12 tape drives, allocated as follows. Is this a safe state? What is the safe sequence?

|  | Maximum Needs | Current Allocation |
|---|---|---|
| **P0** | 10 | 5 |
| **P1** | 4 | 2 |
| **P2** | 9 | 2 |

- What happens to the above table if process P2 requests and is granted one more tape drive?
- Key to the safe state approach is that when a request is made for resources, the request is granted only if the resulting allocation state is a safe one.

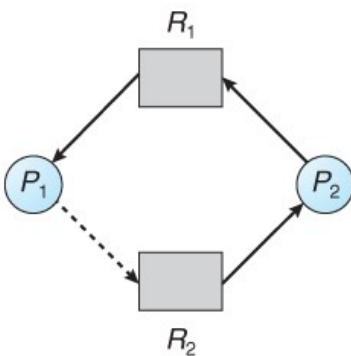**7.5.2 Resource-Allocation Graph Algorithm**

- If resource categories have only single instances of their resources, then deadlock states can be detected by cycles in the resource-allocation graphs.
- In this case, unsafe states can be recognized and avoided by augmenting the resource-allocation graph with *claim edges*, noted by dashed lines, which point from a process to a resource that it may request in the future.
- In order for this technique to work, all claim edges must be added to the graph for any particular process before that process is allowed to request any resources. ( Alternatively, processes may only make requests for resources

for which they have already established claim edges, and claim edges cannot be added to any process that is currently holding resources. )

- When a process makes a request, the claim edge Pi->Rj is converted to a request edge. Similarly when a resource is released, the assignment reverts back to a claim edge.
- This approach works by denying requests that would produce cycles in the resource-allocation graph, taking claim edges into effect.
- Consider for example what happens when process P2 requests resource R2:



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

- The resulting resource-allocation graph would have a cycle in it, and so the request cannot be granted.



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

### 7.5.3 Banker's Algorithm

- For resource categories that contain more than one instance the resource-allocation graph method does not work, and more complex ( and less efficient ) methods must be chosen.
- The Banker's Algorithm gets its name because it is a method that bankers could use to assure that when they lend out resources they will still be able to satisfy all their clients. ( A banker won't loan out a little money to start building a house unless they are assured that they will later be able to loan out the rest of the money to finish the house. )

- When a process starts up, it must state in advance the maximum allocation of resources it may request, up to the amount available on the system.
- When a request is made, the scheduler determines whether granting the request would leave the system in a safe state. If not, then the process must wait until the request can be granted safely.
- The banker's algorithm relies on several key data structures: ( where n is the number of processes and m is the number of resource categories. )
  - Available[ m ] indicates how many resources are currently available of each type.
  - Max[ n ][ m ] indicates the maximum demand of each process of each resource.
  - Allocation[ n ][ m ] indicates the number of each resource category allocated to each process.
  - Need[ n ][ m ] indicates the remaining resources needed of each type for each process. ( Note that Need[ i ][ j ] = Max[ i ][ j ] - Allocation[ i ][ j ] for all i, j. )
- For simplification of discussions, we make the following notations / observations:
  - One row of the Need vector, Need[ i ], can be treated as a vector corresponding to the needs of process i, and similarly for Allocation and Max.
  - A vector X is considered to be <= a vector Y if X[ i ] <= Y[ i ] for all i.

**7.5.3.1 Safety Algorithm**

- In order to apply the Banker's algorithm, we first need an algorithm for determining whether or not a particular state is safe.
- This algorithm determines if the current state of a system is safe, according to the following steps:
  1. Let Work and Finish be vectors of length m and n respectively.
     - Work is a working copy of the available resources, which will be modified during the analysis.
     - Finish is a vector of booleans indicating whether a particular process can finish. ( or has finished so far in the analysis. )
     - Initialize Work to Available, and Finish to false for all elements.
  2. Find an i such that both (A) Finish[ i ] == false, and (B) Need[ i ] < Work. This process has not finished, but could with the given available working set. If no such i exists, go to step 4.
  3. Set Work = Work + Allocation[ i ], and set Finish[ i ] to true. This corresponds to process i finishing up and releasing its resources back into the work pool. Then loop back to step 2.
  4. If finish[ i ] == true for all i, then the state is a safe state, because a safe sequence has been found.

- ( JTB's Modification:
    1. In step 1. instead of making Finish an array of booleans initialized to false, make it an array of ints initialized to 0. Also initialize an int s = 0 as a step counter.
    2. In step 2, look for Finish[ i ] == 0.
    3. In step 3, set Finish[ i ] to ++s. S is counting the number of finished processes.
    4. For step 4, the test can be either Finish[ i ] > 0 for all i, or s >= n. The benefit of this method is that if a safe state exists, then Finish[ ] indicates one safe sequence ( of possibly many. ) )

**7.5.3.2 Resource-Request Algorithm ( The Bankers Algorithm )**

- Now that we have a tool for determining if a particular state is safe or not, we are now ready to look at the Banker's algorithm itself.
- This algorithm determines if a new request is safe, and grants it only if it is safe to do so.
- When a request is made ( that does not exceed currently available resources ), pretend it has been granted, and then see if the resulting state is a safe one. If so, grant the request, and if not, deny the request, as follows:
  1. Let Request[ n ][ m ] indicate the number of resources of each type currently requested by processes. If Request[ i ] > Need[ i ] for any process i, raise an error condition.
  2. If Request[ i ] > Available for any process i, then that process must wait for resources to become available. Otherwise the process can continue to step 3.
  3. Check to see if the request can be granted safely, by pretending it has been granted and then seeing if the resulting state is safe. If so, grant the request, and if not, then the process must wait until its request can be granted safely.The procedure for granting a request ( or pretending to for testing purposes ) is:
     - Available = Available - Request
     - Allocation = Allocation + Request
     - Need = Need - Request

**7.5.3.3 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 |

- 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 |  |

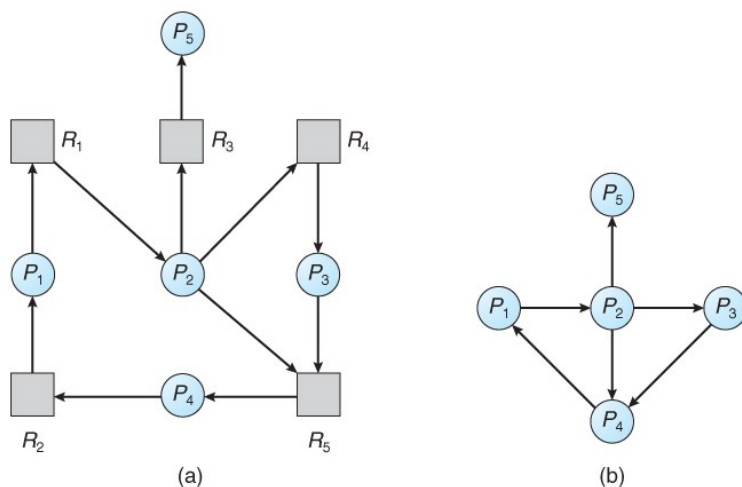- What about requests of ( 3, 3,0 ) by P4? or ( 0, 2, 0 ) by P0? Can these be safely granted? Why or why not?

## 7.6 Deadlock Detection

- If deadlocks are not avoided, then another approach is to detect when they have occurred and recover somehow.
- In addition to the performance hit of constantly checking for deadlocks, a policy / algorithm must be in place for recovering from deadlocks, and there is potential for lost work when processes must be aborted or have their resources preempted.

### 7.6.1 Single Instance of Each Resource Type

- If each resource category has a single instance, then we can use a variation of the resource-allocation graph known as a *wait-for graph*.
- A wait-for graph can be constructed from a resource-allocation graph by eliminating the resources and collapsing the associated edges, as shown in the figure below.
- An arc from Pi to Pj in a wait-for graph indicates that process Pi is waiting for a resource that process Pj is currently holding.



**Figure 7.9 - (a) Resource allocation graph. (b) Corresponding wait-for graph**

- As before, cycles in the wait-for graph indicate deadlocks.

- This algorithm must maintain the wait-for graph, and periodically search it for cycles.

**7.6.2 Several Instances of a Resource Type**

- The detection algorithm outlined here is essentially the same as the Banker's algorithm, with two subtle differences:
  - In step 1, the Banker's Algorithm sets Finish[ i ] to false for all i. The algorithm presented here sets Finish[ i ] to false only if Allocation[ i ] is not zero. If the currently allocated resources for this process are zero, the algorithm sets Finish[ i ] to true. This is essentially assuming that IF all of the other processes can finish, then this process can finish also. Furthermore, this algorithm is specifically looking for which processes are involved in a deadlock situation, and a process that does not have any resources allocated cannot be involved in a deadlock, and so can be removed from any further consideration.
  - Steps 2 and 3 are unchanged
  - In step 4, the basic Banker's Algorithm says that if Finish[ i ] == true for all i, that there is no deadlock. This algorithm is more specific, by stating that if Finish[ i ] == false for any process Pi, then that process is specifically involved in the deadlock which has been detected.
- ( Note: An alternative method was presented above, in which Finish held integers instead of booleans. This vector would be initialized to all zeros, and then filled with increasing integers as processes are detected which can finish. If any processes are left at zero when the algorithm completes, then there is a deadlock, and if not, then the integers in finish describe a safe sequence. To modify this algorithm to match this section of the text, processes with allocation = zero could be filled in with N, N - 1, N - 2, etc. in step 1, and any processes left with Finish = 0 in step 4 are the deadlocked processes. )
- Consider, for example, the following state, and determine if it is currently deadlocked:

|       | Allocation A B C | Request A B C | Available 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         |                 |

- Now suppose that process P2 makes a request for an additional instance of type C, yielding the state shown below. Is the system now deadlocked?

|       | 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 1   |           |
| $P_3$ | 2 1 1      | 1 0 0   |           |
| $P_4$ | 0 0 2      | 0 0 2   |           |

### 7.6.3 Detection-Algorithm Usage

- When should the deadlock detection be done? Frequently, or infrequently?
- The answer may depend on how frequently deadlocks are expected to occur, as well as the possible consequences of not catching them immediately. ( If deadlocks are not removed immediately when they occur, then more and more processes can "back up" behind the deadlock, making the eventual task of unblocking the system more difficult and possibly damaging to more processes. )
- There are two obvious approaches, each with trade-offs:
  1. Do deadlock detection after every resource allocation which cannot be immediately granted. This has the advantage of detecting the deadlock right away, while the minimum number of processes are involved in the deadlock. ( One might consider that the process whose request triggered the deadlock condition is the "cause" of the deadlock, but realistically all of the processes in the cycle are equally responsible for the resulting deadlock. ) The down side of this approach is the extensive overhead and performance hit caused by checking for deadlocks so frequently.
  2. Do deadlock detection only when there is some clue that a deadlock may have occurred, such as when CPU utilization reduces to 40% or some other magic number. The advantage is that deadlock detection is done much less frequently, but the down side is that it becomes impossible to detect the processes involved in the original deadlock, and so deadlock recovery can be more complicated and damaging to more processes.
  3. ( As I write this, a third alternative comes to mind: Keep a historical log of resource allocations, since that last known time of no deadlocks. Do deadlock checks periodically ( once an hour or when CPU usage is low?), and then use the historical log to trace through and determine when the deadlock occurred and what processes caused the initial deadlock. Unfortunately I'm not certain that breaking the original deadlock would then free up the resulting log jam. )

## 7.7 Recovery From Deadlock

- There are three basic approaches to recovery from deadlock:

1. Inform the system operator, and allow him/her to take manual intervention.
2. Terminate one or more processes involved in the deadlock
3. Preempt resources.

## 7.7.1 Process Termination

- Two basic approaches, both of which recover resources allocated to terminated processes:
  - Terminate all processes involved in the deadlock. This definitely solves the deadlock, but at the expense of terminating more processes than would be absolutely necessary.
  - Terminate processes one by one until the deadlock is broken. This is more conservative, but requires doing deadlock detection after each step.
- In the latter case there are many factors that can go into deciding which processes to terminate next:
  - Process priorities.
  - How long the process has been running, and how close it is to finishing.
  - How many and what type of resources is the process holding. ( Are they easy to preempt and restore? )
  - How many more resources does the process need to complete.
  - How many processes will need to be terminated
  - Whether the process is interactive or batch.
  - ( Whether or not the process has made non-restorable changes to any resource. )

## 7.7.2 Resource Preemption

- When preempting resources to relieve deadlock, there are three important issues to be addressed:
  1. **Selecting a victim** - Deciding which resources to preempt from which processes involves many of the same decision criteria outlined above.
  2. **Rollback** - Ideally one would like to roll back a preempted process to a safe state prior to the point at which that resource was originally allocated to the process. Unfortunately it can be difficult or impossible to determine what such a safe state is, and so the only safe rollback is to roll back all the way back to the beginning. ( I.e. abort the process and make it start over. )
  3. **Starvation** - How do you guarantee that a process won't starve because its resources are constantly being preempted? One option would be to use a priority system, and increase the priority of a process every time its resources get preempted. Eventually it should get a high enough priority that it won't get preempted any more.

## 7.5.3.2 Resource-Request Algorithm ( The Bankers Algorithm )

- Now that we have a tool for determining if a particular state is safe or not, we are now ready to look at the Banker's algorithm itself.
- This algorithm determines if a new request is safe, and grants it only if it is safe to do so.
- When a request is made ( that does not exceed currently available resources ), pretend it has been granted, and then see if the resulting state is a safe one. If so, grant the request, and if not, deny the request, as follows:
    1. Let Request[ n ][ m ] indicate the number of resources of each type currently requested by processes. If Request[ i ] > Need[ i ] for any process i, raise an error condition.
    2. If Request[ i ] > Available for any process i, then that process must wait for resources to become available. Otherwise the process can continue to step 3.
    3. Check to see if the request can be granted safely, by pretending it has been granted and then seeing if the resulting state is safe. If so, grant the request, and if not, then the process must wait until its request can be granted safely.The procedure for granting a request ( or pretending to for testing purposes ) is:
        - Available = Available - Request
        - Allocation = Allocation + Request
        - Need = Need - Request

### 7.5.3.3 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 |

- 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 |          |

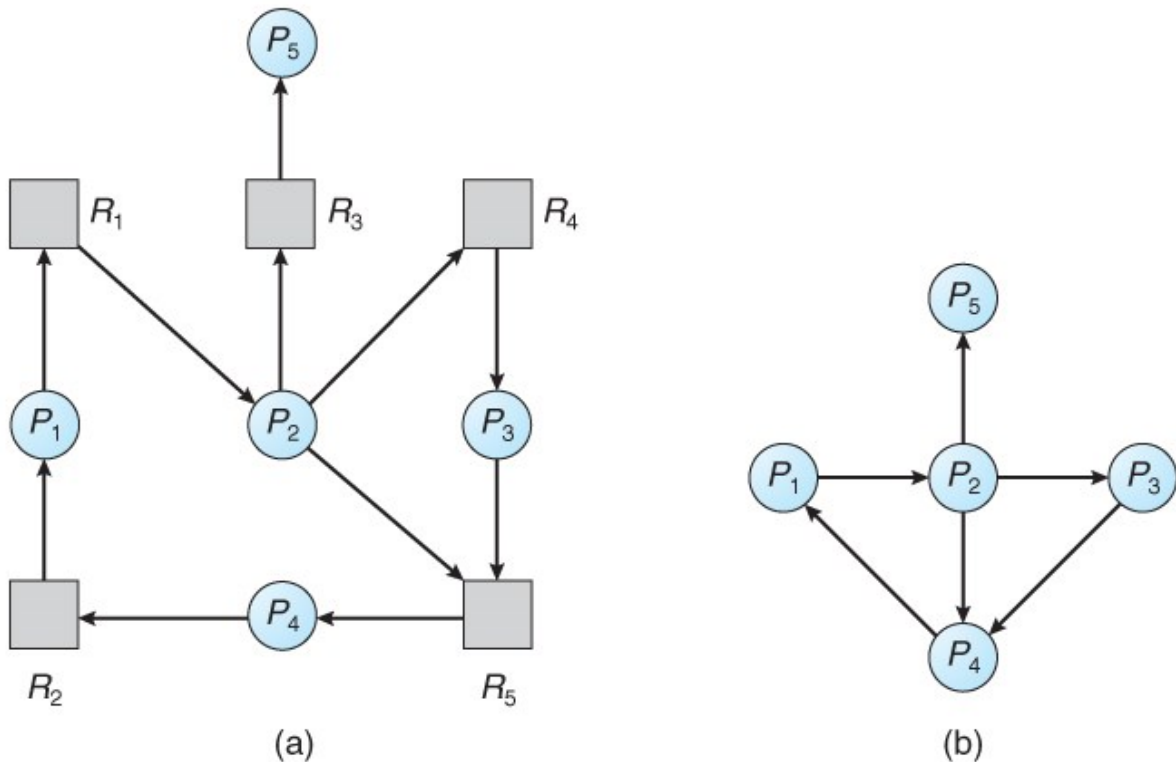- What about requests of ( 3, 3,0 ) by P4? or ( 0, 2, 0 ) by P0? Can these be safely granted? Why or why not?

## 7.6 Deadlock Detection

- If deadlocks are not avoided, then another approach is to detect when they have occurred and recover somehow.
- In addition to the performance hit of constantly checking for deadlocks, a policy / algorithm must be in place for recovering from deadlocks, and there is potential for lost work when processes must be aborted or have their resources preempted.

### 7.6.1 Single Instance of Each Resource Type

- If each resource category has a single instance, then we can use a variation of the resource-allocation graph known as a *wait-for graph*.
- A wait-for graph can be constructed from a resource-allocation graph by eliminating the resources and collapsing the associated edges, as shown in the figure below.
- An arc from Pi to Pj in a wait-for graph indicates that process Pi is waiting for a resource that process Pj is currently holding.

**Figure 7.9 - (a) Resource allocation graph. (b) Corresponding wait-for graph**

- As before, cycles in the wait-for graph indicate deadlocks.
- This algorithm must maintain the wait-for graph, and periodically search it for cycles.

### 7.6.2 Several Instances of a Resource Type

- The detection algorithm outlined here is essentially the same as the Banker's algorithm, with two subtle differences:
  - In step 1, the Banker's Algorithm sets Finish[ i ] to false for all i. The algorithm presented here sets Finish[ i ] to false only if Allocation[ i ] is not zero. If the currently allocated resources for this process are zero, the algorithm sets Finish[ i ] to true. This is essentially assuming that IF all of the other processes can finish, then this process can finish also. Furthermore, this algorithm is specifically looking for which processes are involved in a deadlock situation, and a process that does not have any resources allocated cannot be involved in a deadlock, and so can be removed from any further consideration.
  - Steps 2 and 3 are unchanged
  - In step 4, the basic Banker's Algorithm says that if Finish[ i ] == true for all i, that there is no deadlock. This algorithm is more specific, by stating that if Finish[ i ] == false for any process Pi, then that process is specifically involved in the deadlock which has been detected.
- ( Note: An alternative method was presented above, in which Finish held integers instead of booleans. This vector would be initialized to all zeros, and then filled with increasing integers as processes are detected which can

105

finish. If any processes are left at zero when the algorithm completes, then there is a deadlock, and if not, then the integers in finish describe a safe sequence. To modify this algorithm to match this section of the text, processes with allocation = zero could be filled in with N, N - 1, N - 2, etc. in step 1, and any processes left with Finish = 0 in step 4 are the deadlocked processes. )

- Consider, for example, the following state, and determine if it is currently deadlocked:

|  | 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 |  |

- Now suppose that process P2 makes a request for an additional instance of type C, yielding the state shown below. Is the system now deadlocked?

|  | 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 1 |  |
| $P_3$ | 2 1 1 | 1 0 0 |  |
| $P_4$ | 0 0 2 | 0 0 2 |  |

### 7.6.3 Detection-Algorithm Usage

- When should the deadlock detection be done? Frequently, or infrequently?
- The answer may depend on how frequently deadlocks are expected to occur, as well as the possible consequences of not catching them immediately. ( If deadlocks are not removed immediately when they occur, then more and more processes can "back up" behind the deadlock, making the eventual task of unblocking the system more difficult and possibly damaging to more processes. )
- There are two obvious approaches, each with trade-offs:

1. Do deadlock detection after every resource allocation which cannot be immediately granted. This has the advantage of detecting the deadlock right away, while the minimum number of processes are involved in the deadlock. ( One might consider that the process whose request triggered the deadlock condition is the "cause" of the deadlock, but realistically all of the processes in the cycle are equally responsible for the resulting deadlock. ) The down side of this approach is the extensive overhead and performance hit caused by checking for deadlocks so frequently.
2. Do deadlock detection only when there is some clue that a deadlock may have occurred, such as when CPU utilization reduces to 40% or some other magic number. The advantage is that deadlock detection is done much less frequently, but the down side is that it becomes impossible to detect the processes involved in the original deadlock, and so deadlock recovery can be more complicated and damaging to more processes.
3. ( As I write this, a third alternative comes to mind: Keep a historical log of resource allocations, since that last known time of no deadlocks. Do deadlock checks periodically ( once an hour or when CPU usage is low?), and then use the historical log to trace through and determine when the deadlock occurred and what processes caused the initial deadlock. Unfortunately I'm not certain that breaking the original deadlock would then free up the resulting log jam. )

## 7.7 Recovery From Deadlock

- There are three basic approaches to recovery from deadlock:
  1. Inform the system operator, and allow him/her to take manual intervention.
  2. Terminate one or more processes involved in the deadlock
  3. Preempt resources.

### 7.7.1 Process Termination

- Two basic approaches, both of which recover resources allocated to terminated processes:
  o Terminate all processes involved in the deadlock. This definitely solves the deadlock, but at the expense of terminating more processes than would be absolutely necessary.
  o Terminate processes one by one until the deadlock is broken. This is more conservative, but requires doing deadlock detection after each step.
- In the latter case there are many factors that can go into deciding which processes to terminate next:
  1. Process priorities.
  2. How long the process has been running, and how close it is to finishing.

3. How many and what type of resources is the process holding. ( Are they easy to preempt and restore? )
4. How many more resources does the process need to complete.
5. How many processes will need to be terminated
6. Whether the process is interactive or batch.
7. ( Whether or not the process has made non-restorable changes to any resource. )

**7.7.2 Resource Preemption**

- When preempting resources to relieve deadlock, there are three important issues to be addressed:
    1. **Selecting a victim** - Deciding which resources to preempt from which processes involves many of the same decision criteria outlined above.
    2. **Rollback** - Ideally one would like to roll back a preempted process to a safe state prior to the point at which that resource was originally allocated to the process. Unfortunately it can be difficult or impossible to determine what such a safe state is, and so the only safe rollback is to roll back all the way back to the beginning. ( I.e. abort the process and make it start over. )
    3. **Starvation** - How do you guarantee that a process won't starve because its resources are constantly being preempted? One option would be to use a priority system, and increase the priority of a process every time its resources get preempted. Eventually it should get a high enough priority that it won't get preempted any more.

# UNIT III

## 8.2 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.***

### 8.2.1 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 OSes 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**

### 8.2.2 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.

## 8.3 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

110

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. )

**8.3.1 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**

**8.3.2 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.

### 8.3.3. 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.

### 8.4 Segmentation

**8.4.1 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:



logical address

**Figure 8.7 Programmer's view of a program.**

**8.4.2 Segmentation Hardware**

- A *segment table* maps segment-offset addresses to physical addresses, and simultaneously checks for invalid addresses, using a system similar to the page tables and relocation base registers discussed previously. ( Note that at

113

this point in the discussion of segmentation, each segment is kept in contiguous memory and may be of different sizes, but that segmentation can also be combined with paging as we shall see shortly. )



**Figure 8.8 - Segmentation hardware**

**Figure 8.9 - Example of segmentation**

## 8.5 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.

### 8.5.1 Basic Method

- The basic idea behind paging is to divide physical memory into a number of equal sized blocks called *frames*, and to divide a 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.

| page number | page offset |
|:---:|:---:|
| $p$ | $d$ |
| $m - n$ | $n$ |

- ( 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. )
- 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. )



**Figure 8.13 - Free frames (a) before allocation and (b) after allocation**

**8.5.2 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*.
- **( Eighth Edition Version: )** For example, suppose that it takes 100 nanoseconds to access main memory, and only 20

nanoseconds to search the TLB. So a TLB hit takes 120 nanoseconds total ( 20 to find the frame number and then another 100 to go get the data ), and a TLB miss takes 220 ( 20 to search the TLB, 100 to go get the frame number, and then another 100 to go get the data. ) So with an 80% TLB hit ratio, the average memory access time would be:

0.80 * 120 + 0.20 * 220 = 140 nanoseconds

for a 40% slowdown to get the frame number. A 98% hit rate would yield 122 nanoseconds average access time ( you should verify this ), for a 22% slowdown.

- **( Ninth Edition Version: )** The ninth edition ignores the 20 nanoseconds required to search the TLB, yielding

0.80 * 100 + 0.20 * 200 = 120 nanoseconds

for a 20% slowdown to get the frame number. A 99% hit rate would yield 101 nanoseconds average access time ( you should verify this ), for a 1% slowdown.

## 8.5.3 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.

**Figure 8.15 - Valid (v) or invalid (i) bit in page table**

**8.5.4 Shared Pages**

- Paging systems can make it very easy to share blocks of memory, by simply duplicating page numbers in multiple page frames. This may be done with either code or data.
- If code is *reentrant*, that means that it does not write to or change the code in any way ( it is non self-modifying ), and it is therefore safe to re-enter it. More importantly, it means the code can be shared by multiple processes, so long as each has their own copy of the data and registers, including the instruction register.
- In the example given below, three different users are running the editor simultaneously, but the code is only loaded into memory ( in the page frames ) one time.
- Some systems also implement shared memory in this fashion.

**Figure 8.16 - Sharing of code in a paging environment**

## 8.6 Structure of the Page Table

### 8.6.1 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 amounts 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**

logical address

| $p_1$ | $p_2$ | d |

$p_1\{$
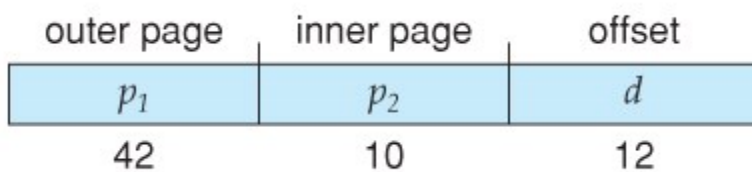
outer page
table

$p_2\{$

page of
page table

$d\{$

**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:

| section | page | offset |
|---------|------|--------|
| $s$ | $p$ | $d$ |
| 2 | 21 | 9 |

- 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.

| 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.
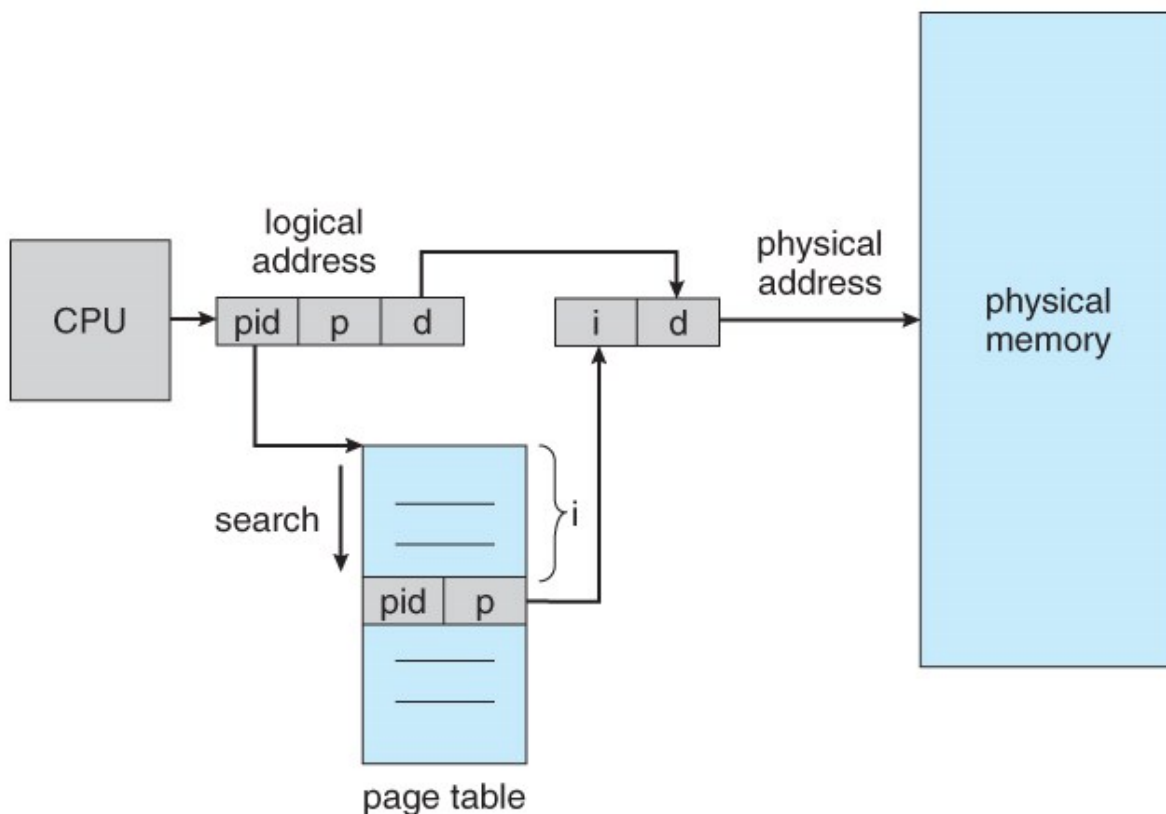
**8.6.2 Hashed Page Tables**

- 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**

**8.6.3 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**

### 8.6.4 Oracle SPARC Solaris ( Optional, New Section in 9th Edition )

Consider as a final example a modern 64-bit CPU and operating system that are tightly integrated to provide low-overhead virtual memory. Solaris running on the SPARC CPU is a fully 64-bit operating system and as such has to solve the problem of virtual memory without using up all of its physical memory by keeping multiple levels of page tables. Its approach is a bit complex but solves the problem efficiently using hashed page tables. There are two hash tables—one for the kernel and one for all user processes. Each maps memory addresses from virtual to physical memory. Each hash-table entry represents a contiguous area of mapped virtual memory, which is more efficient than having a separate hash-table entry for each page. Each entry has a base address and a span indicating the number of pages the entry represents. Virtual-to-physical translation would take too long if each address required searching through a hash table, so the CPU implements a TLB that holds translation table entries (TTEs) for fast hardware lookups. A cache of these TTEs reside in a translation storage buffer (TSB), which includes an entry per recently accessed page. When a virtual address reference occurs, the hardware searches the TLB for a translation. If none is found, the hardware walks through the in-memory TSB looking for the TTE that corresponds to the virtual address that caused the lookup. This TLB walk functionality is found on many modern CPUs. If a match is found in the TSB, the CPU copies the TSB entry into the TLB, and the memory translation completes. If no match is found in the TSB, the kernel is interrupted to search the hash table. The kernel then creates a TTE from the appropriate hash table and stores it in the TSB for automatic loading into the TLB by the CPU memory-management unit. Finally, the interrupt handler returns

control to the MMU, which completes the address translation and retrieves the requested byte or word from main memory.

## 9.2 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.
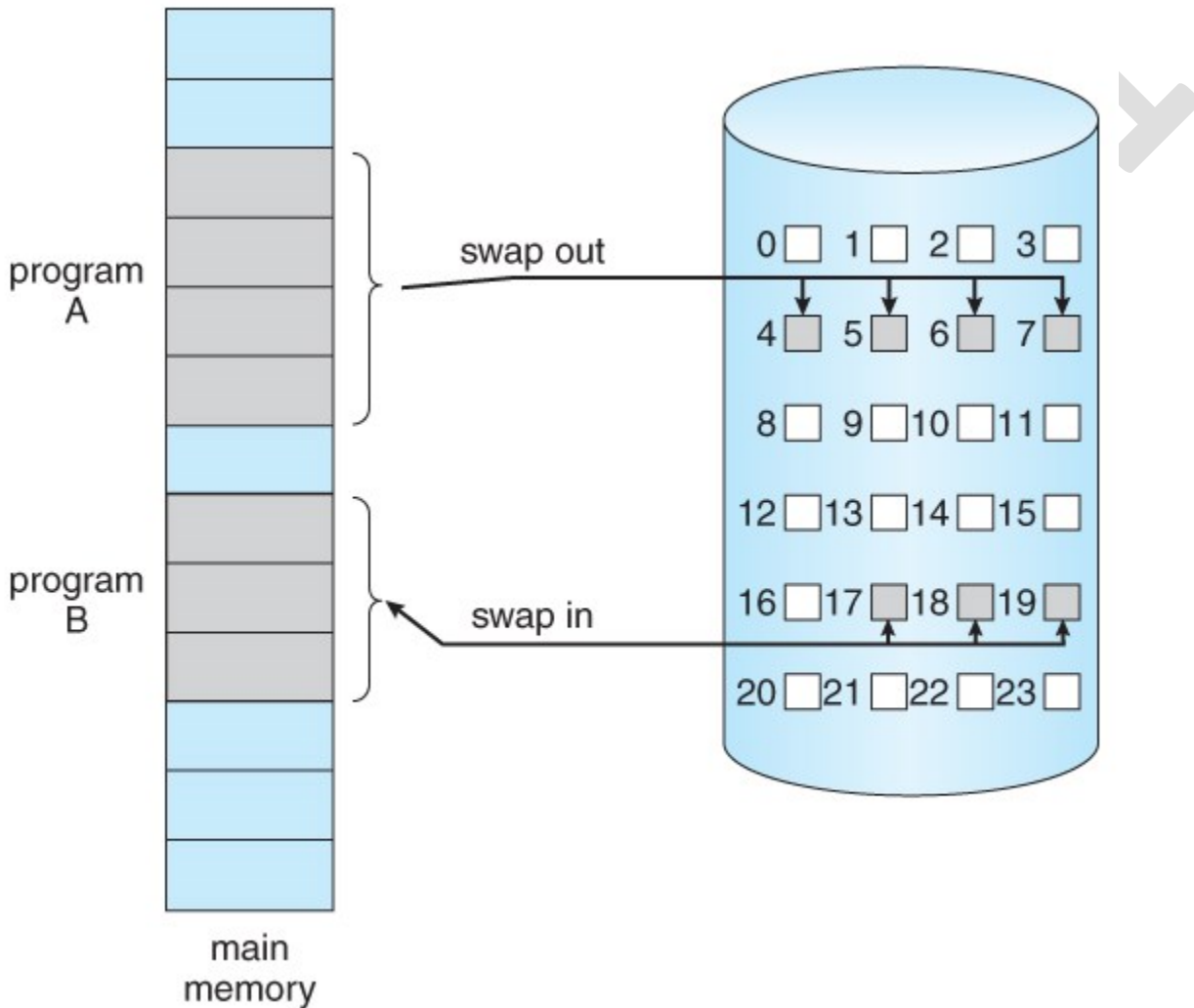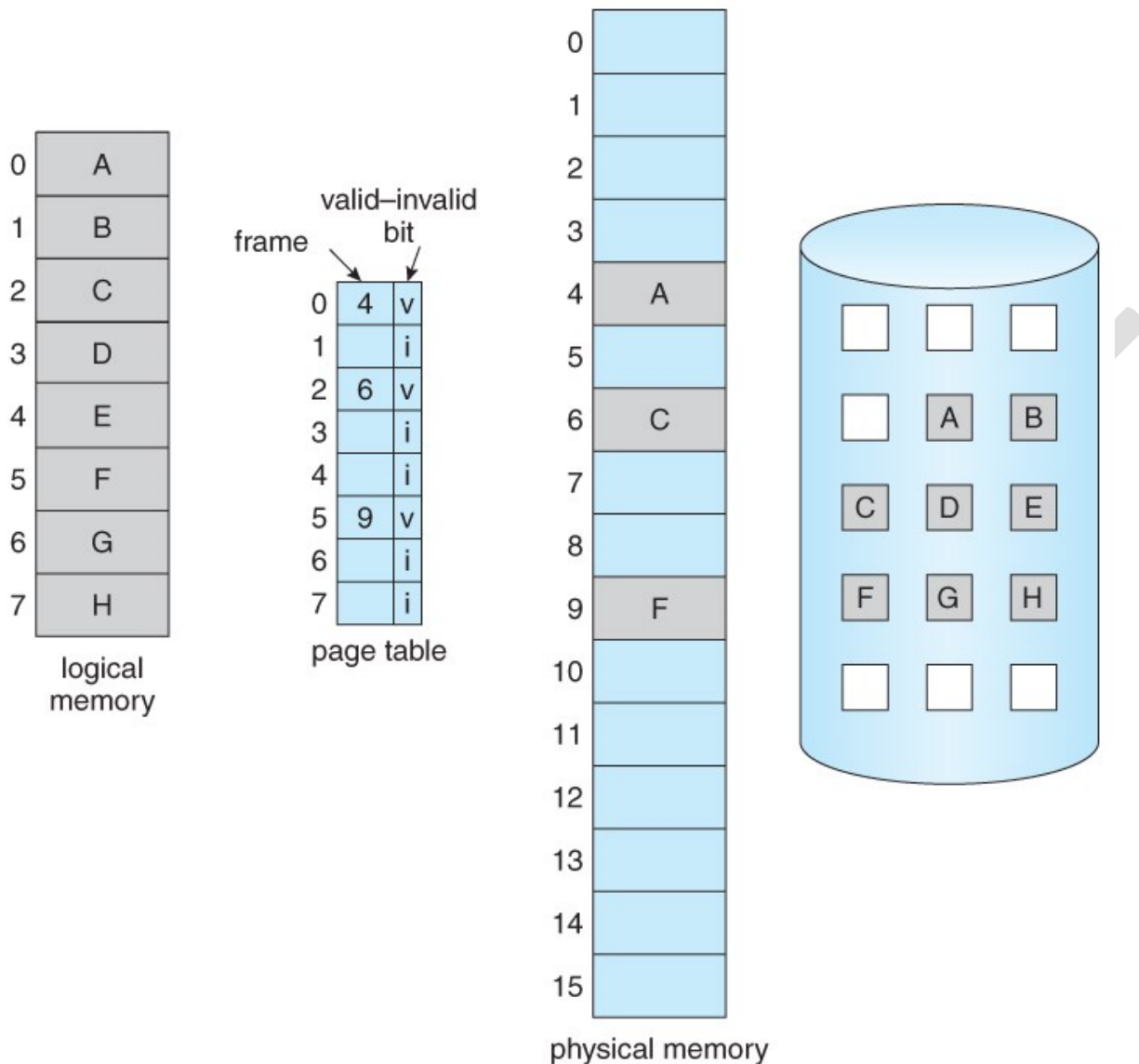


**Figure 9.4 - Transfer of a paged memory to contiguous disk space**

### 9.2.1 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. )
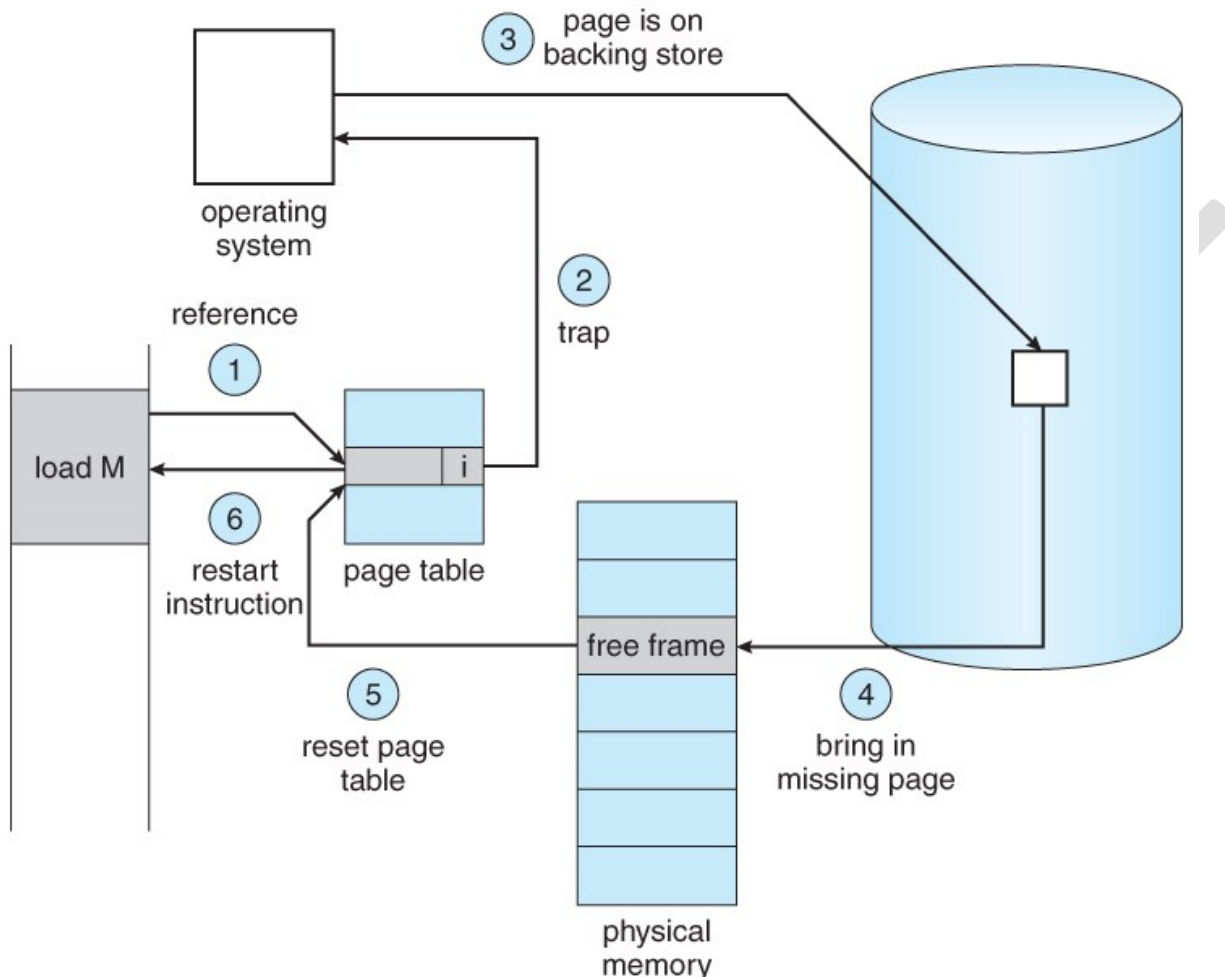
- If the process only ever accesses pages that are loaded in memory ( **memory resident** pages ), then the process runs exactly as if all the pages were loaded in to memory.



**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*, covered in section 9.6.1.
- 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.

**9.2.2 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.

**9.3 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.
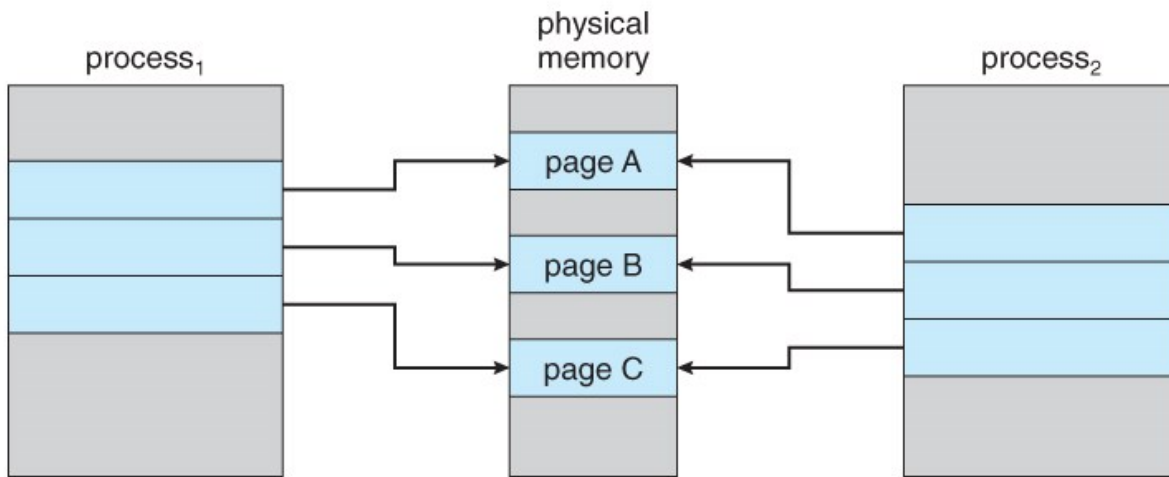
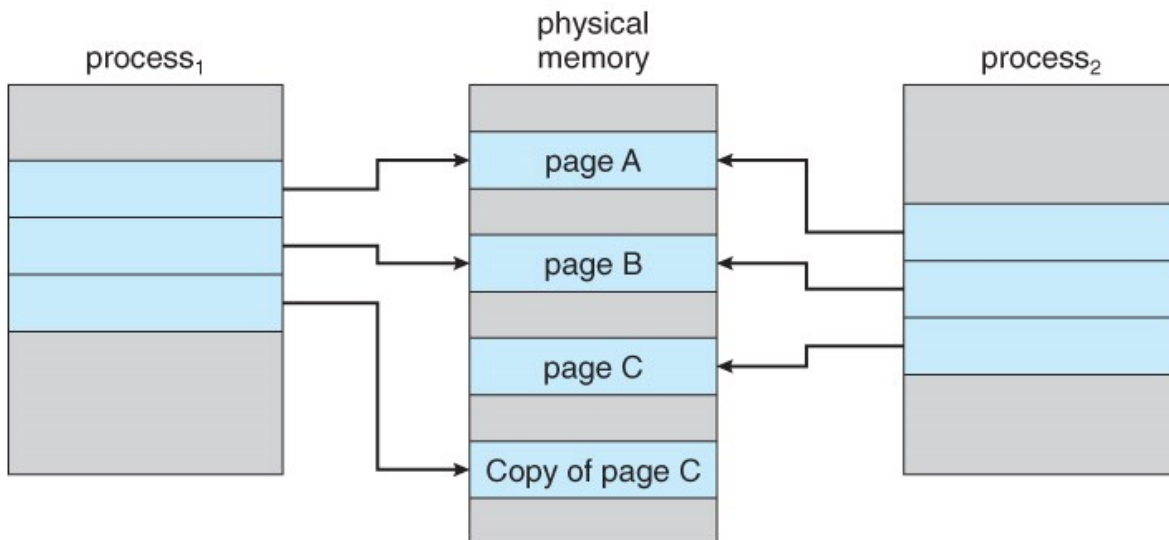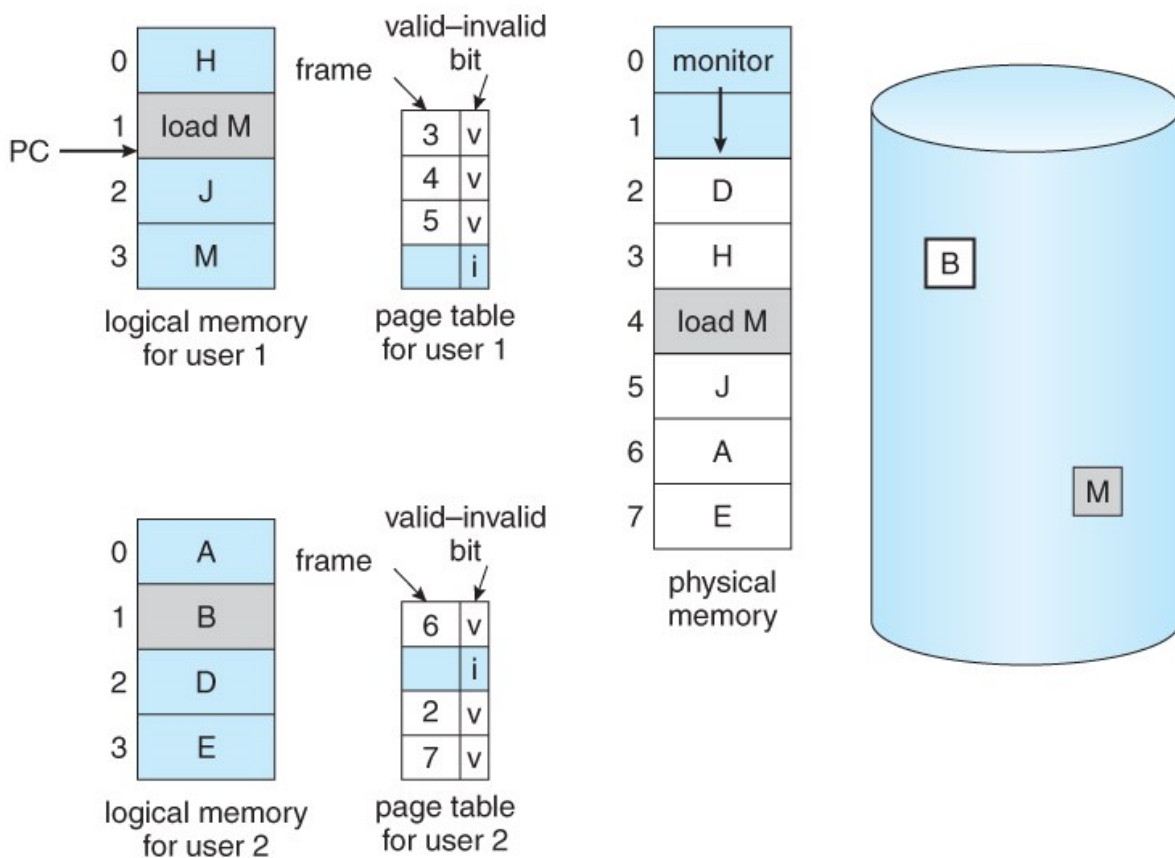**Figure 9.7 - Before process 1 modifies page C.**



**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.
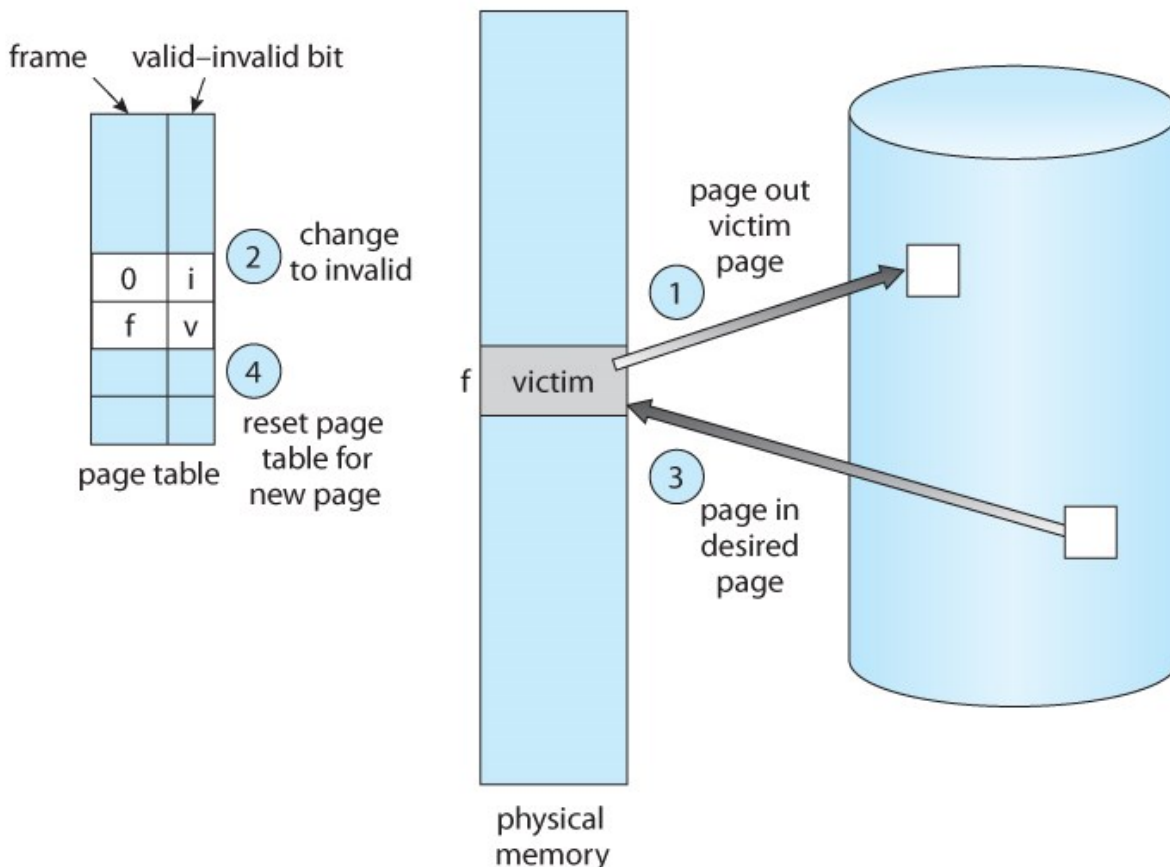
## 9.4 Page Replacement

- 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.**

### 9.4.1 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:
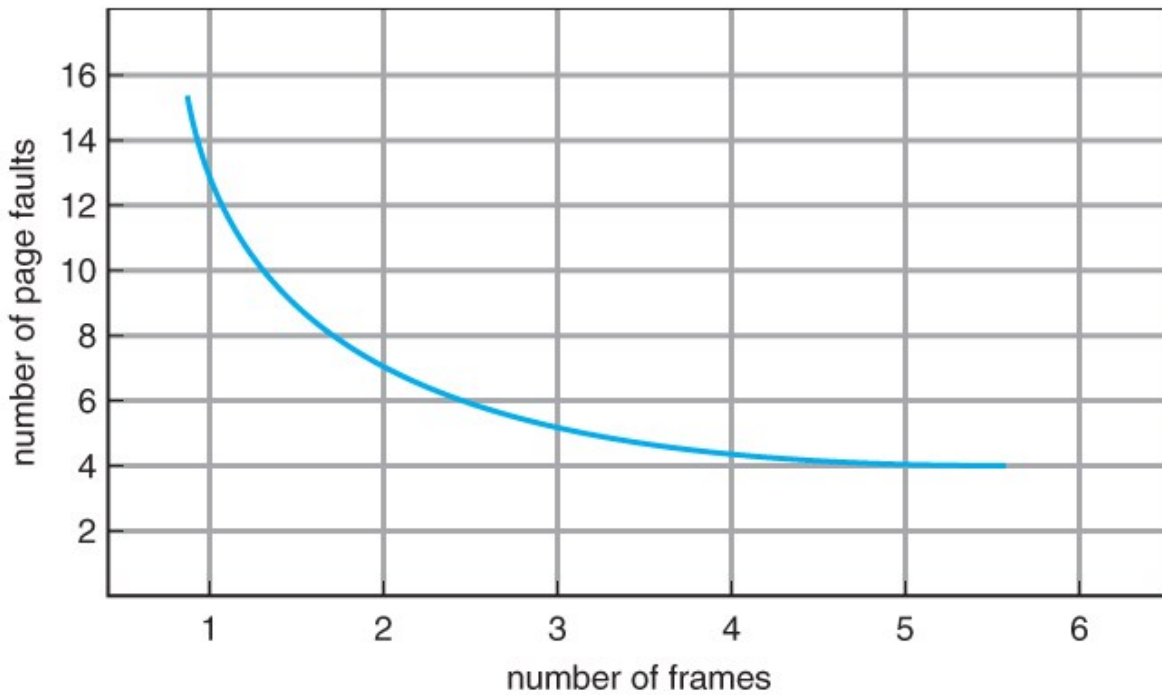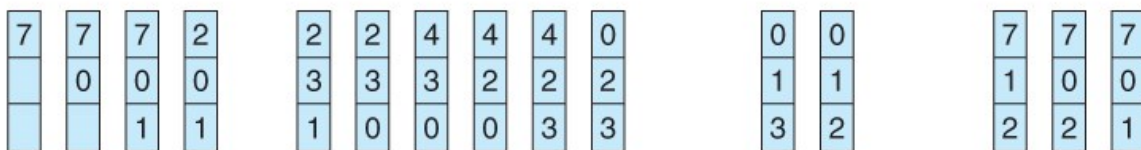
**Figure 9.11 - Graph of page faults versus number of frames.**

### 9.4.2 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:
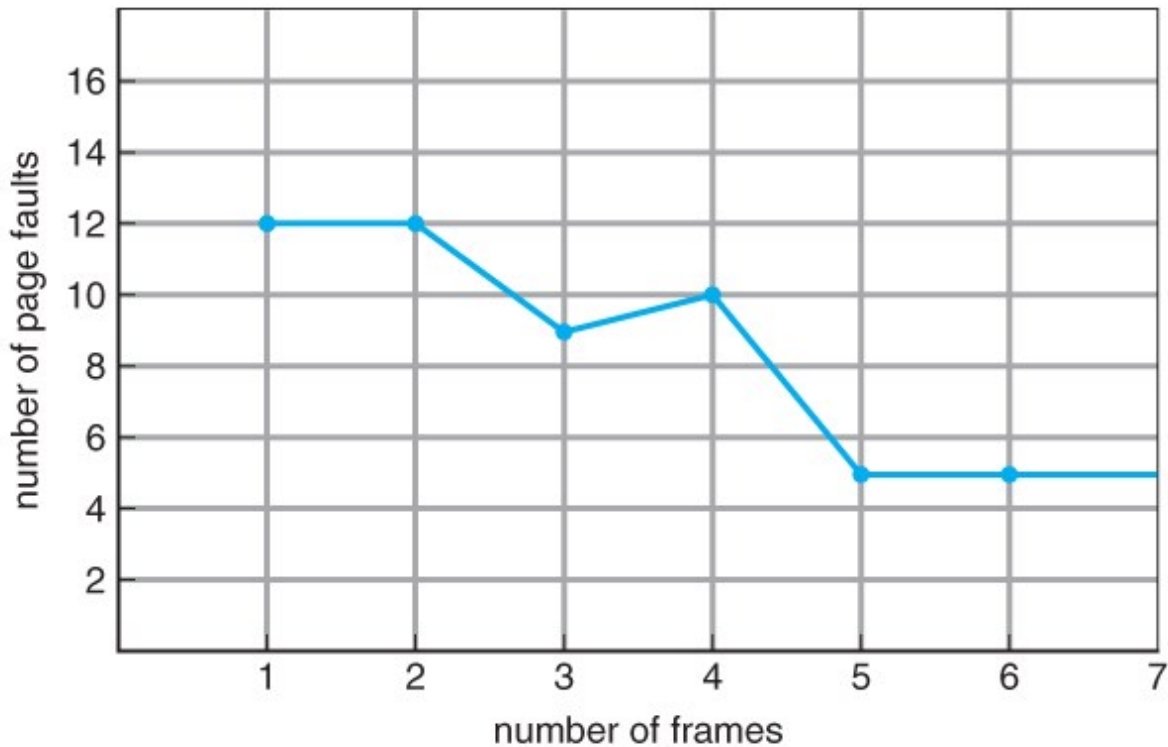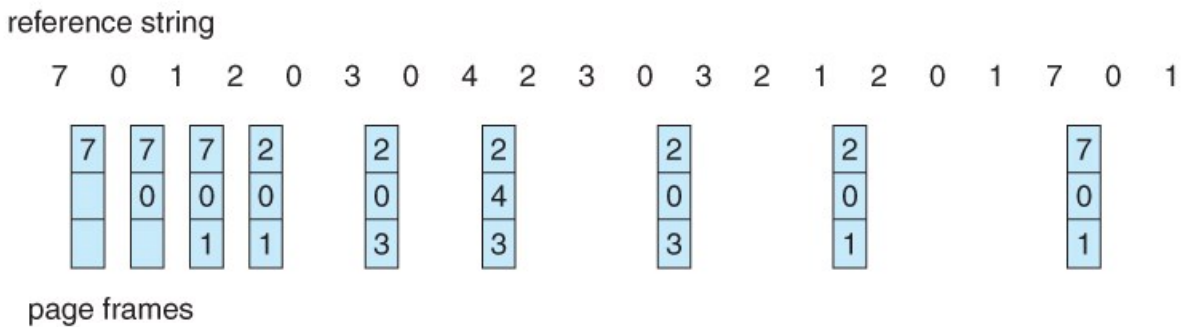


**Figure 9.12 - FIFO page-replacement algorithm.**

- 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:

137

**Figure 9.13 - Page-fault curve for FIFO replacement on a reference string.**
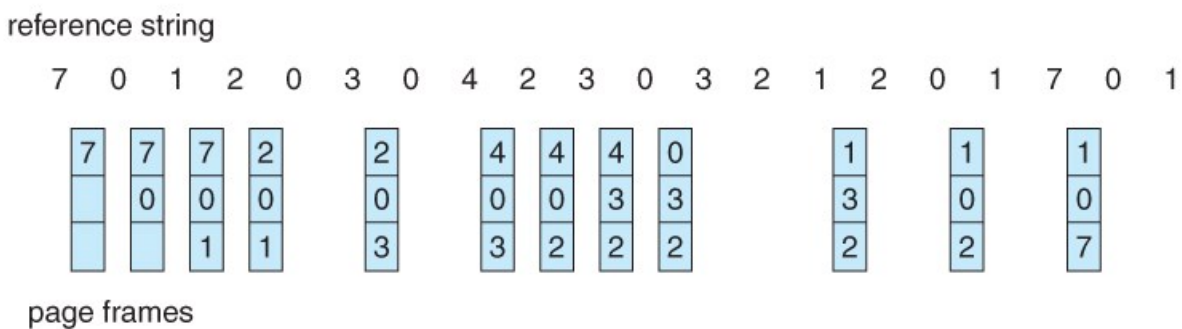
### 9.4.3 Optimal Page Replacement

- 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.
- 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.

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 9.14 - 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. )

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 9.15 - 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.
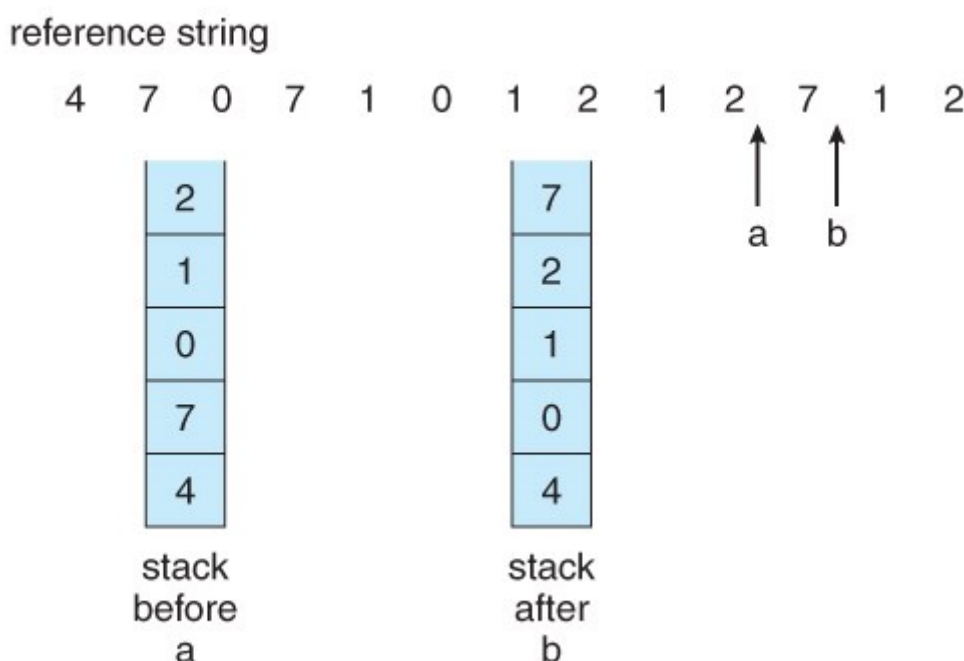
reference string

4   7   0   7   1   0   1   2   1   2   7   1   2

| 2 |
| 1 |
| 0 |
| 7 |
| 4 |

stack
before
a

| 7 |
| 2 |
| 1 |
| 0 |
| 4 |

stack
after
b

a   b

**Figure 9.16 - Use of a stack to record the most recent page references.**

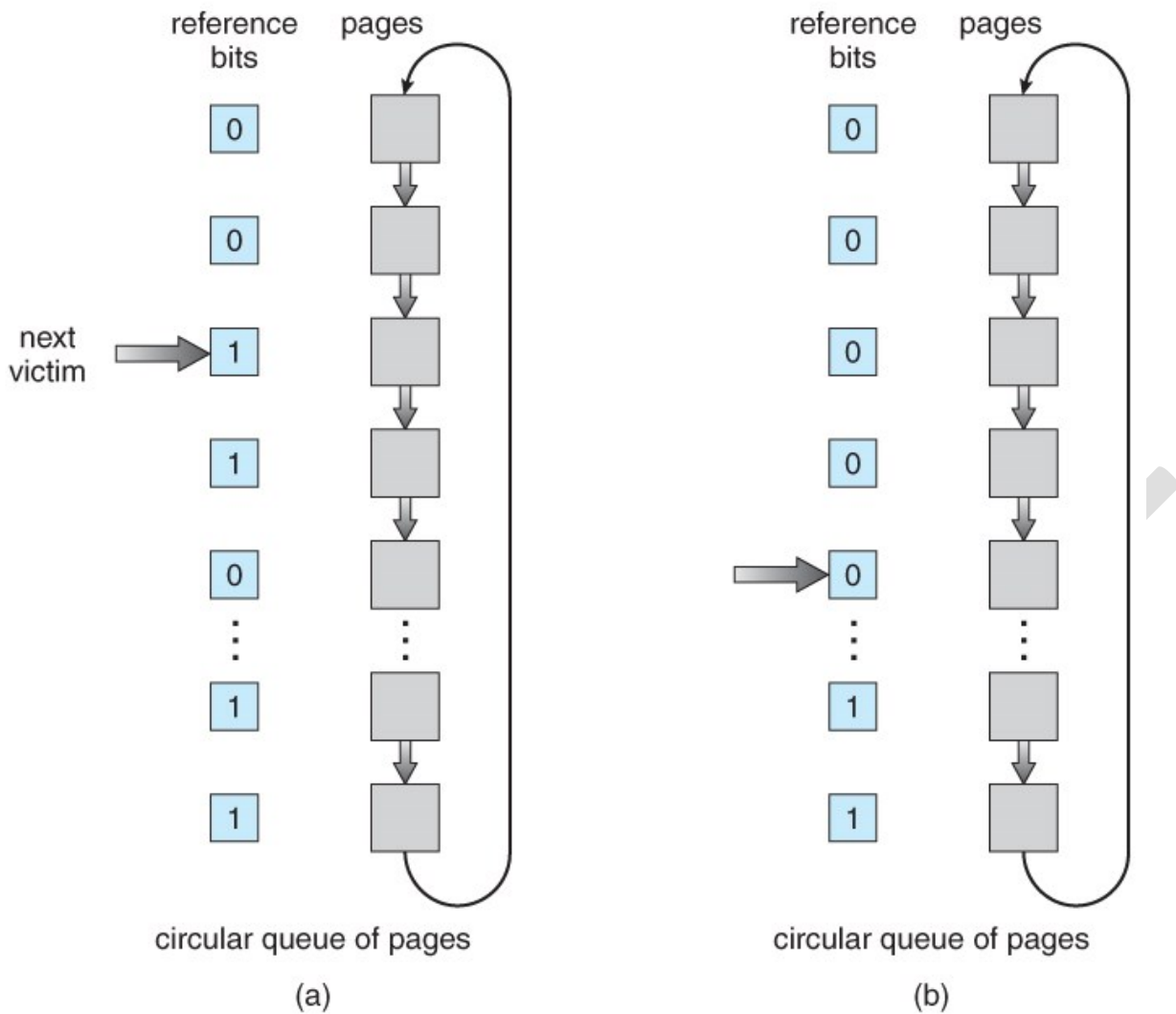**9.4.5 LRU-Approximation Page Replacement**

- 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.
- 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.**

**9.4.5.3 Enhanced Second-Chance Algorithm**

- The **enhanced second chance algorithm** looks at the reference bit and the modify bit ( dirty bit ) as an ordered page, and classifies pages into one of four classes:
  1. ( 0, 0 ) - Neither recently used nor modified.
  2. ( 0, 1 ) - Not recently used, but modified.
  3. ( 1, 0 ) - Recently used, but clean.
  4. ( 1, 1 ) - Recently used and modified.
- This algorithm searches the page table in a circular fashion ( in as many as four passes ), looking for the first page it can find in the lowest numbered category. I.e. it first makes a pass looking for a ( 0, 0 ), and then if it can't find one, it makes another pass looking for a ( 0, 1 ), etc.
- The main difference between this algorithm and the previous one is the preference for replacing clean pages if possible.

**9.4.6 Counting-Based Page Replacement**

- There are several algorithms based on counting the number of references that have been made to a given page, such as:
  - *Least Frequently Used, LFU:* Replace the page with the lowest reference count. A problem can occur if a page is used frequently initially and then not used any more, as the reference count remains high. A solution to this problem is to right-shift the counters periodically, yielding a time-decaying average reference count.
  - *Most Frequently Used, MFU:* Replace the page with the highest reference count. The logic behind this idea is that pages that have already been referenced a lot have been in the system a long time, and we are probably done with them, whereas pages referenced only a few times have only recently been loaded, and we still need them.
- In general counting-based algorithms are not commonly used, as their implementation is expensive and they do not approximate OPT well.

### 9.4.7 Page-Buffering Algorithms

There are a number of page-buffering algorithms that can be used in conjunction with the afore-mentioned algorithms, to improve overall performance and sometimes make up for inherent weaknesses in the hardware and/or the underlying page-replacement algorithms:

- Maintain a certain minimum number of free frames at all times. When a page-fault occurs, go ahead and allocate one of the free frames from the free list first, to get the requesting process up and running again as quickly as possible, and then select a victim page to write to disk and free up a frame as a second step.
- Keep a list of modified pages, and when the I/O system is otherwise idle, have it write these pages out to disk, and then clear the modify bits, thereby increasing the chance of finding a "clean" page for the next potential victim.
- Keep a pool of free frames, but remember what page was in it before it was made free. Since the data in the page is not actually cleared out when the page is freed, it can be made an active page again without having to load in any new data from disk. This is useful when an algorithm mistakenly replaces a page that in fact is needed again soon.

### 9.4.8 Applications and Page Replacement

- Some applications ( most notably database programs ) understand their data accessing and caching needs better than the general-purpose OS, and should therefore be given reign to do their own memory management.
- Sometimes such programs are given a *raw disk partition* to work with, containing raw data blocks and no file system structure. It is then up to the application to use this disk partition as extended memory or for whatever other reasons it sees fit.

### 9.5 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.

### 9.5.1 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 . . . 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.

### 9.5.2 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. )

### 9.5.3 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.
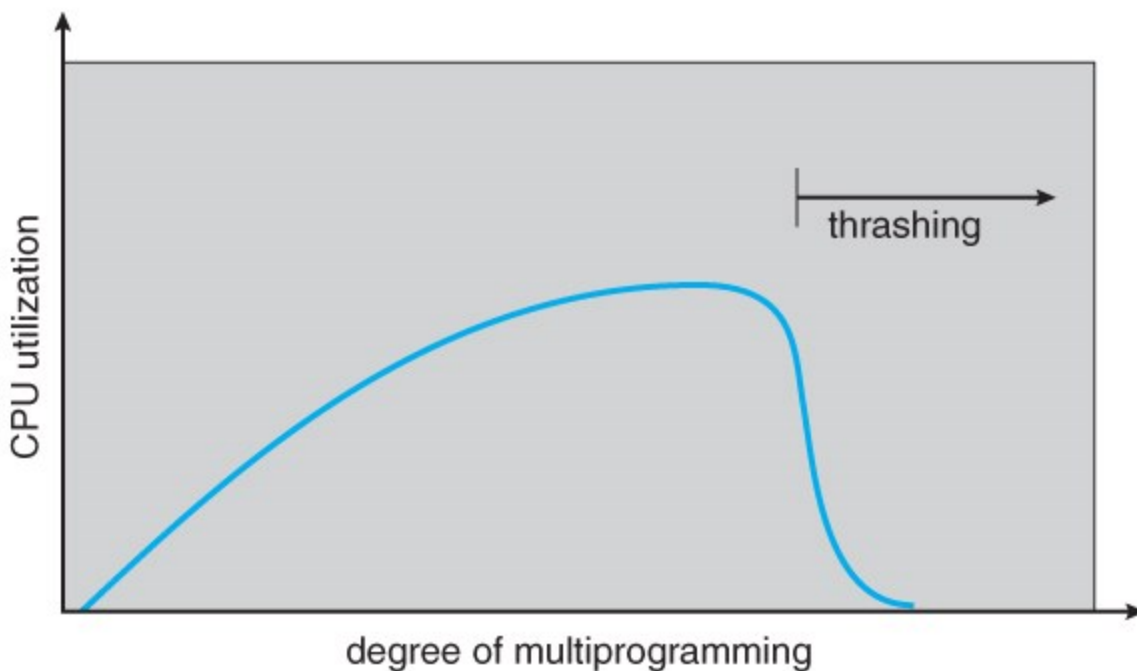
### 9.5.4 Non-Uniform Memory Access

- 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. )

## 9.6 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.
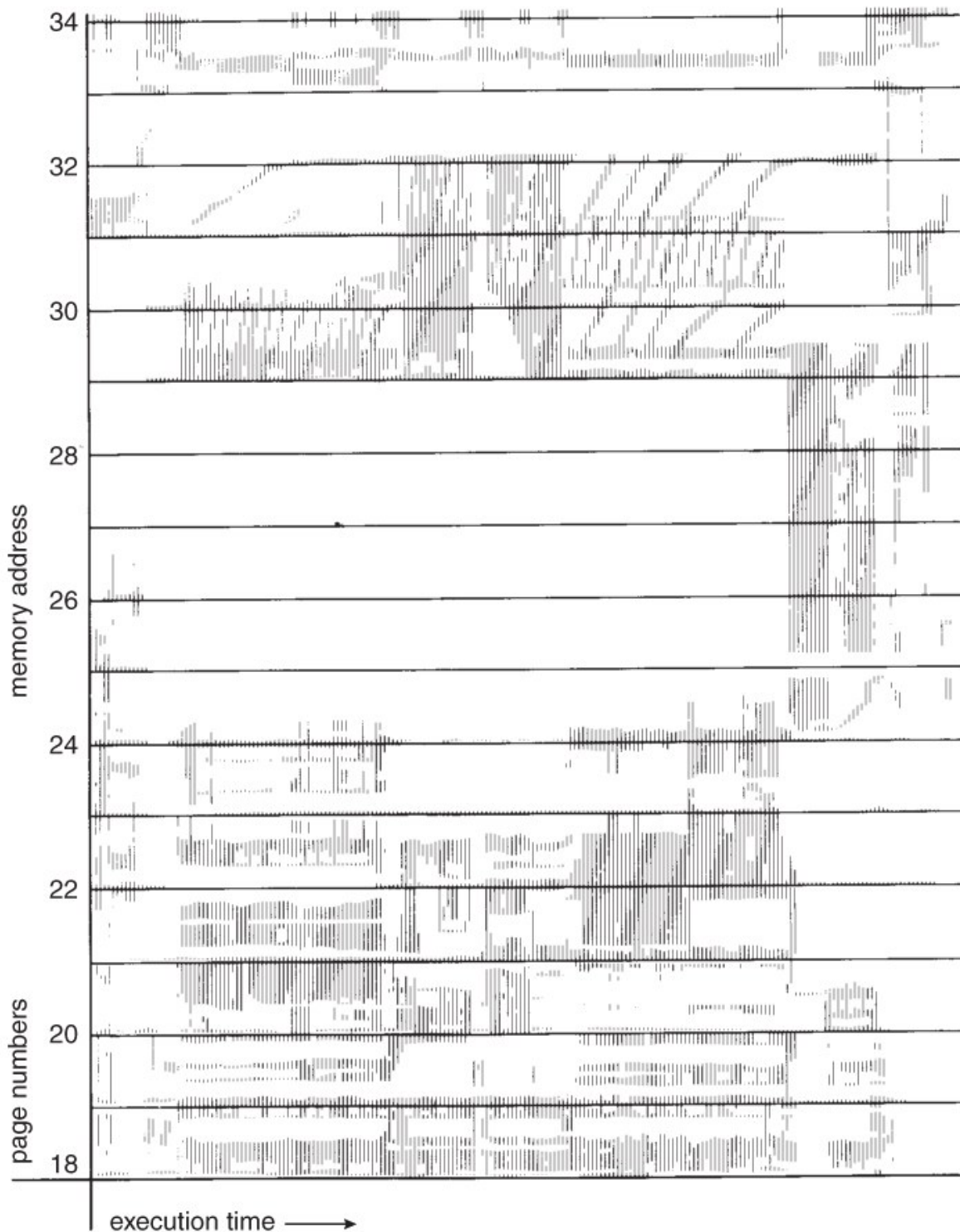- A process that is spending more time paging than executing is said to be *thrashing.*

### 9.6.1 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.
- 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. )
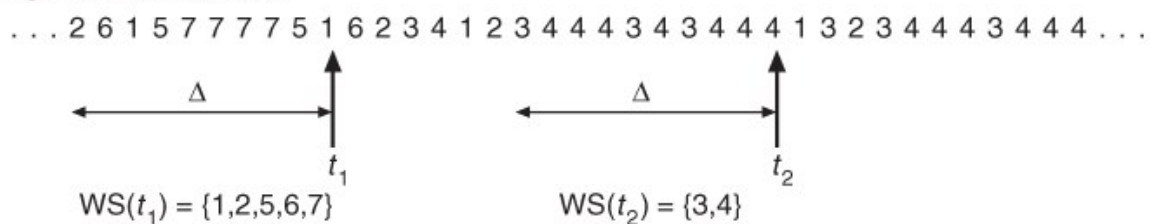
**Figure 9.19 - Locality in a memory-reference pattern.**

### 9.6.2 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

. . . 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 . . .

$\Delta$         $\Delta$

$t_1$         $t_2$

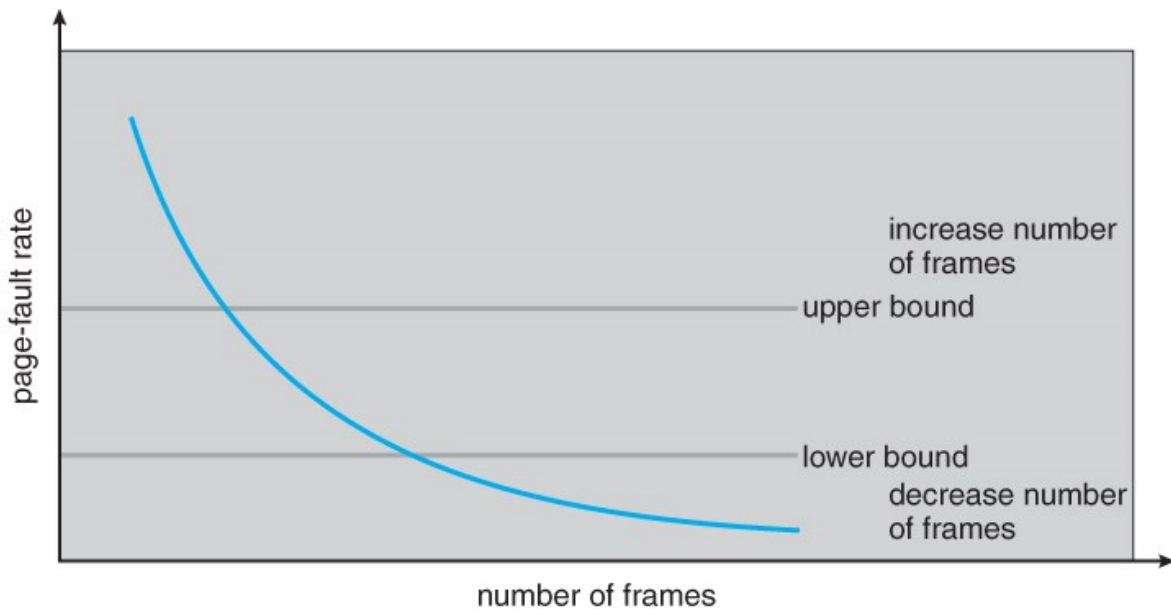$WS(t_1) = \{1,2,5,6,7\}$         $WS(t_2) = \{3,4\}$

**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.
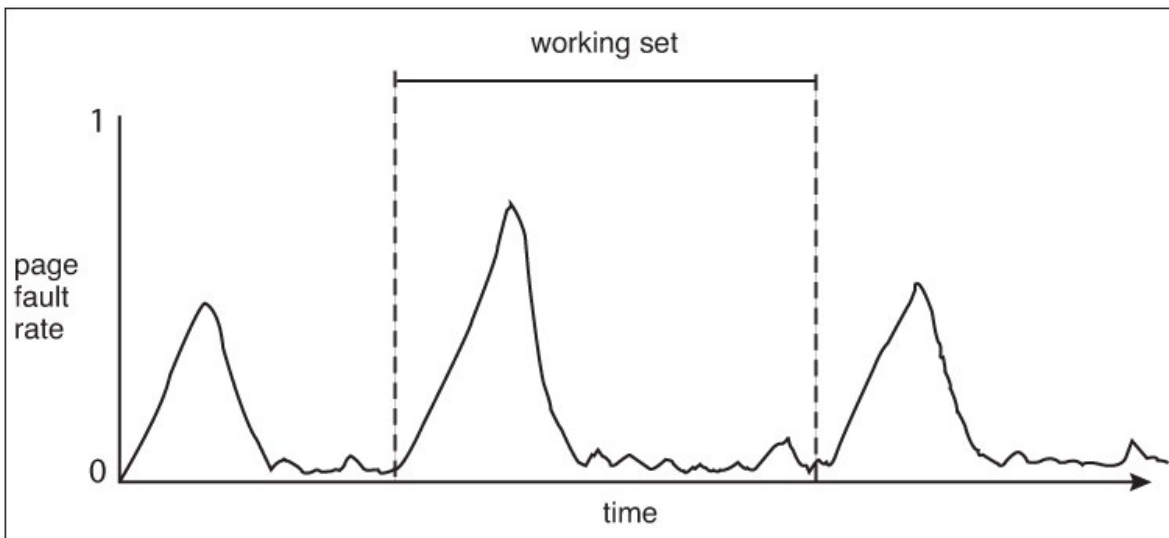
**9.6.3 Page-Fault Frequency**

- 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:
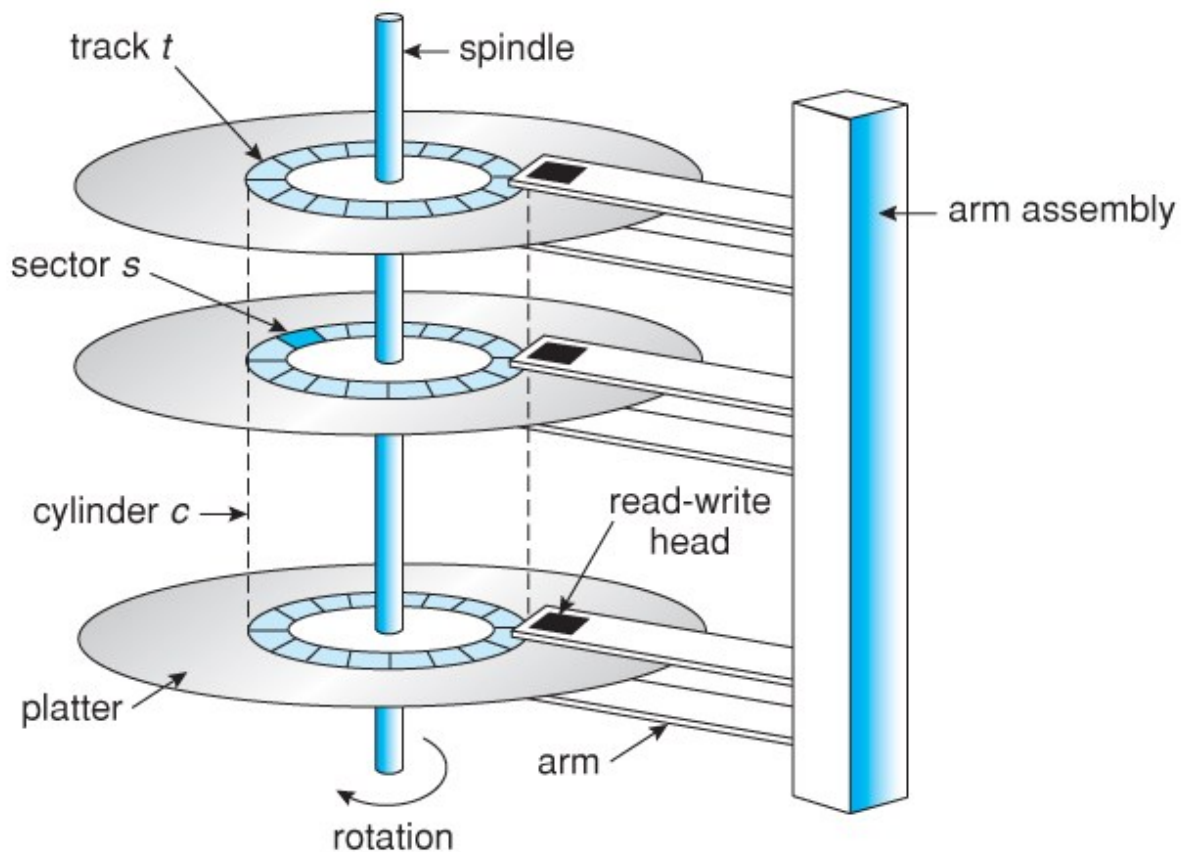


**Unnumbered side bar in Ninth Edition**

# UNIT IV

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:
  - 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.
  - 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. )
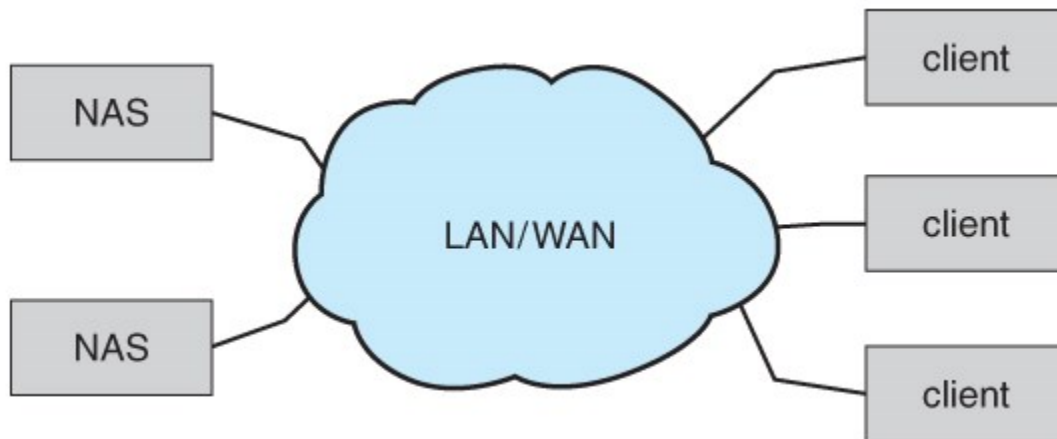
## 10.3 Disk Attachment

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:
    - 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.
    - 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. )
    - The SCSI standard also supports multiple host adapters in a single computer, i.e. multiple SCSI busses.
    - Modern advancements in SCSI include "fast" and "wide" versions, as well as SCSI-2.
    - SCSI cables may be either 50 or 68 conductors. SCSI devices may be external as well as internal.
    - 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:
    - 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.
    - The *arbitrated loop, FC-AL,* that can address up to 126 devices ( drives and controllers. )
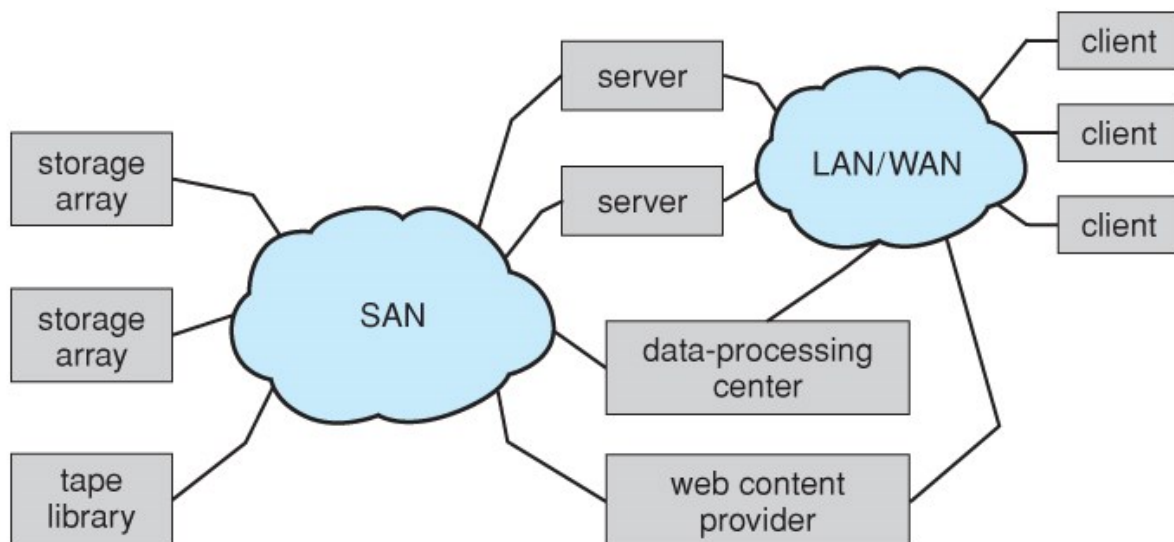
### 10.3.2 Network-Attached Storage

- 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:
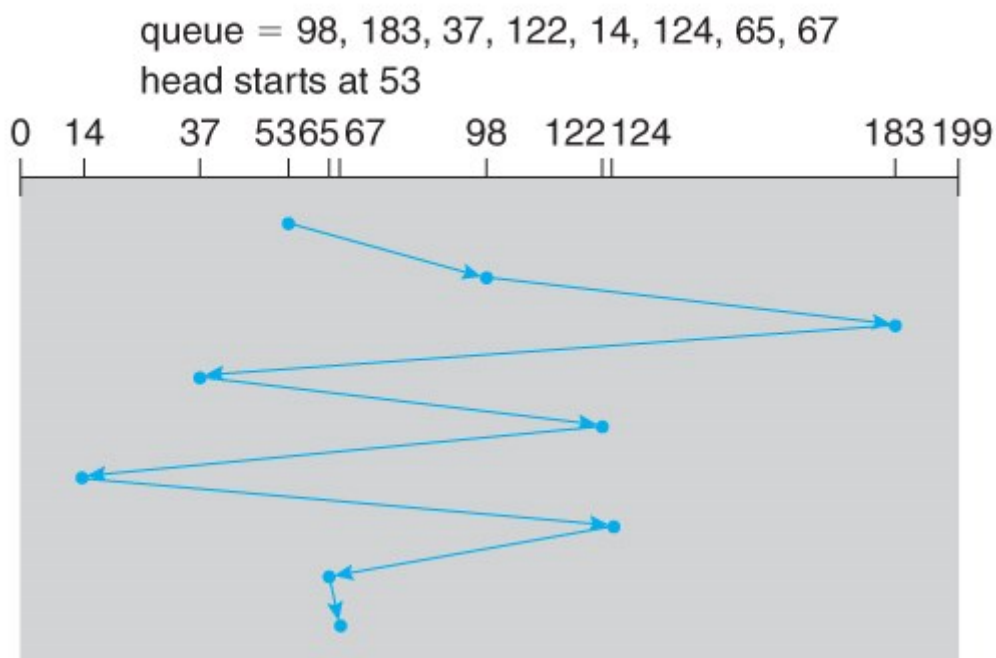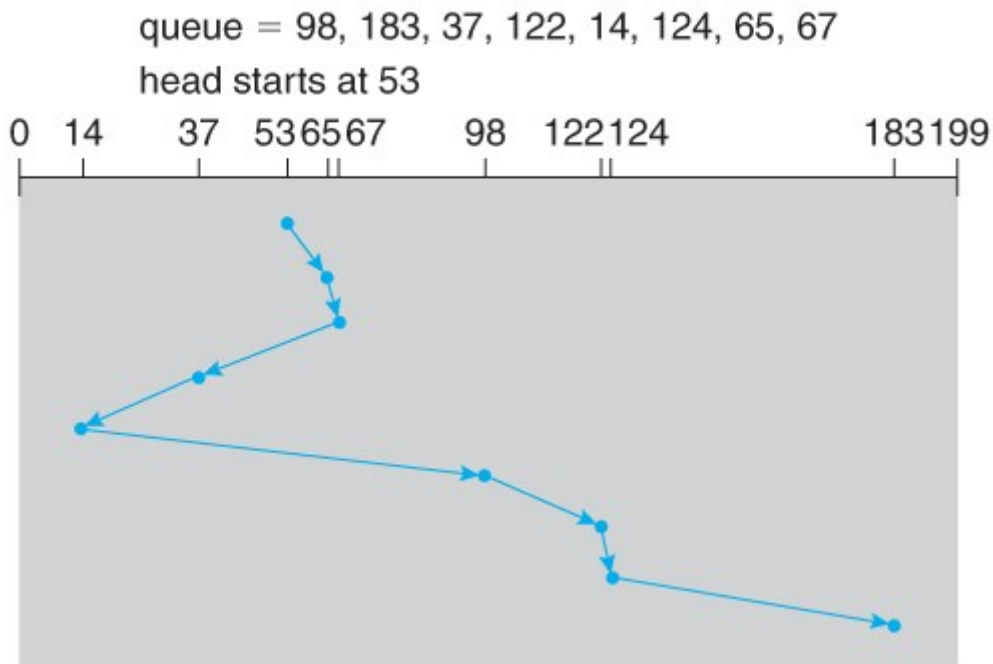
queue = 98, 183, 37, 122, 14, 124, 65, 67
head starts at 53

**Figure 10.4 - FCFS disk scheduling.**
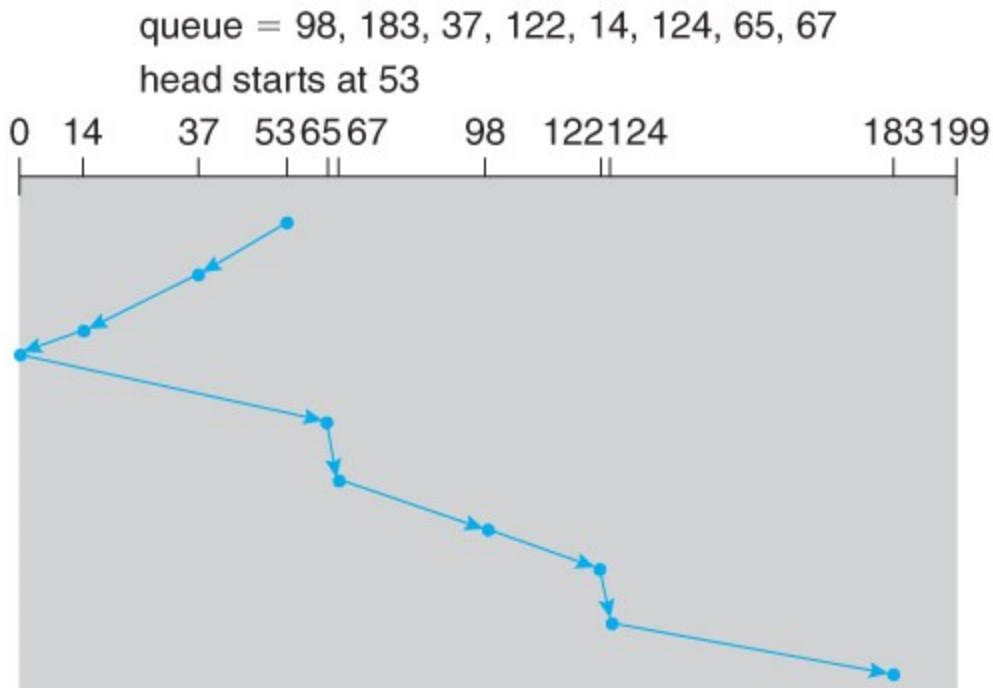
### 10.4.2 SSTF Scheduling

- **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.

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53

0   14      37   5365 67      98   122124                    183199



**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.

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53

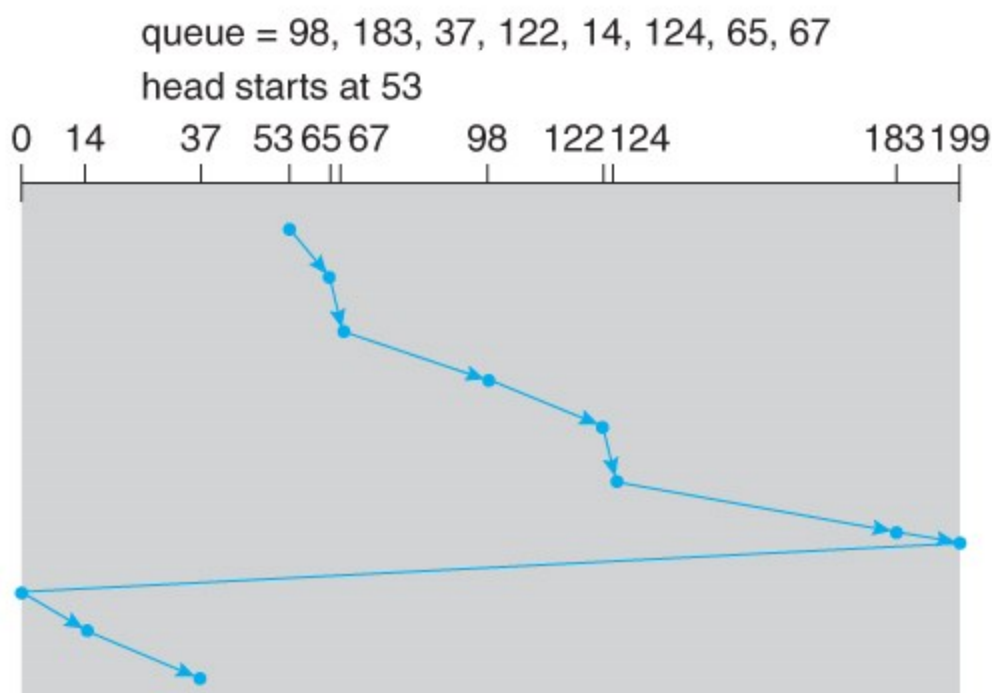0   14      37   5365 67      98   122124                    183199



**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.
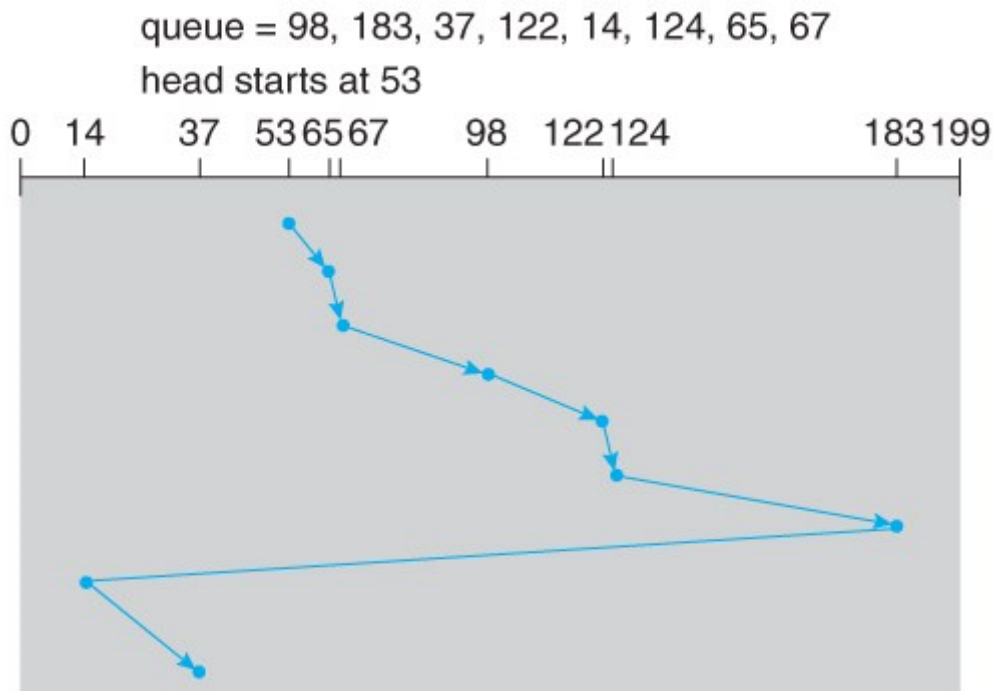
### 10.4.4 C-SCAN Scheduling

- 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.**

### 12.4.5 LOOK 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:

queue = 98, 183, 37, 122, 14, 124, 65, 67
head starts at 53

**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. )
  - o 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.

- o 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.
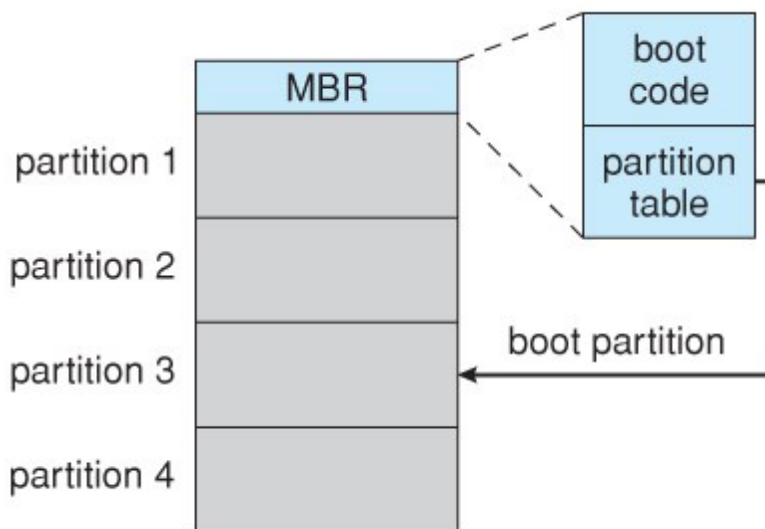
## 10.5 Disk Management

### 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

162

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. )
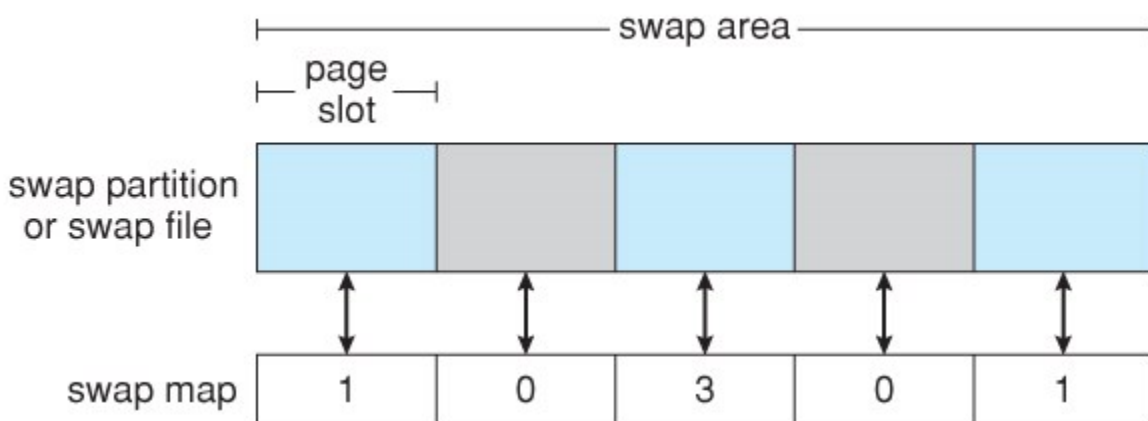


**Figure 10.10 - The data structures for swapping on Linux systems.**

## 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.
    - o 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.

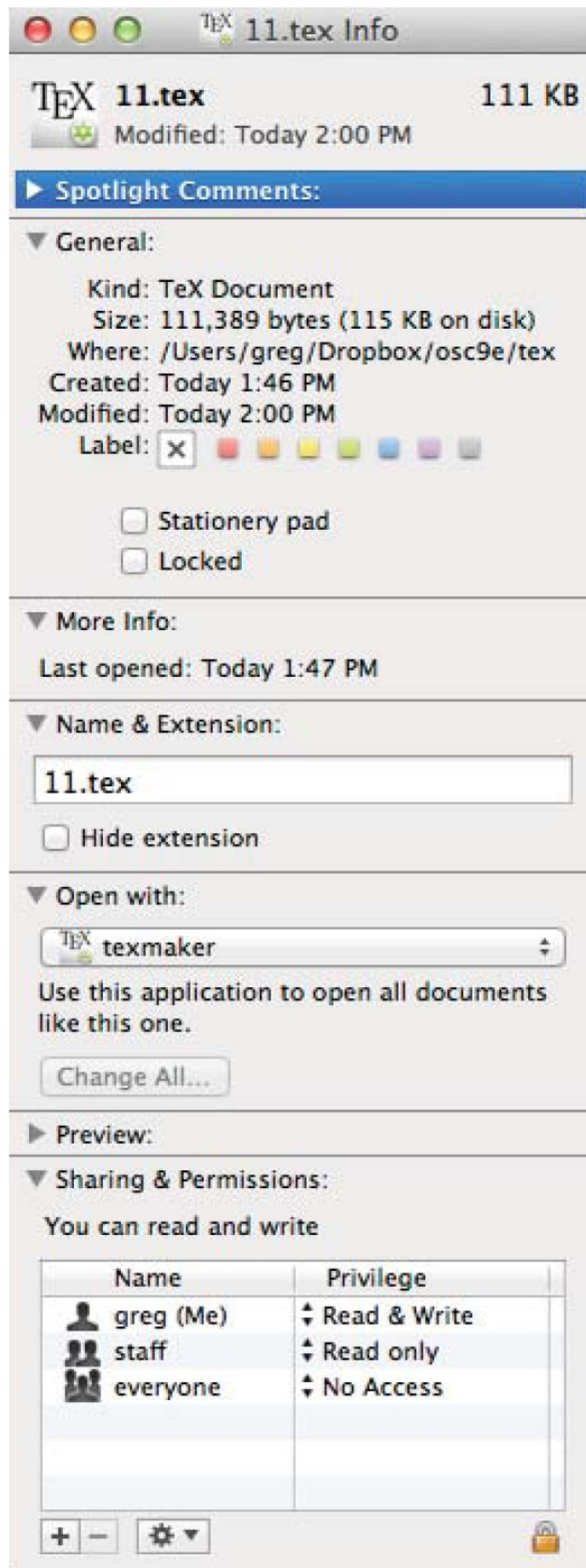### 10.7.2 Improvement in Performance via Parallelism

- 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.
    - o 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.
    - o *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.

## File Concept

### 11.1.1 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**

### 11.1.2 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 OSes 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 LOCKING IN JAVA (Cont.)

```java
import java.io.*;
import java.nio.channels.*;

public class LockingExample {
  public static final boolean EXCLUSIVE = false;
  public static final boolean SHARED = true;

  public static void main(String arsg[]) throws IOException {
    FileLock sharedLock = null;
    FileLock exclusiveLock = null;

    try {
      RandomAccessFile raf = new RandomAccessFile("file.txt","rw");

      // get the channel for the file
      FileChannel ch = raf.getChannel();

      // this locks the first half of the file - exclusive
      exclusiveLock = ch.lock(0, raf.length()/2, EXCLUSIVE);

      /** Now modify the data . . . */

      // release the lock
      exclusiveLock.release();

      // this locks the second half of the file - shared
      sharedLock = ch.lock(raf.length()/2+1,raf.length(),SHARED);

      /** Now read the data . . . */

      // release the lock
      exclusiveLock.release();
    } catch (java.io.IOException ioe) {
      System.err.println(ioe);
    }
    finally {
      if (exclusiveLock != null)
          exclusiveLock.release();
      if (sharedLock != null)
          sharedLock.release();
    }
  }
}
```

**Figure 11.2 - File-locking example in Java.**

### 11.1.3 File Types

- Windows ( and some other systems ) use special file extensions to indicate the type of each file:

| 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 |

**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 )

### 11.1.4 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.

**11.1.5 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.

## 11.2 Access Methods

### 11.2.1 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.
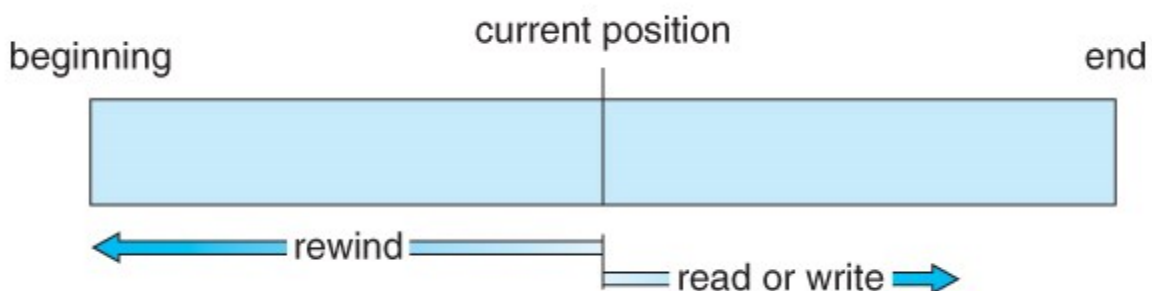


**Figure 11.4 - Sequential-access file.**
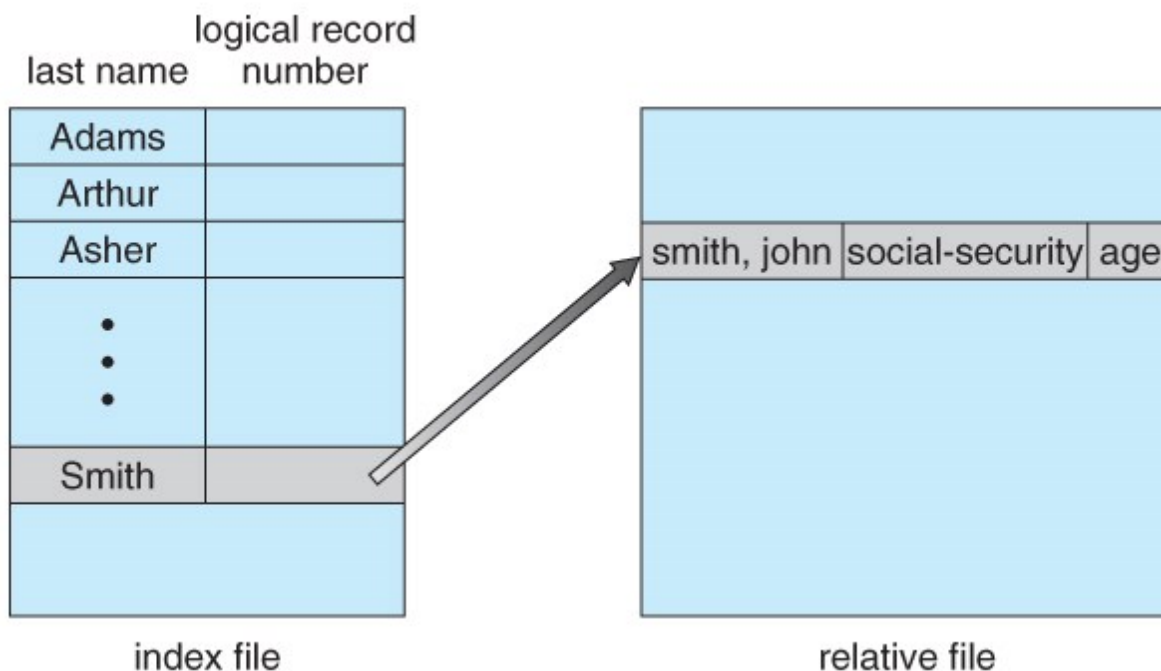
### 11.2.2 Direct Access

- 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.**

### 11.2.3 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.
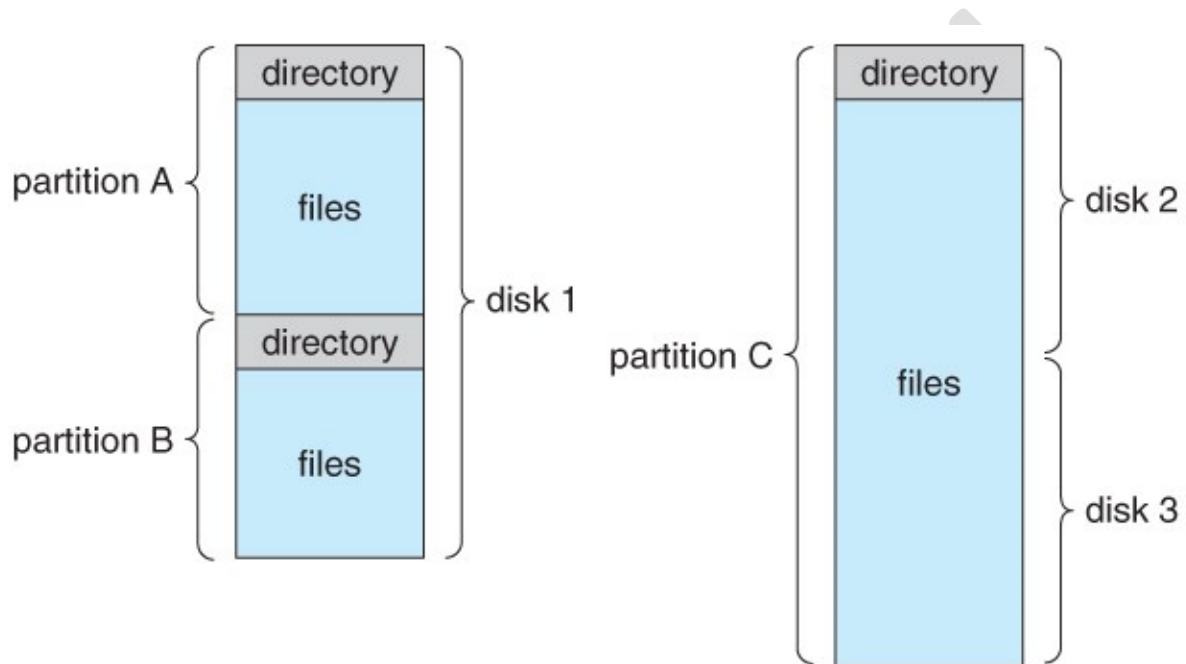


**Figure 11.6 - Example of index and relative files.**

## 11.3 Directory Structure

### 11.3.1 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.



**Figure 11.7 - A typical file-system organization.**

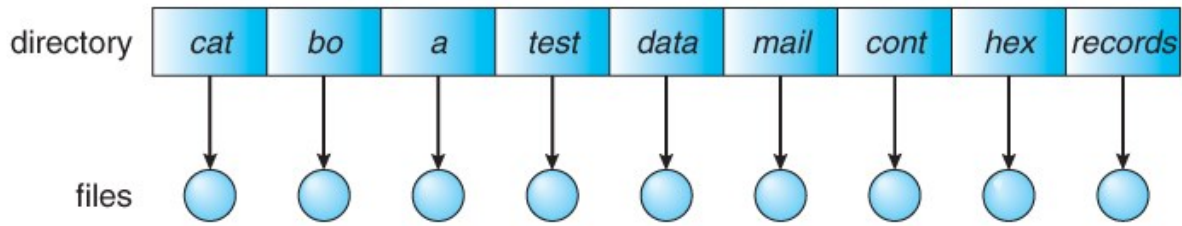| | |
|---|---|
| / | ufs |
| /devices | devfs |
| /dev | dev |
| /system/contract | ctfs |
| /proc | proc |
| /etc/mnttab | mntfs |
| /etc/svc/volatile | tmpfs |
| /system/object | objfs |
| /lib/libc.so.1 | lofs |
| /dev/fd | fd |
| /var | ufs |
| /tmp | tmpfs |
| /var/run | tmpfs |
| /opt | ufs |
| /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.

### 11.3.2 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.

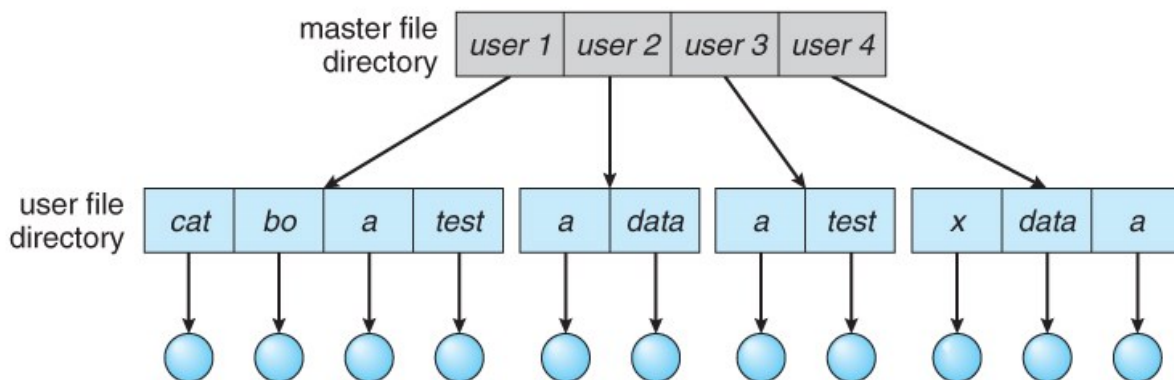### 11.3.3. Single-Level Directory

- Simple to implement, but each file must have a unique name.

**Figure 11.9 - Single-level directory.**

### 11.3.4 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.**

### 11.3.5 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.

174

- 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.
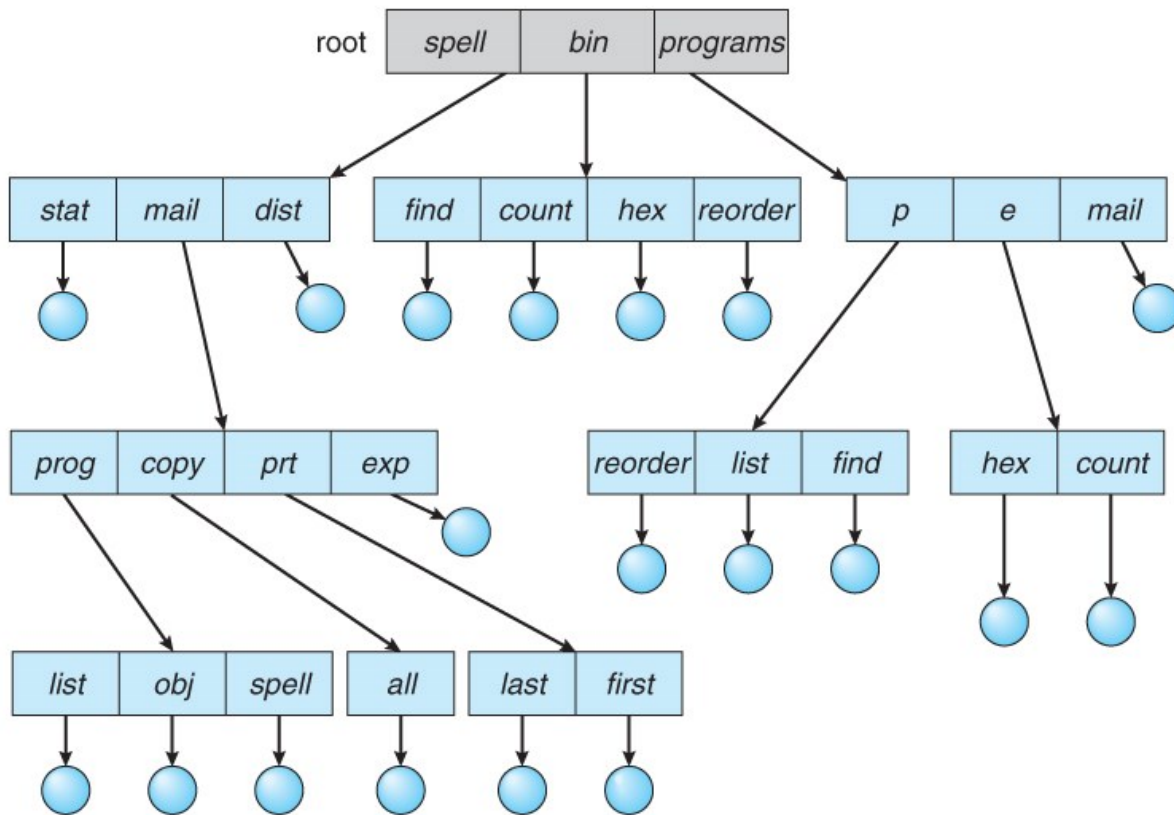


**Figure 11.11 - Tree-structured directory structure.**

### 11.3.6 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:
    - One option is to find all the symbolic links and adjust them also.
    - Another is to leave the symbolic links dangling, and discover that they are no longer valid the next time they are used.
    - What if the original file is removed, and replaced with another file having the same name before the symbolic link is next used?
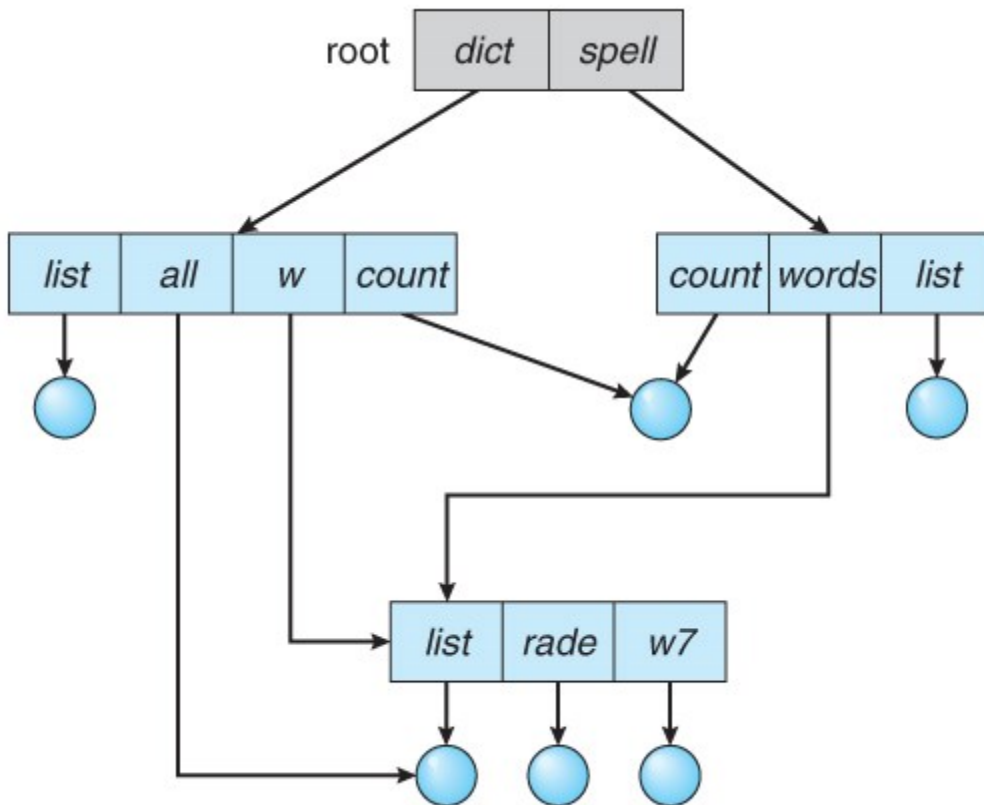


**Figure 11.12 - Acyclic-graph directory structure.**

### 11.3.7 General Graph Directory

- If cycles are allowed in the graphs, then several problems can arise:
    - 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. )
    - 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. )
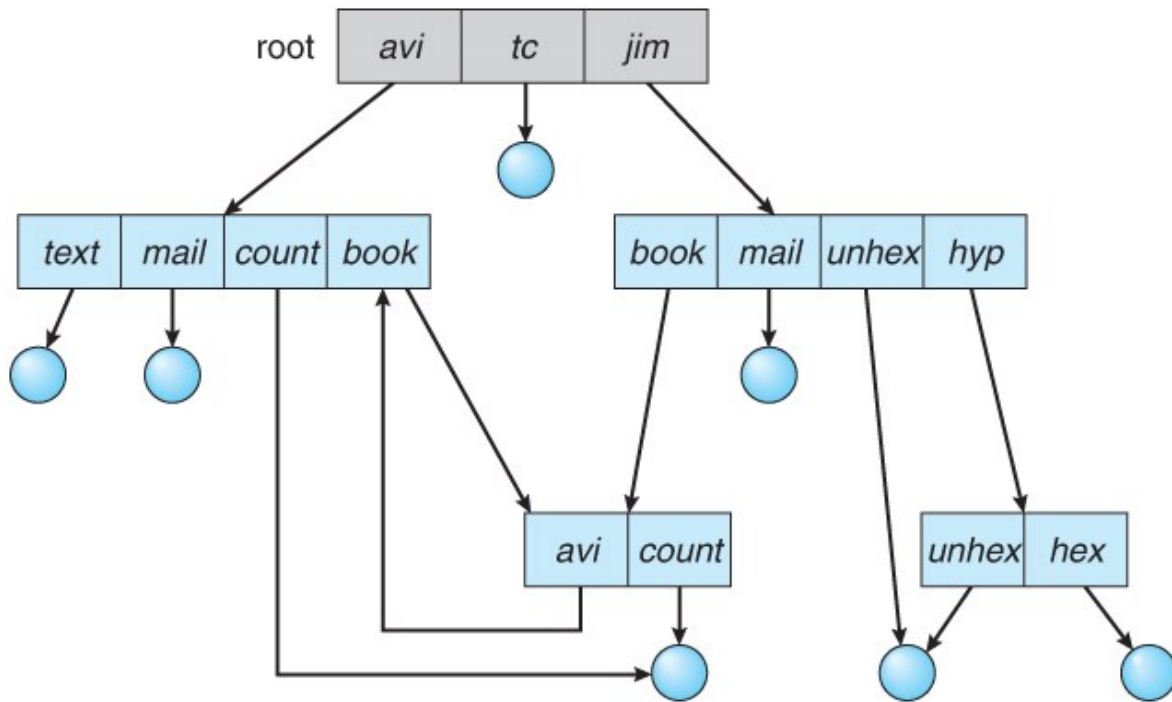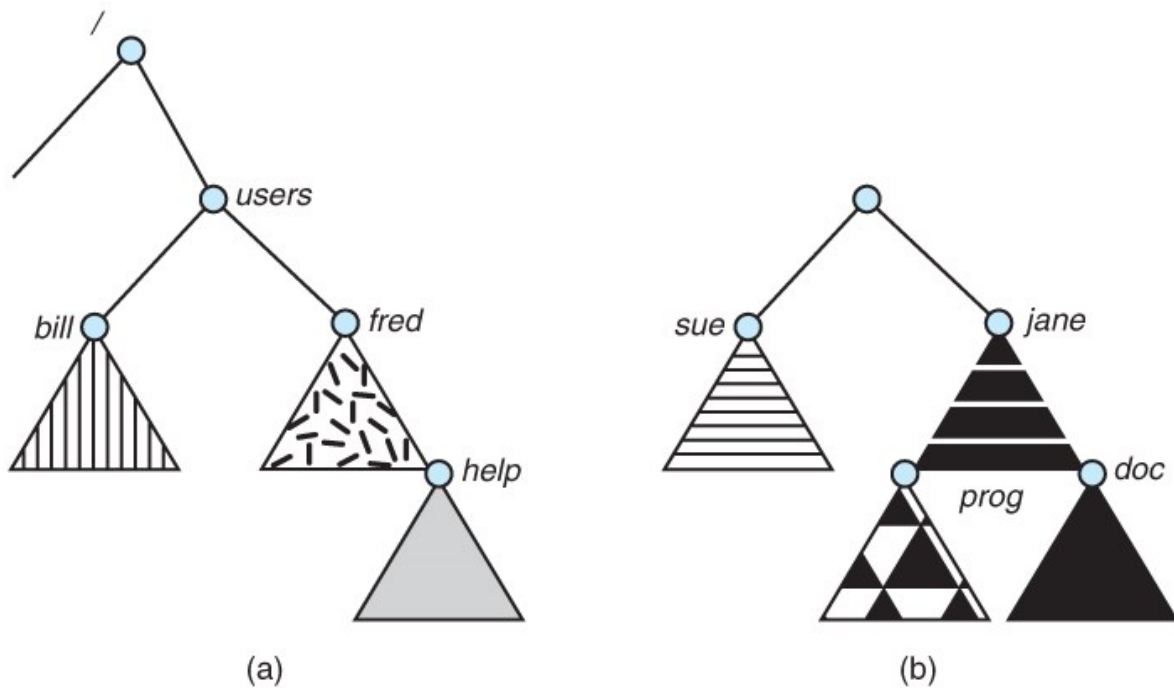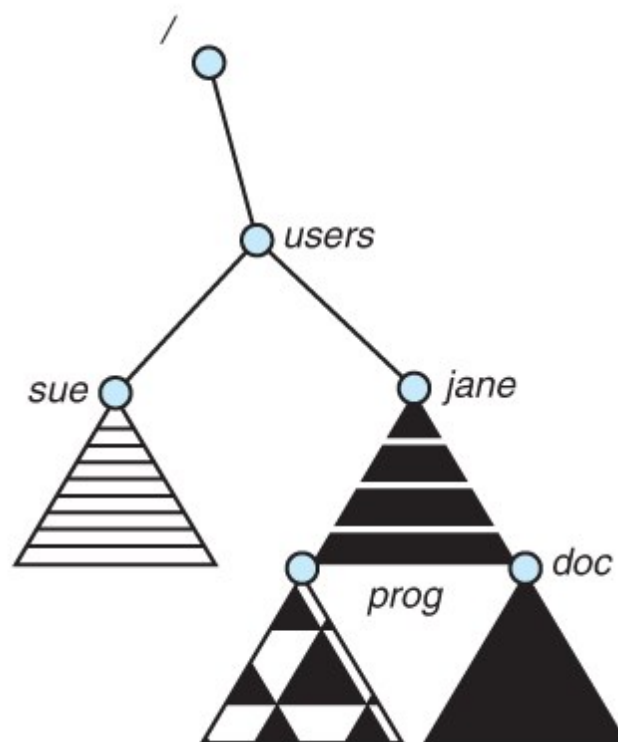
**Figure 11.13 - General graph directory.**

## 11.4 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.

## 11.5 File Sharing

### 11.5.1 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.

### 11.5.2 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. ( The actual files are still transported across the network as needed, possibly using ftp as the underlying transport mechanism. )
  - 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.

### 11.5.2.1 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:

- Servers commonly restrict mount permission to certain trusted systems only. Spoofing ( a computer pretending to be a different computer ) is a potential security risk.
- Servers may restrict remote access to read-only.
- 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.

### 11.5.2.2 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.
- 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.

### 11.5.2.3 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.

### 11.5.3 Consistency Semantics

- *Consistency Semantics* deals with the consistency between the views of shared files on a networked system. When one user changes the file, when do other users see the changes?
- At first glance this appears to have all of the synchronization issues discussed in Chapter 6. Unfortunately the long delays involved in network

operations prohibit the use of atomic operations as discussed in that chapter.

### 11.5.3.1 UNIX Semantics

- The UNIX file system uses the following semantics:
  - o Writes to an open file are immediately visible to any other user who has the file open.
  - o One implementation uses a shared location pointer, which is adjusted for all sharing users.
- The file is associated with a single exclusive physical resource, which may delay some accesses.

### 11.5.3.2 Session Semantics

- The Andrew File System, AFS uses the following semantics:
  - o Writes to an open file are not immediately visible to other users.
  - o When a file is closed, any changes made become available only to users who open the file at a later time.
- According to these semantics, a file can be associated with multiple ( possibly different ) views. Almost no constraints are imposed on scheduling accesses. No user is delayed in reading or writing their personal copy of the file.
- AFS file systems may be accessible by systems around the world. Access control is maintained through ( somewhat ) complicated access control lists, which may grant access to the entire world ( literally ) or to specifically named users accessing the files from specifically named remote environments.

### 11.5.3.3 Immutable-Shared-Files Semantics

- Under this system, when a file is declared as *shared* by its creator, it becomes immutable and the name cannot be re-used for any other resource. Hence it becomes read-only, and shared access is simple.

## 11.6 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.

### 11.6.1 Types of Access

- The following low-level operations are often controlled:

- Read - View the contents of the file
- Write - Change the contents of the file.
- Execute - Load the file onto the CPU and follow the instructions contained therein.
- Append - Add to the end of an existing file.
- Delete - Remove a file from the system.
- 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.

## 11.6.2 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
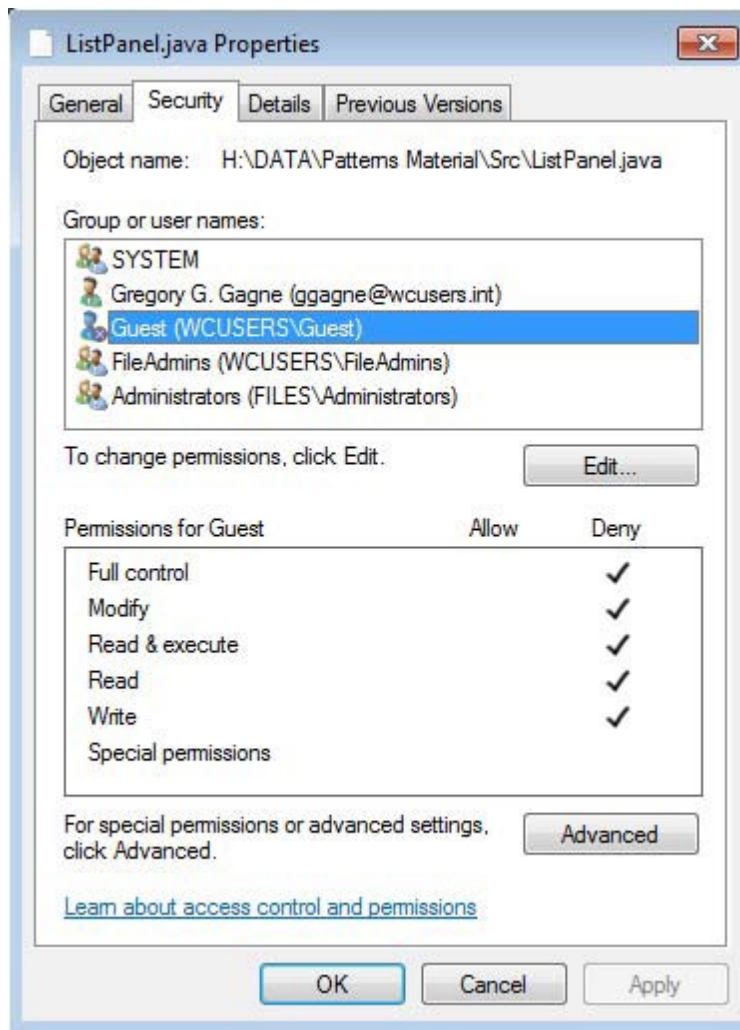
usually restricted to root, and must be done with caution, as it introduces a potential security leak.

- o 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.
- o 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:

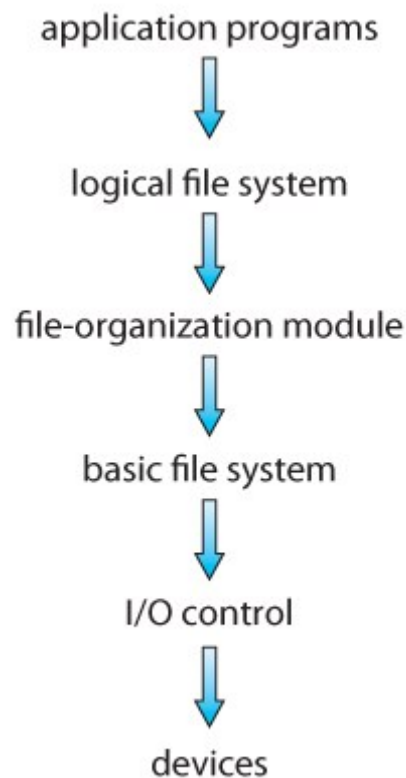**Figure 11.16 - Windows 7 access-control list management.**

### 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.

## 12.1 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.**
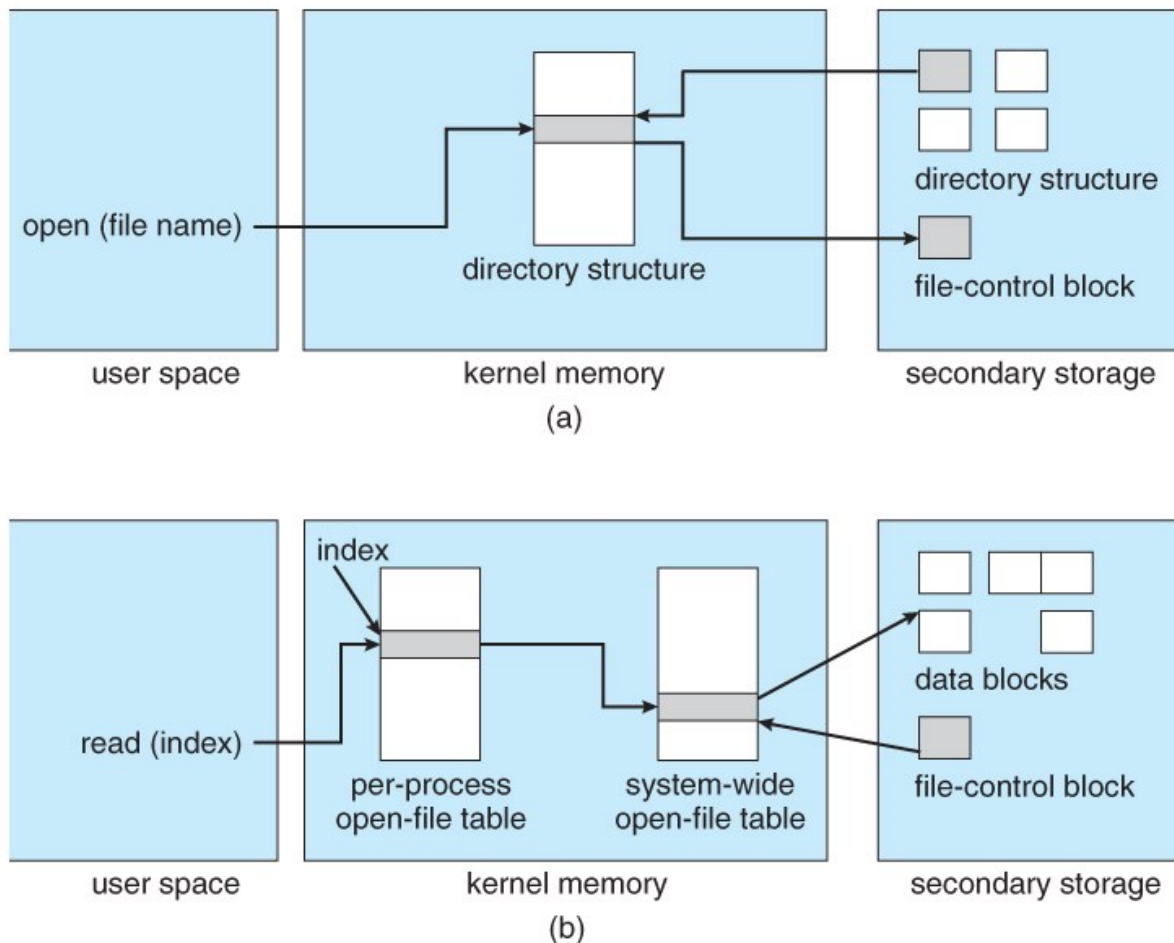
## 12.2 File-System Implementation

### 12.2.1 Overview

- File systems store several important data structures on the disk:
  - 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.
  - 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.
  - 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.**
  - 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:
  - An in-memory mount table.
  - An in-memory directory cache of recently accessed directory information.
  - *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.
  - *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:
  - 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.
  - 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*.
  - 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.
  - 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.

187

**Figure 12.3 - In-memory file-system structures. (a) File open. (b) File read.**
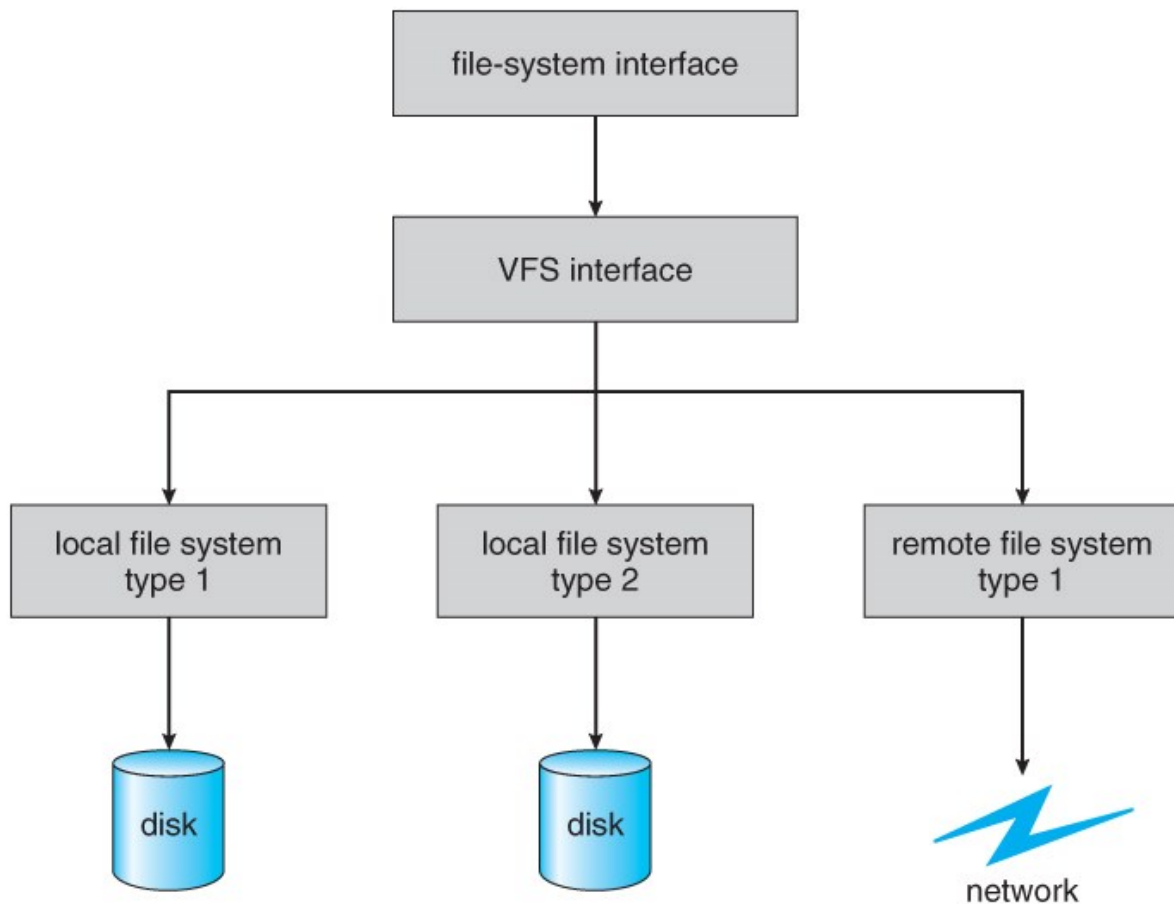
### 12.2.2 Partitions and Mounting

- 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 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.

### 12.2.3 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.**

## 12.3 Directory Implementation

- Directories need to be fast to search, insert, and delete, with a minimum of wasted disk space.

### 12.3.1 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.
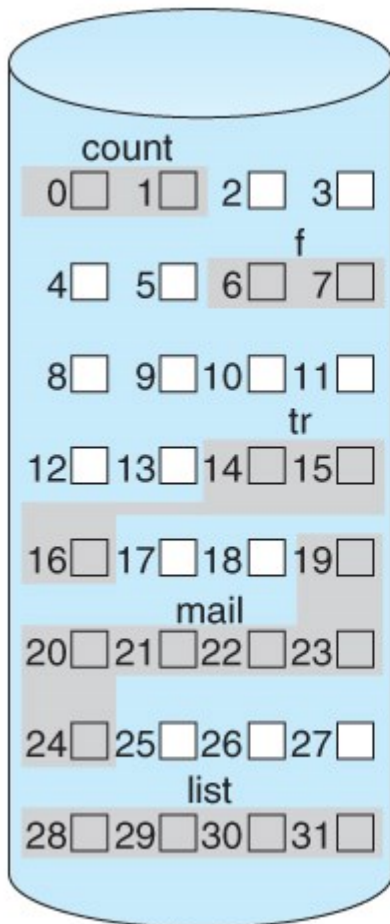
### 12.3.2 Hash Table

- A hash table can also be used to speed up searches.
- Hash tables are generally implemented *in addition to* a linear or other structure

## 12.4 Allocation Methods

- There are three major methods of storing files on disks: contiguous, linked, and indexed.

### 12.4.1 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:

  - Over-estimation of the file's final size increases external fragmentation and wastes disk space.

  - Under-estimation may require that a file be moved or a process aborted if the file grows beyond its originally allocated space.

  - 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.

directory

| file | start | length |
|------|-------|--------|
| count | 0 | 2 |
| tr | 14 | 3 |
| mail | 19 | 6 |
| list | 28 | 4 |
| f | 6 | 2 |

**Figure 12.5 - Contiguous allocation of disk space.**

**12.4.2 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.
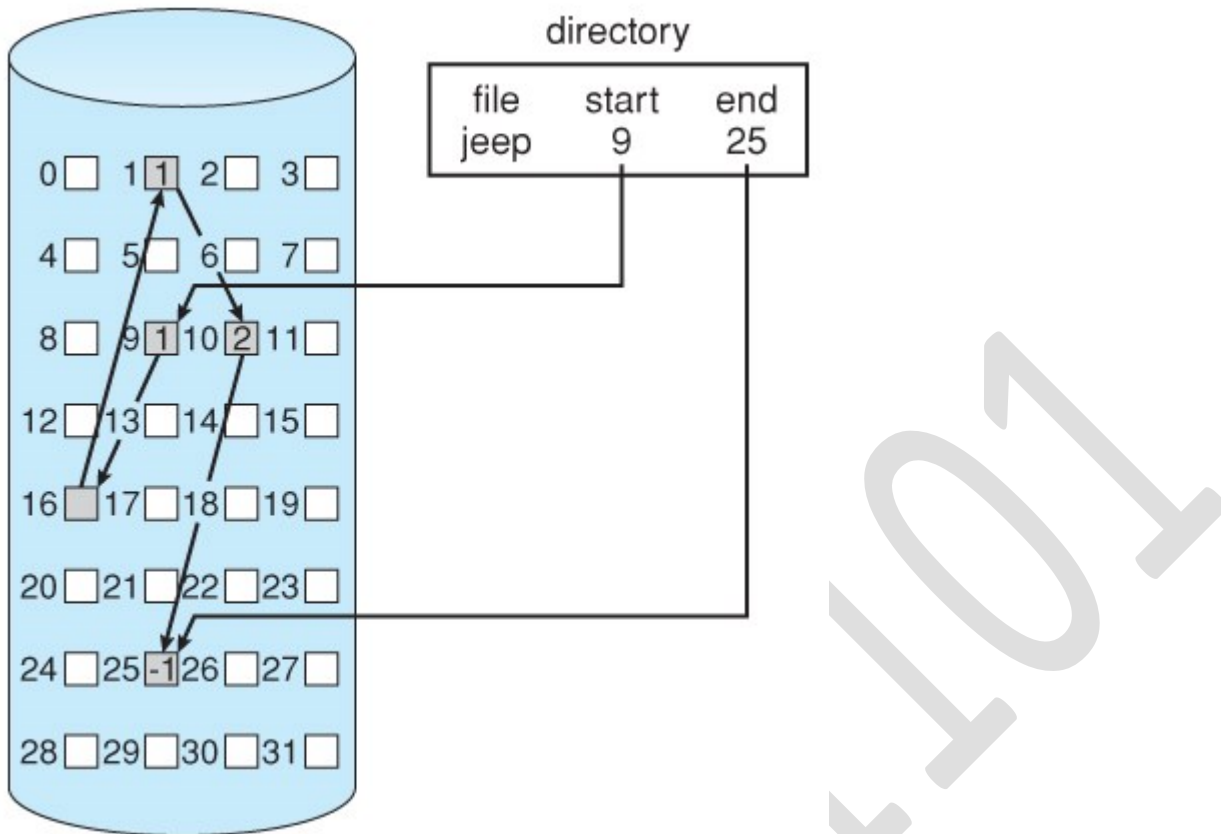
**Figure 12.6 - Linked allocation of disk space.**

- 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.
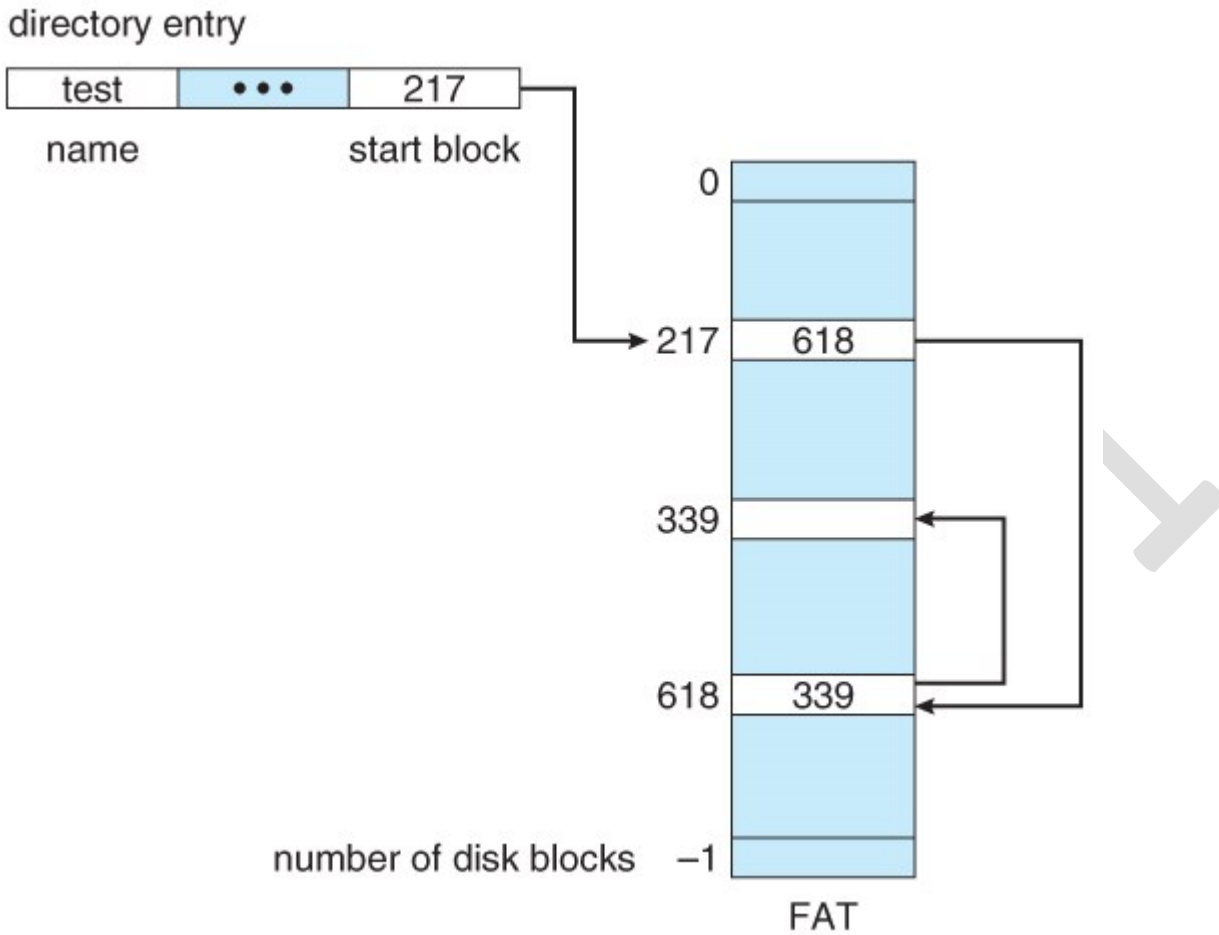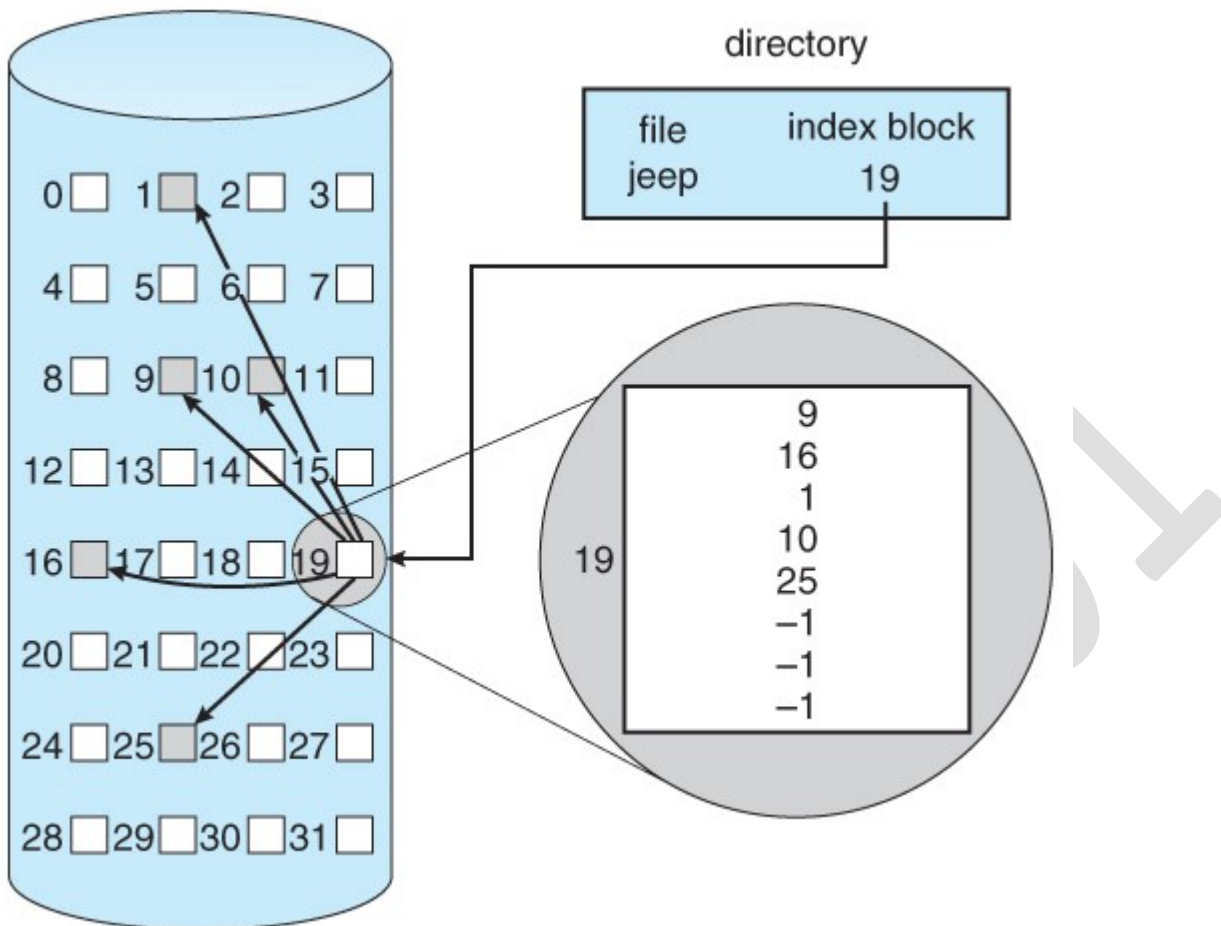
**Figure 12.7 File-allocation table.**

**12.4.3 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 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. )
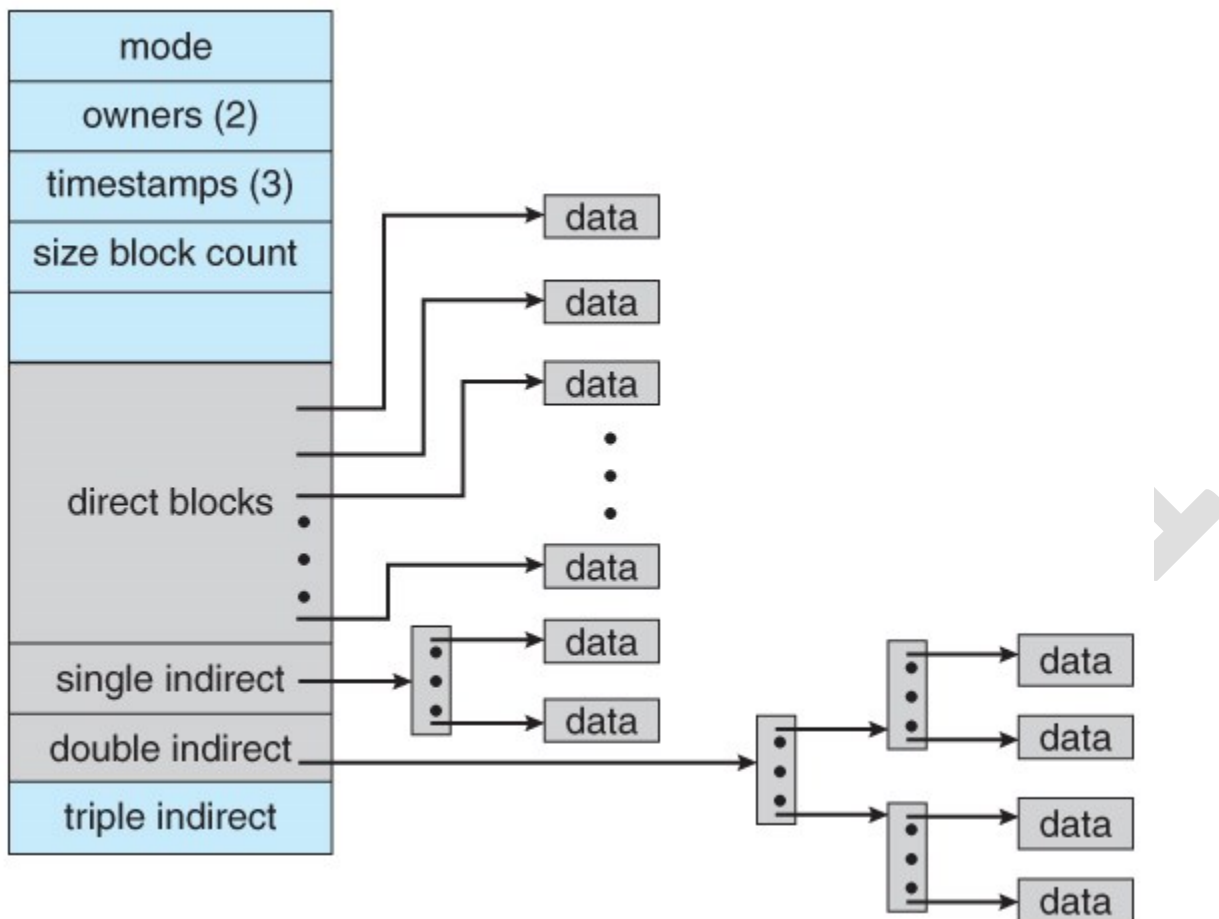
**Figure 12.9 - The UNIX inode.**

### 12.4.4 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.

## 12.5 Free-Space Management

- Another important aspect of disk management is keeping track of and allocating free space.

### 12.5.1 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

196

- The down side is that a 40GB disk requires over 5MB just to store the bitmap. ( For example. )

### 12.5.2 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.
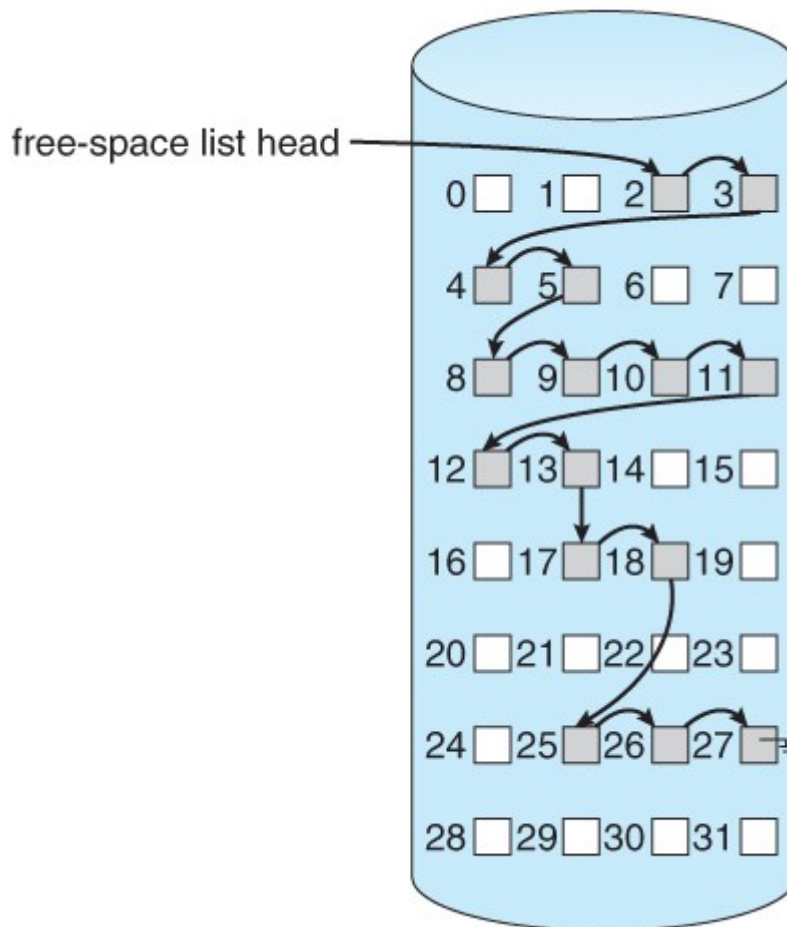


**Figure 12.10 - Linked free-space list on disk.**

### 12.5.3 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.

### 12.5.4 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. )
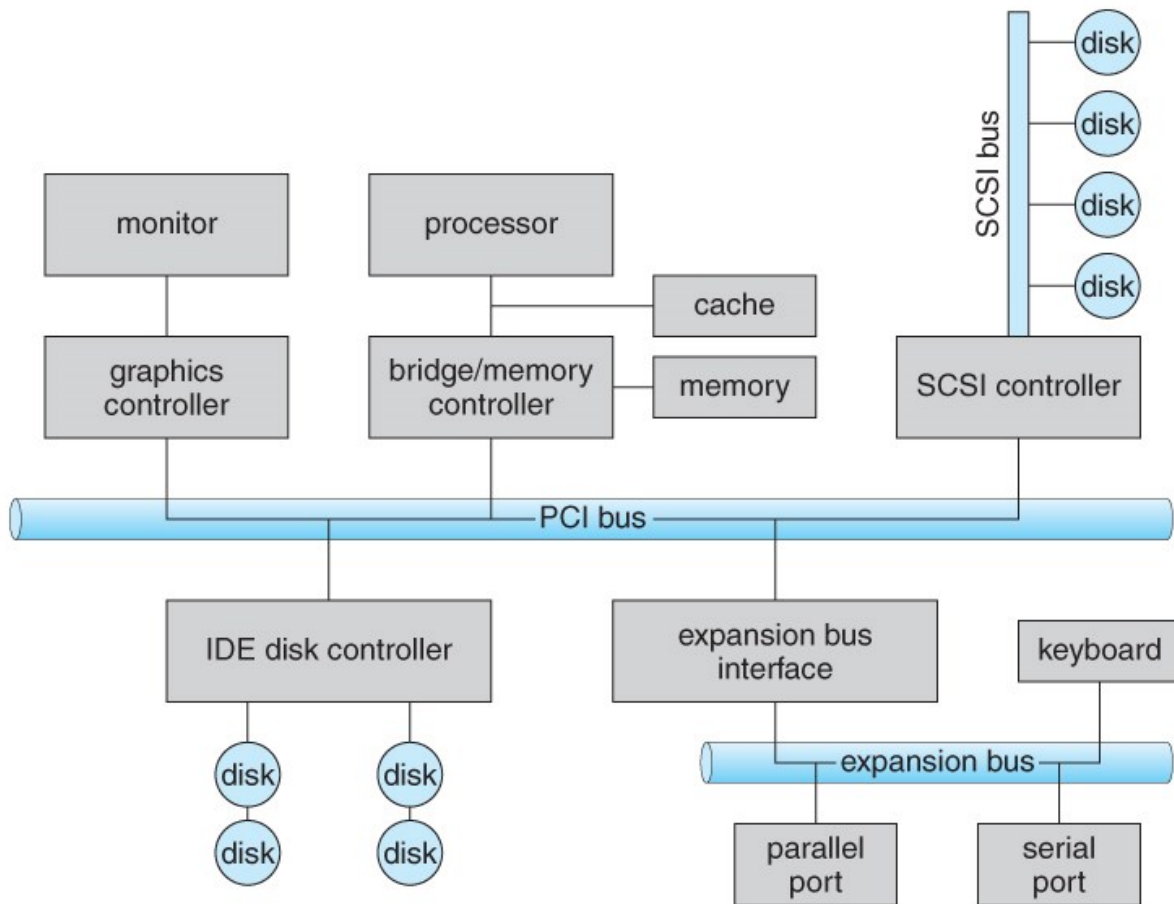
### 12.5.5 Space Maps

- 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.

## 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.*
  - o Buses include rigid protocols for the types of messages that can be sent across the bus and the procedures for resolving contention issues.
  - o 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.*
  - 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.
    - o 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.
    - o 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.

- o 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.
- o 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.
- o 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 |

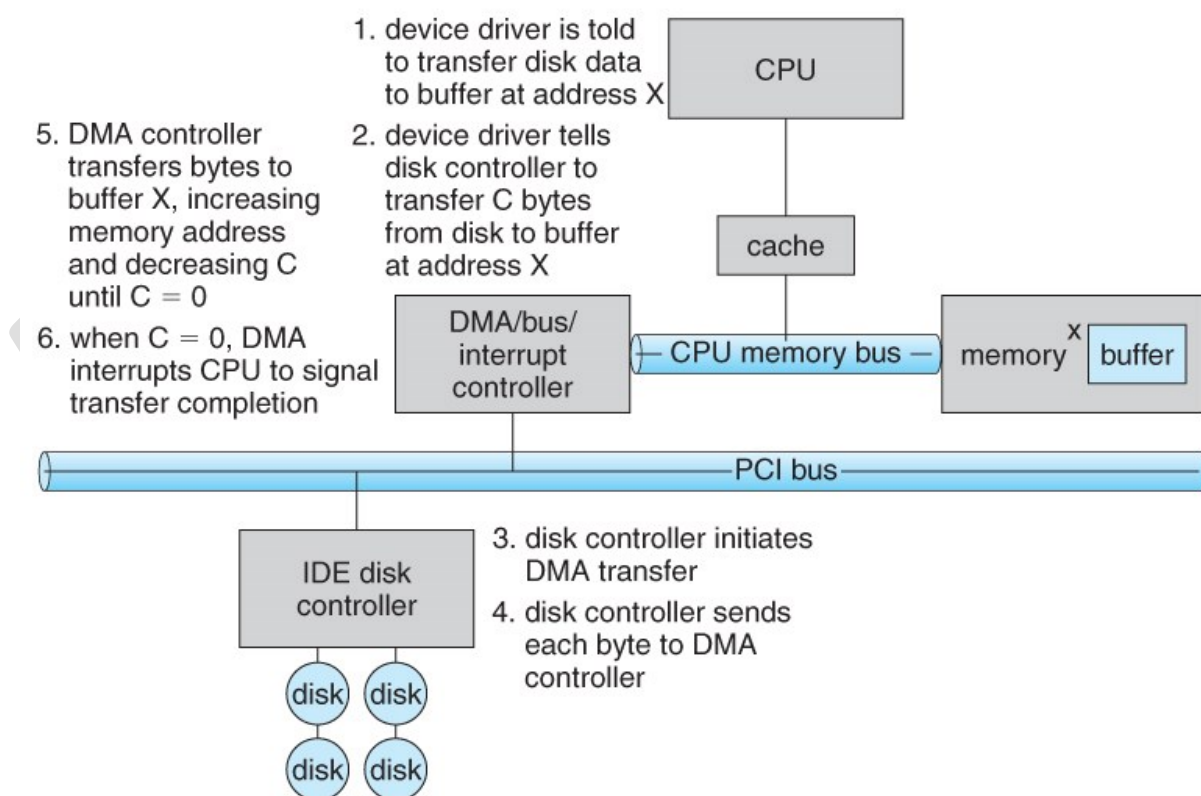**Figure 13.4 - Intel Pentium processor event-vector table.**

- 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. )
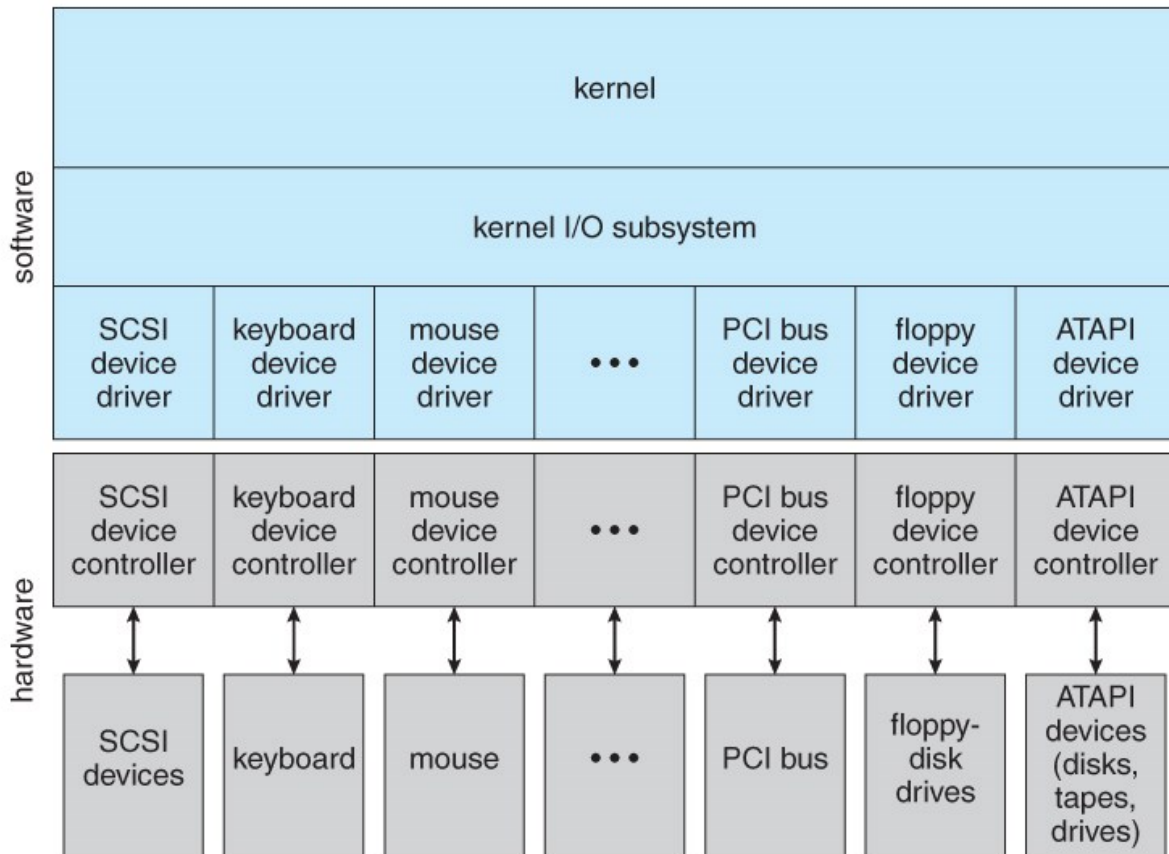- Figure 13.5 below illustrates the DMA process.



**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( ).
  - 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. )
  - 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.
  - 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.

207

- 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:
  - o Get the current time of day.
  - o Get the elapsed time ( system or wall clock ) since a previous event.
  - o 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.
  - o The scheduler uses a PIT to trigger interrupts for ending time slices.
  - o The disk system may use a PIT to schedule periodic maintenance cleanup, such as flushing buffers to disk.
  - o 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.
  - o 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. )
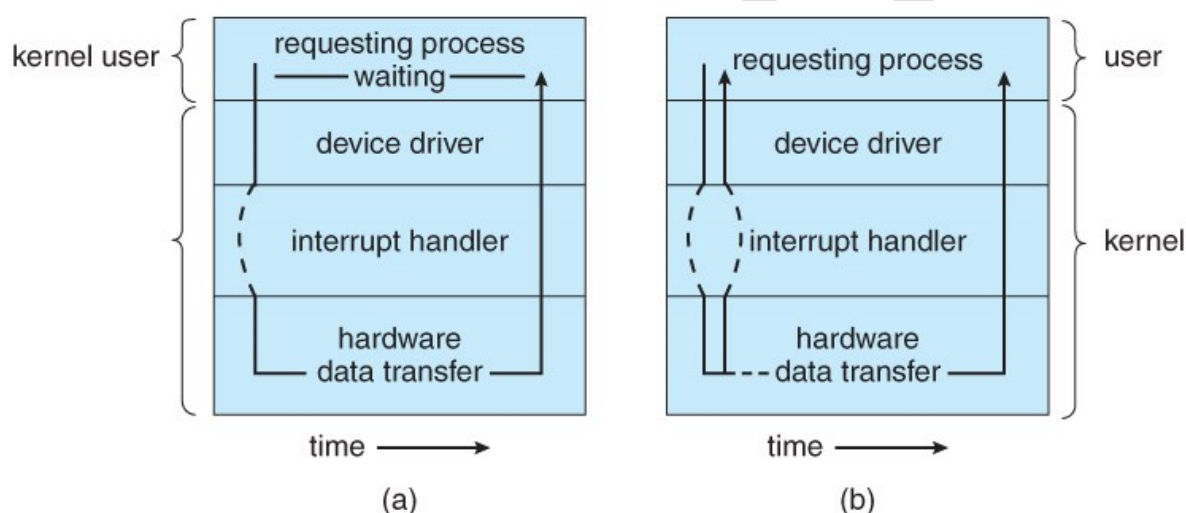


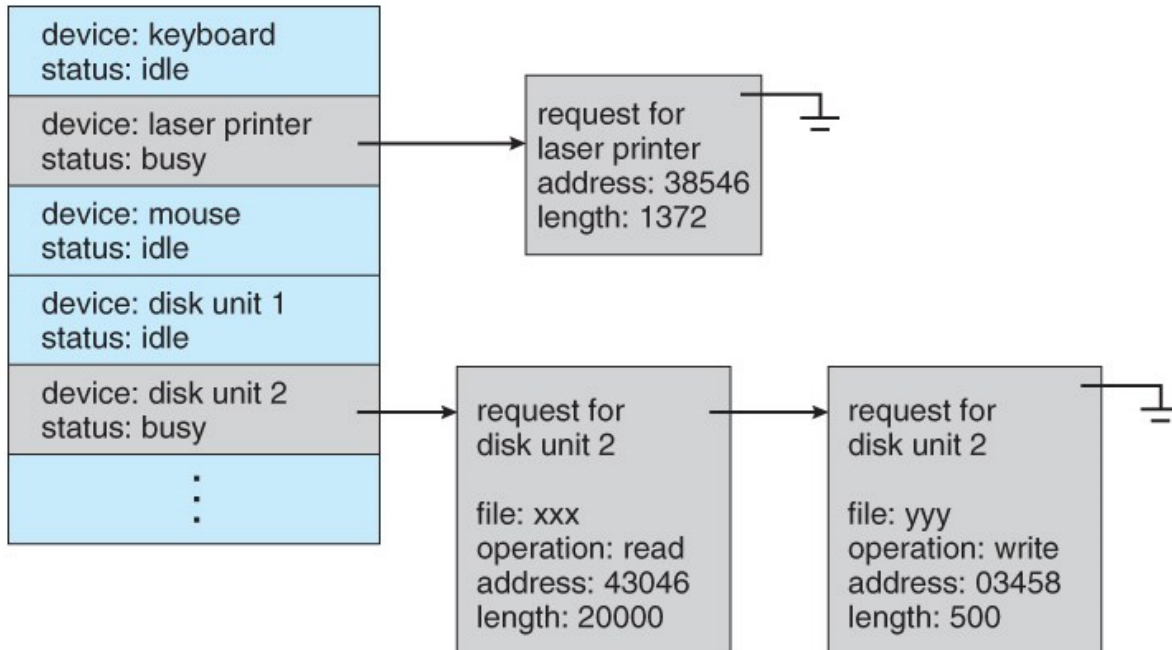**Figure 13.8 - Two I/O methods: (a) synchronous and (b) asynchronous.**

**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.
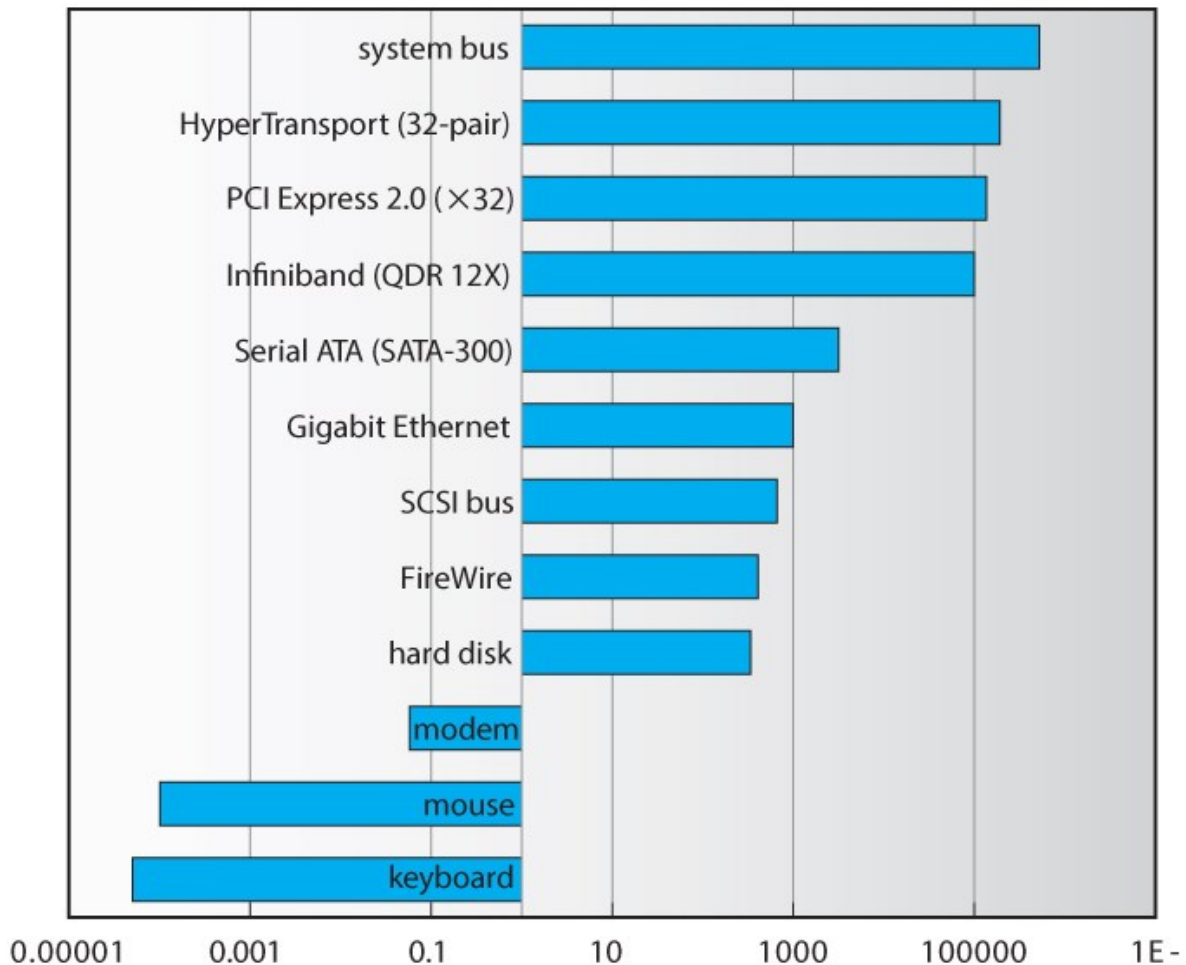
- On systems with many devices, separate request queues are often kept for each device:



**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 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.
- 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,

211

then the application sees that print job as complete, and the print scheduler sends each file to the appropriate printer one at a time.
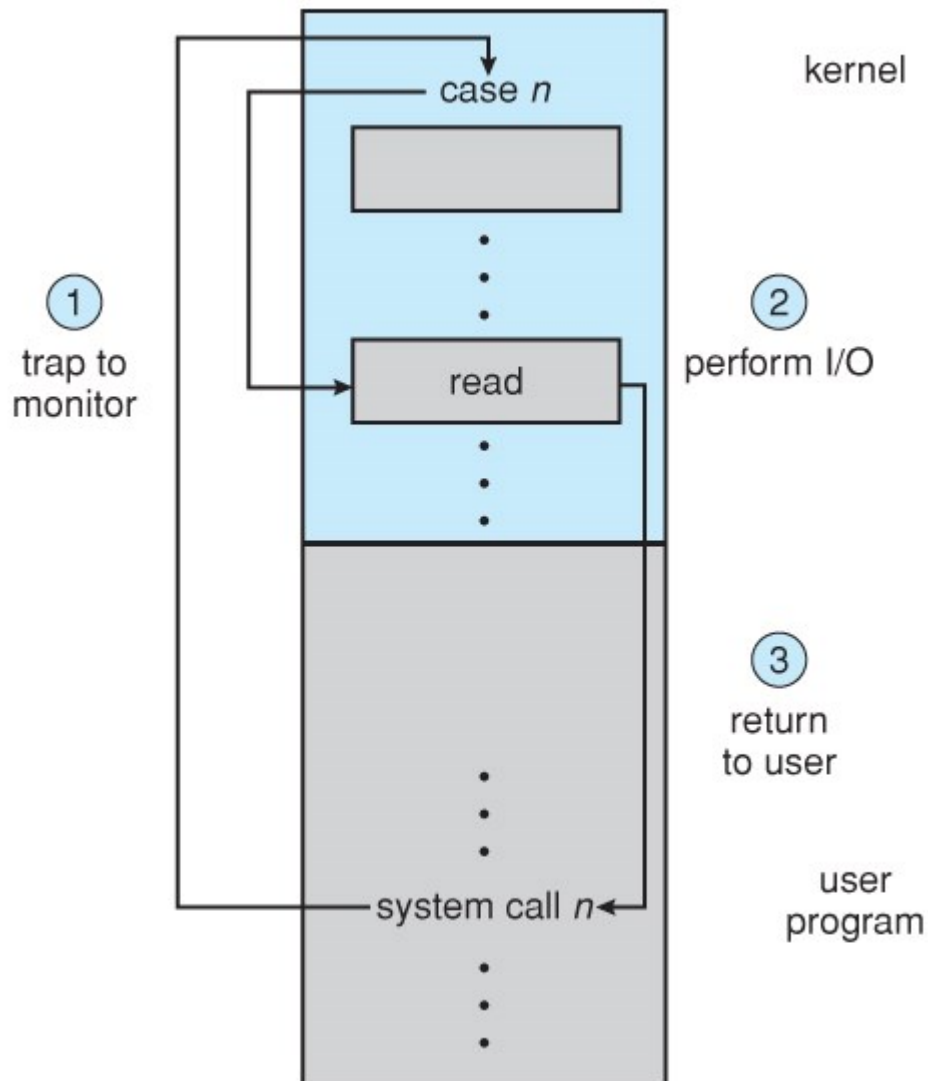
- Support is provided for viewing the spool queues, removing jobs from the queues, moving jobs from one queue to another queue, and in some cases changing the priorities of jobs in the queues.
- Spool queues can be general ( any laser printer ) or specific ( printer number 42. )
- OSes can also provide 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.
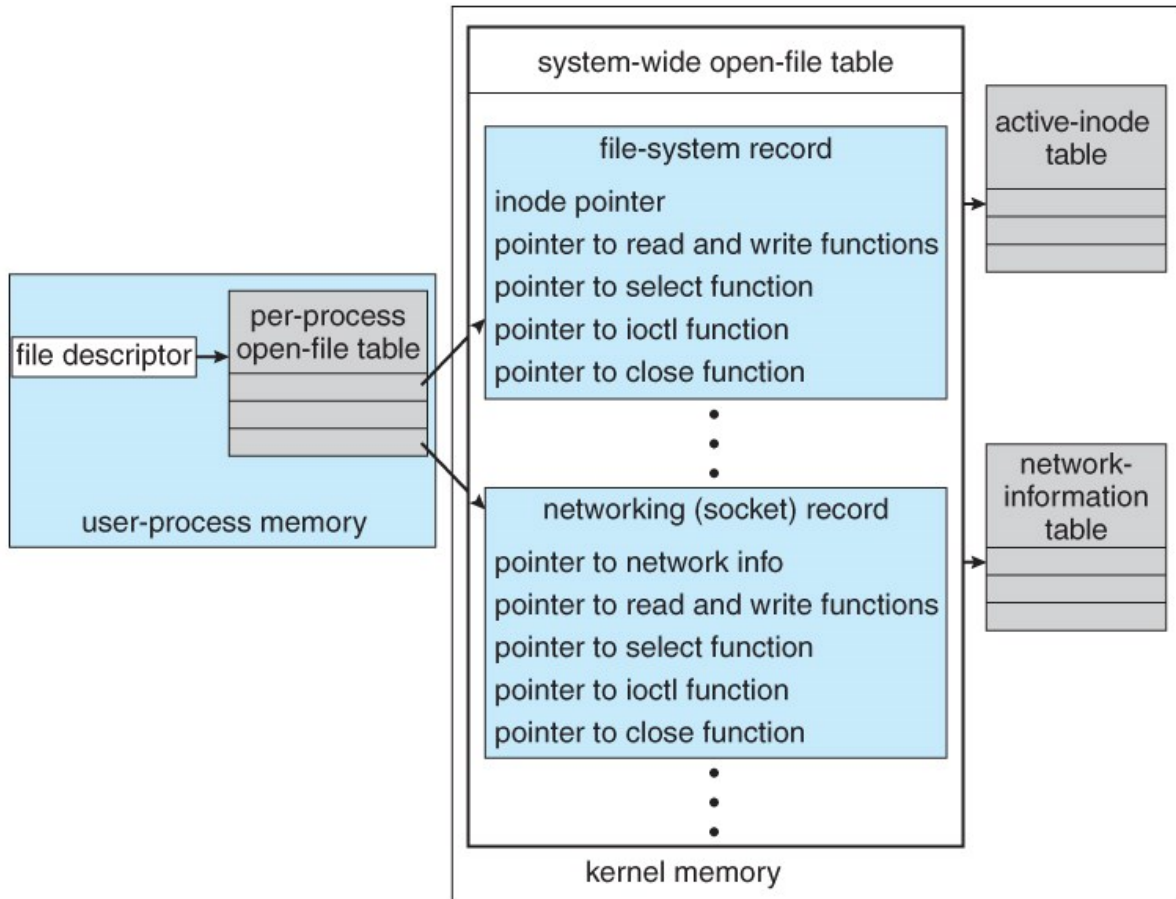
## 13.4.6 I/O Protection

- 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.
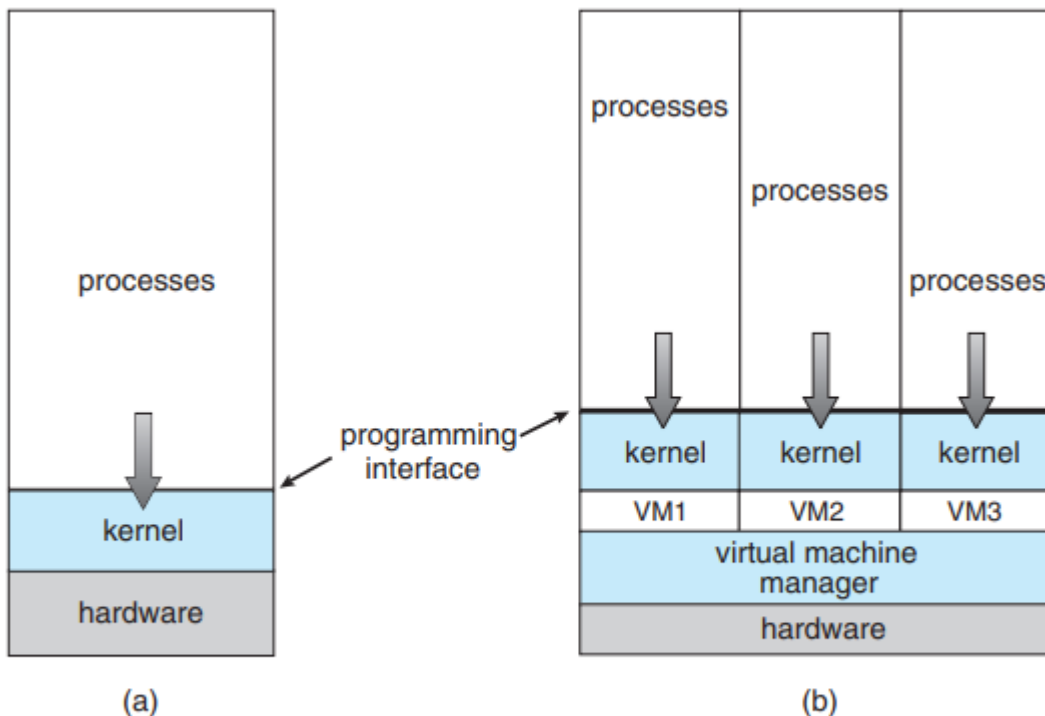
**Figure 13.12 - UNIX I/O kernel structure.**

# UNIT V

Overview The fundamental idea behind a virtual machine is to abstract the hardware of a single computer (the CPU, memory, disk drives, network interface cards, and so forth) into several different execution environments, thereby creating the illusion that each separate environment is running on its own private computer. This concept may seem similar to the layered approach of operating system implementation (see Section 2.8.2), and in some ways it is. In the case of virtualization, there is a layer that creates a virtual system on which operating systems or applications can run. Virtual machine implementations involve several components. At the base is the host, the underlying hardware system that runs the virtual machines. The virtual machine manager (VMM) (also known as a hypervisor) creates and runs virtual machines by providing an interface that is identical to the host (except in the case of paravirtualization, discussed later). Each guest process is provided with a virtual copy of the host (Figure 18.1). Usually, the guest process is in fact an operating system. A single physical machine can thus run multiple operating systems concurrently, each in its own virtual machine. Take a moment to note that with virtualization, the definition of "operating system" once again blurs. For example, consider VMM software such as VMware ESX. This virtualization software is installed on the hardware, runs when the hardware boots, and provides services to applications. The services include traditional ones, such as scheduling and memory management, along with new types, such as migration of applications between systems. Furthermore, the applications are, in fact, guest operating systems. Is the VMware ESX VMM an operating system that, in turn, runs other operating systems? Certainly it acts like an operating system. For clarity, however, we call the component that provides virtual environments a VMM

The implementation of VMMs varies greatly. Options include the following: • Hardware-based solutions that provide support for virtual machine creation and management via firmware. These VMMs, which are commonly found in mainframe and large to midsized servers, are generally known as type 0 hypervisors. IBM LPARs and Oracle LDOMs are examples



**Figure 18.1** System models. (a) Nonvirtual machine. (b) Virtual machine.

• Operating-system-like software built to provide virtualization, including VMware ESX (mentioned above), Joyent SmartOS, and Citrix XenServer. These VMMs are known as type 1 hypervisors. • General-purpose operating systems that provide standard functions as well as VMM functions,

including Microsoft Windows Server with HyperV and Red Hat Linux with the KVM feature. Because such systems have a feature set similar to type 1 hypervisors, they are also known as type 1.

• Applications that run on standard operating systems but provide VMM features to guest operating systems. These applications, which include VMware Workstation and Fusion, Parallels Desktop, and Oracle VirtualBox, are type 2 hypervisors.

• Paravirtualization, a technique in which the guest operating system is modified to work in cooperation with the VMM to optimize performance.

• Programming-environment virtualization, in which VMMs do not virtualize real hardware but instead create an optimized virtual system. This technique is used by Oracle Java and Microsoft.Net.

• Emulators that allow applications written for one hardware environment to run on a very different hardware environment, such as a different type of CPU.

• Application containment, which is not virtualization at all but rather provides virtualization-like features by segregating applications from the operating system. Oracle Solaris Zones, BSD Jails, and IBM AIX WPARs "contain" applications, making them more secure and manageable.

**History**

Virtual machines first appeared commercially on IBM mainframes in 1972. Virtualization was provided by the IBM VM operating system. This system has evolved and is still available. In addition, many of its original concepts are found in other systems, making it worth exploring IBM VM/370 divided a mainframe into multiple virtual machines, each running its own operating system. A major difficulty with the VM approach involved disk systems. Suppose that the physical machine had three disk drives but wanted to support seven virtual machines. Clearly, it could not allocate a disk drive to each virtual machine. The solution was to provide virtual disks— termed minidisks in IBM's VM operating system. The minidisks were identical to the system's hard disks in all respects except size. The system implemented each minidisk by allocating as many tracks on the physical disks as the minidisk needed. Once the virtual machines were created, users could run any of the operating systems or software packages that were available on the underlying machine. For the IBM VM system, a user normally ran CMS—a single-user interactive operating system. For many years after IBM introduced this technology, virtualization remained in its domain. Most systems could not support virtualization. However, a formal definition of virtualization helped to establish system requirements and a target for functionality. The virtualization requirements called for:

• Fidelity. A VMM provides an environment for programs that is essentially identical to the original machine.

• Performance. Programs running within that environment show only minor performance decreases.

• Safety. The VMM is in complete control of system resources.

**Benefits and Features**

Several advantages make virtualization attractive. Most of them are fundamentally related to the ability to share the same hardware yet run several different execution environments (that is, different operating systems) concurrently. One important advantage of virtualization is that the host system is protected from the virtual machines, just as the virtual machines are protected from each other. A virus inside a guest operating system might damage that operating system but is unlikely to affect the host or the other guests. Because each virtual machine is almost completely isolated from all other virtual machines, there are almost no protection problems. A potential disadvantage of isolation is that it can prevent sharing of resources. Two approaches to providing sharing have been implemented. First, it is possible to share a file-system volume and thus to share files. Second, it is possible to define a network of virtual machines, each of which can send information over the virtual communications network. The network is modeled after physical communication networks but is implemented in software. Of course, the VMM is free to allow any number of its guests to use physical resources, such as a physical network connection (with sharing provided by the VMM), in which case the allowed guests could communicate with each other via the physical network.

A virtual-machine system can eliminate much of this latter problem. System programmers are given their own virtual machine, and system development is done on the virtual machine instead of on a
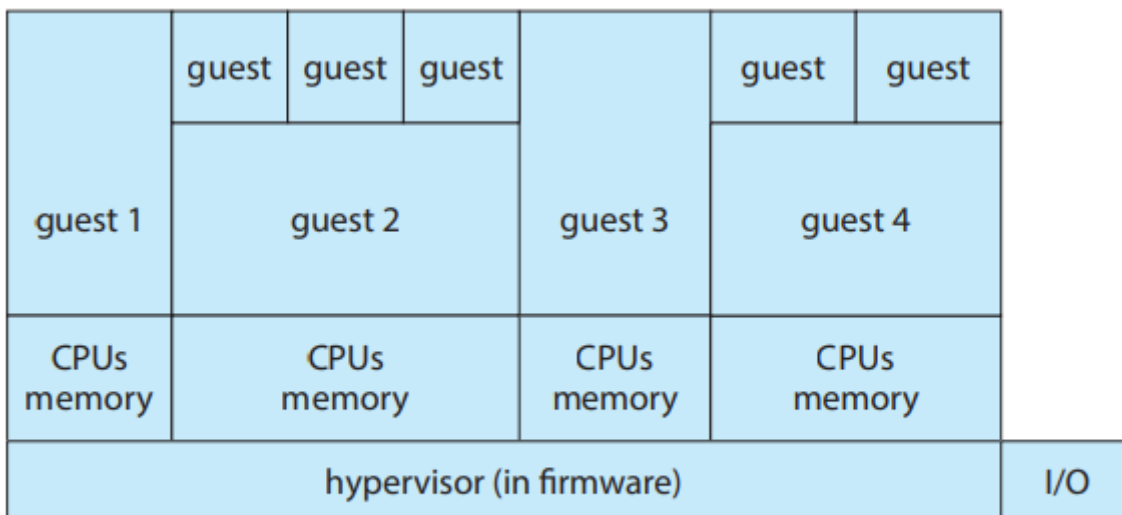
physical machine. Normal system operation is disrupted only when a completed and tested change is ready to be put into production.

A major advantage of virtual machines in production data-center use is system consolidation, which involves taking two or more separate systems and running them in virtual machines on one system. Such physical-to-virtual conversions result in resource optimization, since many lightly used systems can be combined to create one more heavily used system. Consider, too, that management tools that are part of the VMM allow system administrators to manage many more systems than they otherwise could. A virtual environment might include 100 physical servers, each running 20 virtual servers. Without virtualization, 2,000 servers would require several system administrators. With virtualization and its tools, the same work can be managed by one or two administrators. One of the tools that make this possible is templating, in which one standard virtual machine image, including an installed and configured guest operating system and applications, is saved and used as a source for multiple running VMs. Other features include managing the patching of all guests, backing up and restoring the guests, and monitoring their resource use. Virtualization can improve not only resource utilization but also resource management. Some VMMs include a live migration feature that moves a running guest from one physical server to another without interrupting its operation or active network connections. If a server is overloaded, live migration can thus free resources on the source host while not disrupting the guest. Similarly, when host hardware must be repaired or upgraded, guests can be migrated to other servers, the evacuated host can be maintained, and then the guests can be migrated back. This operation occurs without downtime and without interruption to users.

## Types of VMs and Their Implementations

Type 0 Hypervisor Type 0 hypervisors have existed for many years under many names, including "partitions" and "domains." They are a hardware feature, and that brings its own positives and negatives. Operating systems need do nothing special to take advantage of their features. The VMM itself is encoded in the firmware and loaded at boot time. In turn, it loads the guest images to run in each partition. The feature set of a type 0 hypervisor tends to be smaller than those of the other types because it is implemented in hardware. For example, a system might be split into four virtual systems, each with dedicated CPUs, memory, and I/O devices. Each guest believes that it has dedicated hardware because it does, simplifying many implementation details.

| guest 1 | guest | guest | guest | guest 3 | guest | guest |
| | guest 2 | | | | guest 4 | |
| CPUs memory | CPUs memory | | CPUs memory | | CPUs memory | |
| hypervisor (in firmware) | | | | | | I/O |

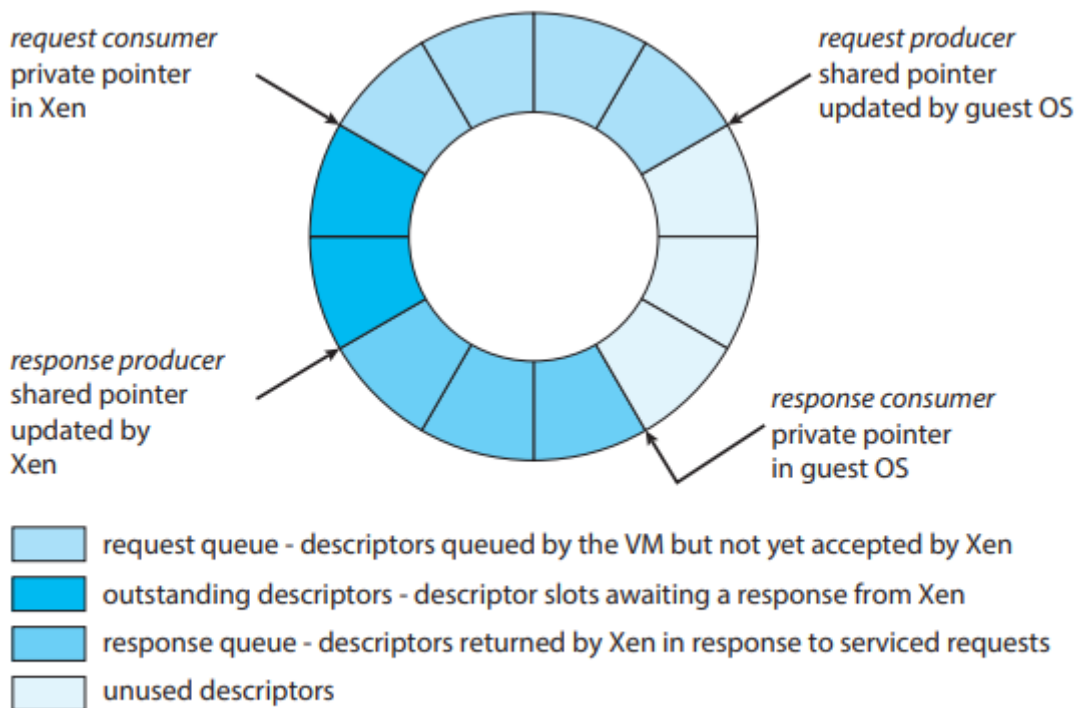**Figure 18.5**  Type 0 hypervisor.

I/O presents some difficulty, because it is not easy to dedicate I/O devices to guests if there are not enough. What if a system has two Ethernet ports and more than two guests, for example? Either all guests must get their own I/O devices, or the system must provided I/O device sharing. In these cases,

the hypervisor manages shared access or grants all devices to a control partition.

**Type 1 Hypervisor**

Type 1 hypervisors are commonly found in company data centers and are, in a sense, becoming "the data-center operating system." They are special-purpose operating systems that run natively on the hardware, but rather than providing system calls and other interfaces for running programs, they create, run, and manage guest operating systems. In addition to running on standard hardware, they can run on type 0 hypervisors, but not on other type 1 hypervisors. Whatever the platform, guests generally do not know they are running on anything but the native hardware. Type 1 hypervisors run in kernel mode, taking advantage of hardware protection. Where the host CPU allows, they use multiple modes to give guest operating systems their own control and improved performance. They implement device drivers for the hardware they run on, since no other component could do so. Because they are operating systems, they must also provide CPU scheduling, memory management, I/O management, protection, and even security. Frequently, they provide APIs, but those APIs support applications in guests or external applications that supply features like backups, monitoring, and security. Many type 1 hypervisors are closed-source commercial offerings, such as VMware ESX, while some are open source or hybrids of open and closed source, such as Citrix XenServer and its open Xen counterpart.

**Type 2 Hypervisor**



request consumer
private pointer
in Xen

request producer
shared pointer
updated by guest OS

response producer
shared pointer
updated by
Xen

response consumer
private pointer
in guest OS

- request queue - descriptors queued by the VM but not yet accepted by Xen
- outstanding descriptors - descriptor slots awaiting a response from Xen
- response queue - descriptors returned by Xen in response to serviced requests
- unused descriptors

**Figure 18.6** Xen I/O via shared circular buffer.[1]

This type of VMM is simply another process run and managed by the host, and even the host does not know that virtualization is happening within the VMM. Type 2 hypervisors have limits not associated with some of the other types. For example, a user needs administrative privileges to access many of the hardware assistance features of modern CPUs. If the VMM is being run by a standard user without additional privileges, the VMM cannot take advantage of these features. Due to this limitation, as well as the extra overhead of running a general-purpose operating system as well as guest operating systems, type 2 hypervisors tend to have poorer overall performance than type 0 or

type 1. As is often the case, the limitations of type 2 hypervisors also provide some benefits. They run on a variety of general-purpose operating systems, and running them requires no changes to the host operating system. A student can use a type 2 hypervisor, for example, to test a non-native operating system without replacing the native operating system. In fact, on an Apple laptop, a student could have versions of Windows, Linux, Unix, and less common operating systems all available for learning and experimentation.

## Paravirtualization

The Xen VMM became the leader in paravirtulization by implementing several techniques to optimize the performance of guests as well as of the host system. For example, as mentioned earlier, some VMMs present virtual devices to guests that appear to be real devices. Instead of taking that approach, the Xen VMM presented clean and simple device abstractions that allow efficient I/O as well as good I/O-related communication between the guest and the VMM. For memory management, Xen did not implement nested page tables. Rather, each guest had its own set of page tables, set to read-only. Xen required the guest to use a specific mechanism, a hypercall from the guest to the hypervisor VMM, when a page-table change was needed. This meant that the guest operating system's kernel code must have been changed from the default code to these Xen-specific methods. To optimize performance, Xen allowed the guest to queue up multiple page-table changes asynchronously via hypercalls and then checked to ensure that the changes were complete before continuing operation
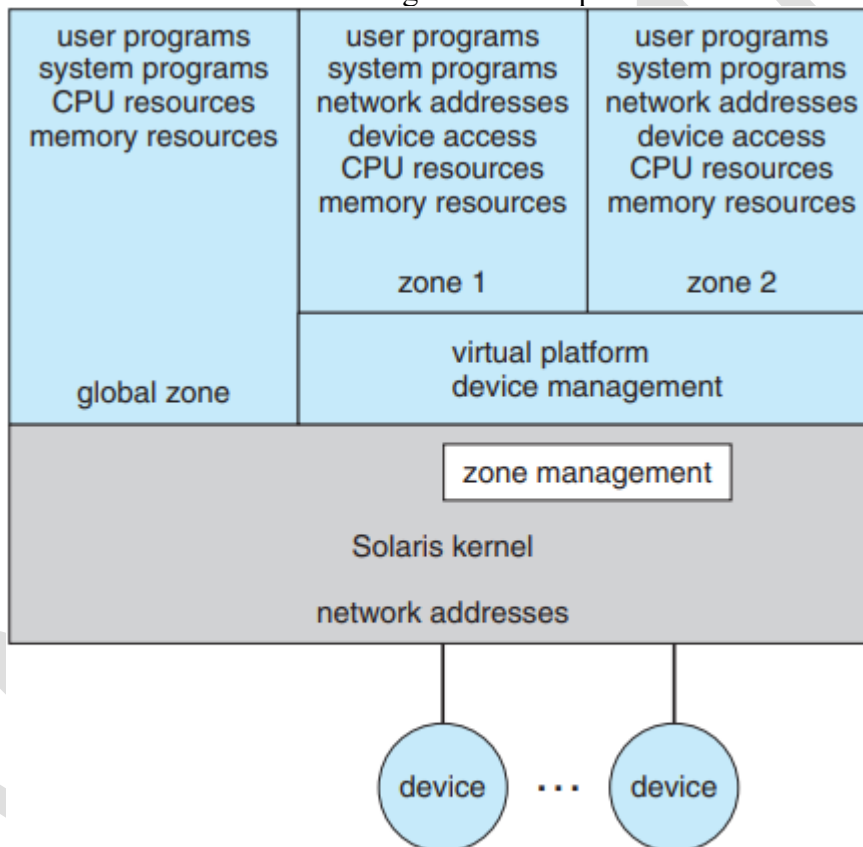
## Programming-Environment Virtualization

Another kind of virtualization, based on a different execution model, is the virtualization of programming environments. Here, a programming language is designed to run within a custom-built virtualized environment. For example, Oracle's Java has many features that depend on its running in the Java virtual machine (JVM), including specific methods for security and memory management. If we define virtualization as including only duplication of hardware, this is not really virtualization at all. But we need not limit ourselves to that definition. Instead, we can define a virtual environment, based on APIs, that provides a set of features we want to have available for a particular language and programs written in that language. Java programs run within the JVM environment, and the JVM is compiled to be a native program on systems on which it runs. This arrangement means that Java programs are written once and then can run on any system (including all of the major operating systems) on which a JVM is available. The same can be said of interpreted languages, which run inside programs that read each instruction and interpret it into native operations.

## Emulation

Virtualization is probably the most common method for running applications designed for one operating system on a different operating system, but on the same CPU. This method works relatively efficiently because the applications were compiled for the instruction set that the target system uses. But what if an application or operating system needs to run on a different CPU? Here, it is necessary to translate all of the source CPU's instructions so that they are turned into the equivalent instructions of the target CPU. Such an environment is no longer virtualized but rather is fully emulated. Emulation is useful when the host system has one system architecture and the guest system was compiled for a different architecture. For example, suppose a company has replaced its outdated computer system with a new system but would like to continue to run certain important programs that were compiled for the old system. The programs could be run in an emulator that translates each of the outdated system's instructions into the native instruction set of the new system. Emulation can increase the life of programs and allow us to explore old architectures without having an actual old machine

## Application Containment

The goal of virtualization in some instances is to provide a method to segregate applications, manage their performance and resource use, and create an easy way to start, stop, move, and manage them. In such cases, perhaps full-fledged virtualization is not needed. If the applications are all compiled for the same operating system, then we do not need complete virtualization to provide these features. We can instead use application containment. Consider one example of application containment. Starting with version 10, Oracle Solaris has included containers, or zones, that create a virtual layer between the operating system and the applications. In this system, only one kernel is installed, and the hardware is not virtualized. Rather, the operating system and its devices are virtualized, providing processes within a zone with the impression that they are the only processes on the system. One or more containers can be created, and each can have its own applications, network stacks, network address and ports, user accounts, and so on. CPU and memory resources can be divided among the zones and the system-wide processes. Each zone, in fact, can run its own scheduler to optimize the performance of its applications on the allotted resources. Figure 18.7 shows a Solaris 10 system with two containers and the standard "global" user space.



**Figure 18.7** Solaris 10 with two zones.

## Virtualization and Operating-System Components
### CPU Scheduling

A system with virtualization, even a single-CPU system, frequently acts like a multiprocessor system. The virtualization software presents one or more virtual CPUs to each of the virtual machines running on the system and then schedules the use of the physical CPUs among the virtual machines. The significant variations among virtualization technologies make it difficult to summarize the effect of virtualization on scheduling. First, let's consider the general case of VMM scheduling. The VMM has a number of physical CPUs available and a number of threads to run on those CPUs. The threads can be VMM threads or guest threads. Guests are configured with a certain number of virtual CPUs at creation time, and that number can be adjusted throughout the life of the

VM. When there are enough CPUs to allocate the requested number to each guest, the VMM can treat the CPUs as dedicated and schedule only a given guest's threads on that guest's CPUs. In this situation, the guests act much like native operating systems running on native CPUs. Of course, in other situations, there may not be enough CPUs to go around. The VMM itself needs some CPU cycles for guest management and I/O management and can steal cycles from the guests by scheduling its threads across all of the system CPUs, but the impact of this action is relatively minor. More difficult is the case of overcommitment, in which the guests are configured for more CPUs than exist in the system. Here, a VMM can use standard scheduling algorithms to make progress on each thread but can also add a fairness aspect to those algorithms. For example, if there are six hardware CPUs and twelve guestallocated CPUs, the VMM can allocate CPU resources proportionally, giving each guest half of the CPU resources it believes it has. The VMM can still present all twelve virtual CPUs to the guests, but in mapping them onto physical CPUs, the VMM can use its scheduler to distribute them appropriately.

## Memory Management

Efficient memory use in general-purpose operating systems is a major key to performance. In virtualized environments, there are more users of memory (the guests and their applications, as well as the VMM), leading to more pressure on memory use. Further adding to this pressure is the fact that VMMs typically overcommit memory, so that the total memory allocated to guests exceeds the amount that physically exists in the system. The extra need for efficient memory use is not lost on the implementers of VMMs, who take extensive measures to ensure the optimal use of memory

1. Recall that a guest believes it controls memory allocation via its pagetable management, whereas in reality the VMM maintains a nested page table that translates the guest page table to the real page table. The VMM can use this extra level of indirection to optimize the guest's use of memory without the guest's knowledge or help. One approach is to provide double paging. Here, the VMM has its own page-replacement algorithms and loads pages into a backing store that the guest believes is physical memory. Of course, the VMM knows less about the guest's memory access patterns than the guest does, so its paging is less efficient, creating performance problems. VMMs do use this method when other methods are not available or are not providing enough free memory. However, it is not the preferred approach.

2. A common solution is for the VMM to install in each guest a pseudo– device driver or kernel module that the VMM controls. (Apseudo–device driver uses device-driver interfaces, appearing to the kernel to be a device driver, but does not actually control a device. Rather, it is an easy way to add kernel-mode code without directly modifying the kernel.) This balloon memory manager communicates with the VMM and is told to allocate or deallocate memory. If told to allocate, it allocates memory and tells the operating system to pin the allocated pages into physical memory

3. Another common method for reducing memory pressure is for the VMM to determine if the same page has been loaded more than once. If this is the case, the VMM reduces the number of copies of the page to one and maps the other users of the page to that one copy. VMware, for example, randomly samples guest memory and creates a hash for each page sampled. That hash value is a "thumbprint" of the page. The hash of every page examined is compared with other hashes stored in a hash table. If there is a match, the pages are compared byte by byte to see if they really are identical. If they are, one page is freed, and its logical address is mapped to the other's physical address.

## I/O

In the area of I/O, hypervisors have some leeway and can be less concerned with how they represent the underlying hardware to their guests. Because of the wide variation in I/O devices, operating systems are used to dealing with varying and flexible I/O mechanisms. For example, an operating system's device-driver mechanism provides a uniform interface to the operating system whatever
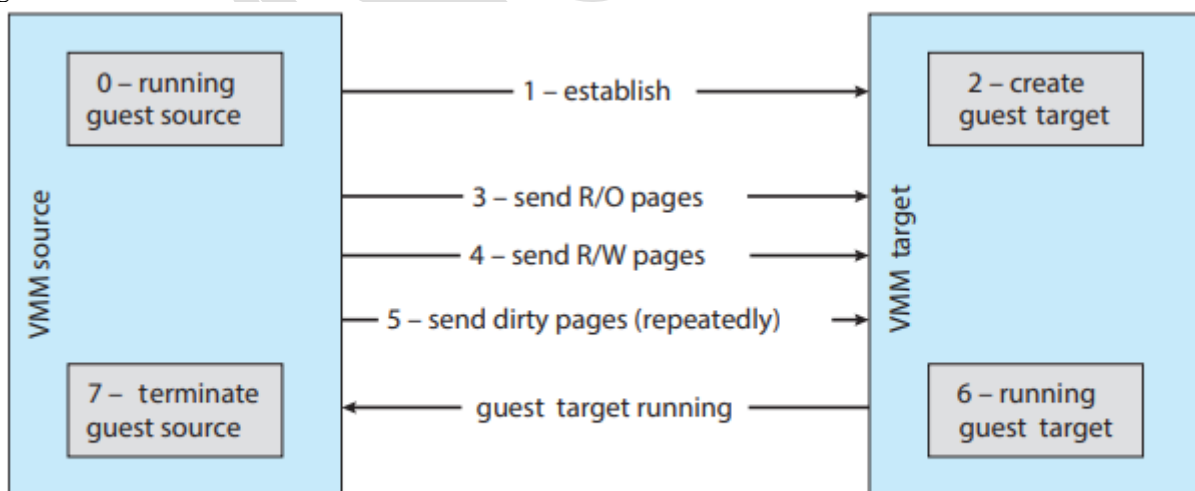
the I/O device. Device-driver interfaces are designed to allow third-party hardware manufacturers to provide device drivers connecting their devices to the operating system. Usually, device drivers can be dynamically loaded and unloaded. Virtualization takes advantage of this built-in flexibility by providing specific virtualized devices to guest operating systems.

## Storage Management

Type 1 hypervisors store the guest root disk (and configuration information) in one or more files in the file systems provided by the VMM. Type 2 hypervisors store the same information in the host operating system's file systems. In essence, a disk image, containing all of the contents of the root di of the guest, is contained in one file in the VMM. Aside from the potential performance problems that causes, this is a clever solution, because it simplifies copying and moving guests. If the administrator wants a duplicate of the guest (for testing, for example), she simply copies the associated disk image of the guest and tells the VMM about the new copy. Booting the new virtual machine brings up an identical guest. Moving a virtual machine from one system to another that runs the same VMM is as simple as halting the guest, copying the image to the other system, and starting the guest there

## Live Migration
1. The source VMM establishes a connection with the target VMM and confirms that it is allowed to send a guest.
2. The target creates a new guest by creating a new VCPU, new nested page table, and other state storage.
3. The source sends all read-only memory pages to the target.
4. The source sends all read–write pages to the target, marking them as clean.
 5. The source repeats step 4, because during that step some pages were probably modified by the guest and are now dirty. These pages need to be sent again and marked again as clean.
 6. When the cycle of steps 4 and 5 becomes very short, the source VMM freezes the guest, sends the VCPU's final state, other state details, and the final dirty pages, and tells the target to start running the guest. Once the target acknowledges that the guest is running, the source terminates the guest.
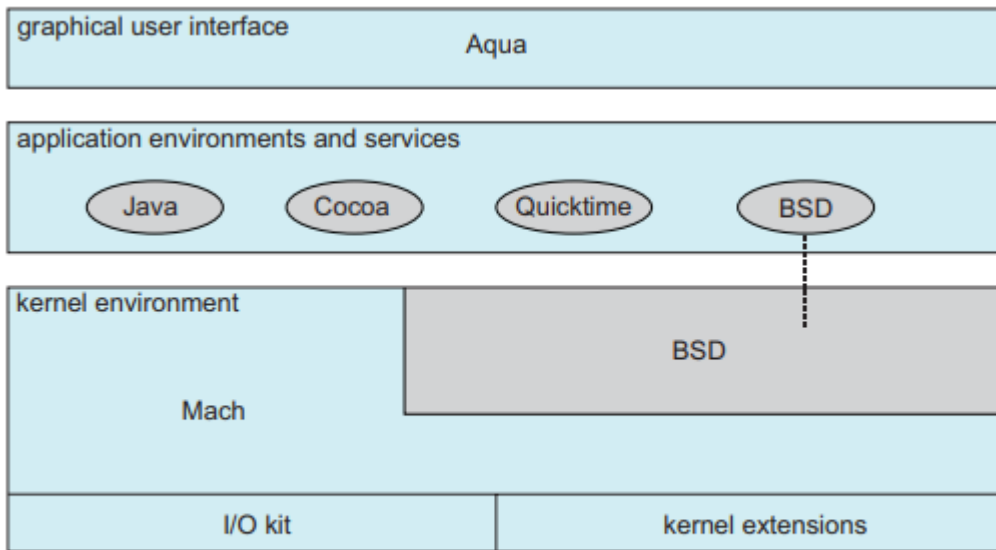


**Figure 18.8**  Live migration of a guest between two servers.

## 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 2.17. 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.
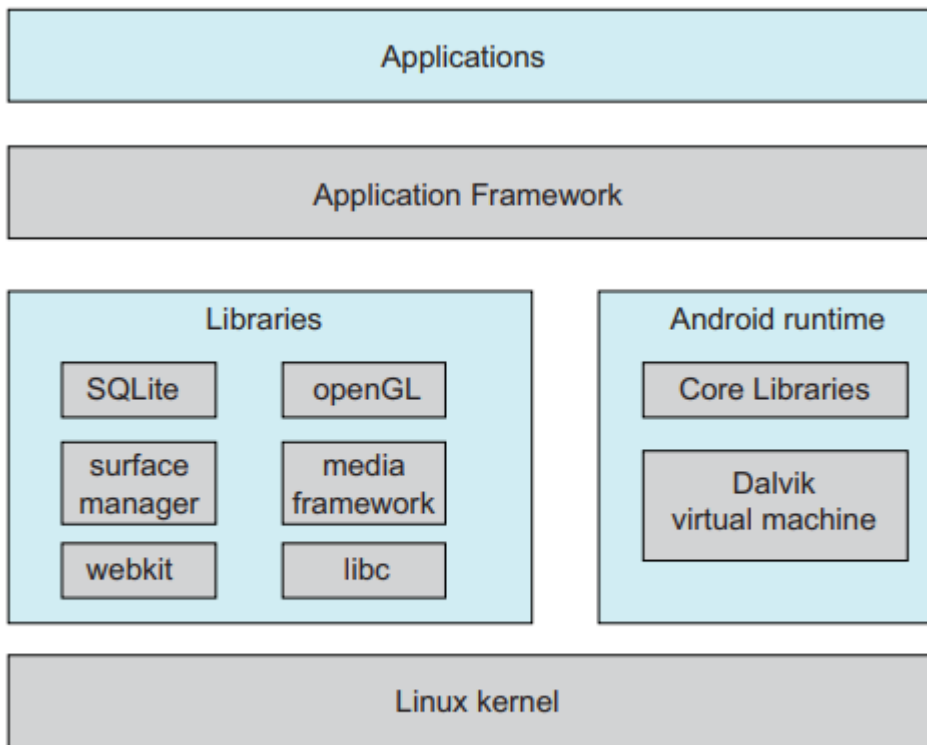


**Figure 2.16** The Mac OS X structure.



**Figure 2.17** Architecture of Apple's iOS. The core services layer provides a variety of features, including support for cloud computing and databases. The bottom layer represents the core operating system, which is based on the kernel environment shown in Figure 2.16

**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.

**Figure 2.18** Architecture of Google's Android.

Linux is used primarily for process, memory, and device-driver support for hardware and has been expanded to include power management. The Android runtime environment includes a core set of libraries as well as the Dalvik virtual machine. 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