

# Unit 1: Introduction

- What is an Operating System?
- Objectives and Functions
- Computer System Architecture
- OS Structure
- OS Operations
- Evolution of Operating System
- Distributed Systems
- Clustered System
- Real -Time Systems
- Special – Purpose Systems
- Operating System Services
- System Calls
- System Programs
- Operating System Design and Implementation
- OS Structure
- Virtual Machines

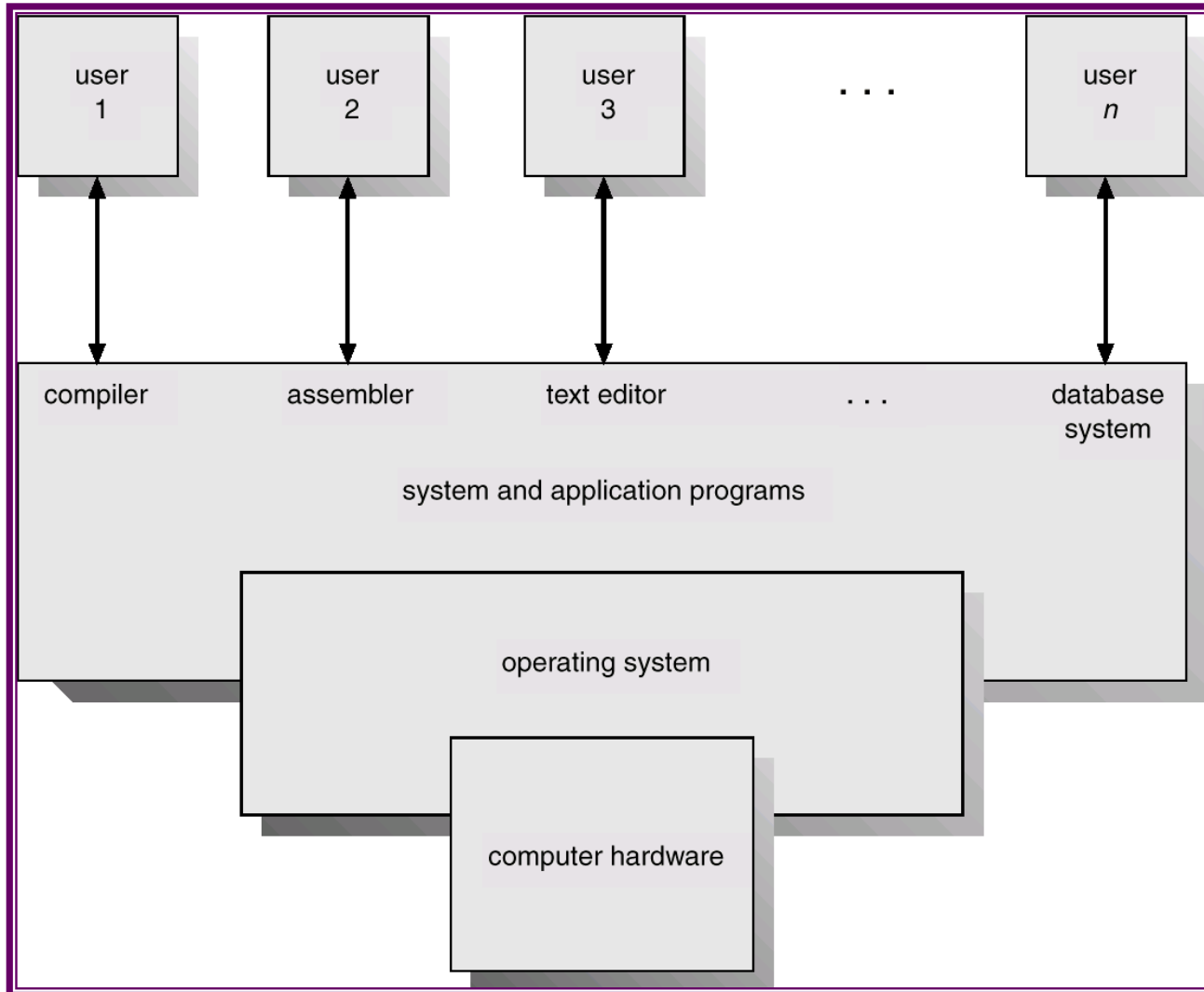
# What is an Operating System?

- 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 Components

1. Hardware – provides basic computing resources (CPU, memory, I/O devices).
2. Operating system – controls and coordinates the use of the hardware among the various application programs for the various users.
3. Applications programs – define the ways in which the system resources are used to solve the computing problems of the users (compilers, database systems, video games, business programs).
4. Users (people, machines, other computers).

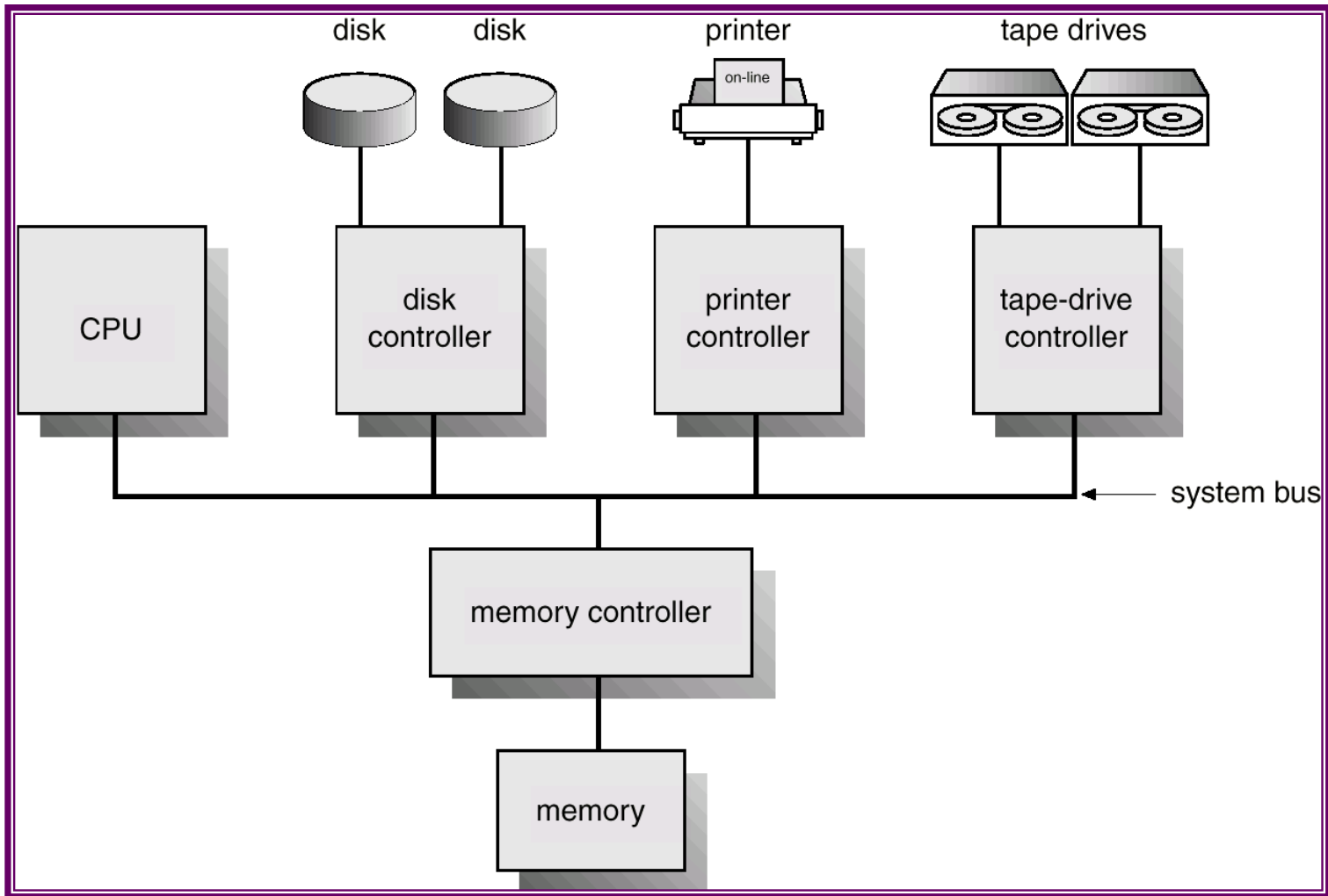
# Abstract View of System Components



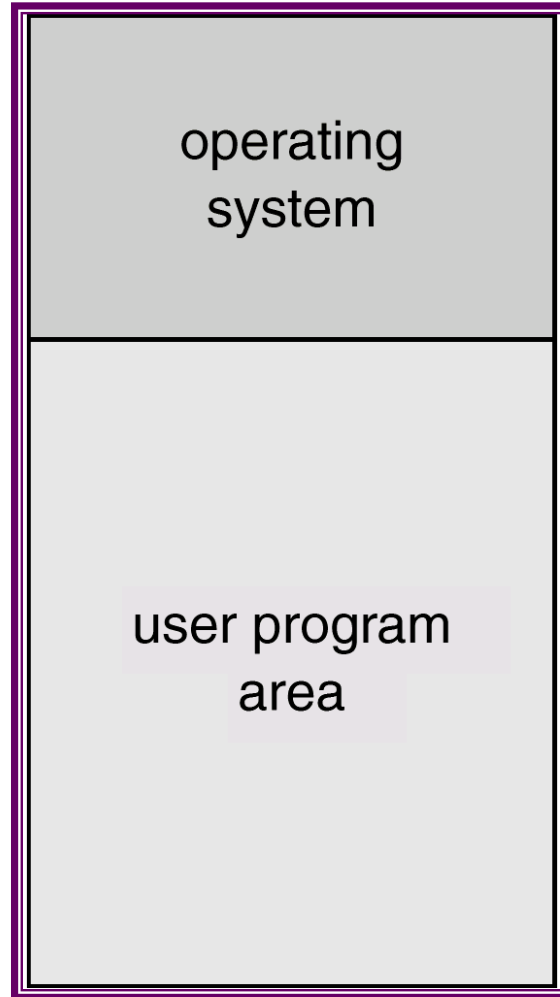
# Operating System Definitions

- Resource allocator – manages and allocates resources.
- Control program – controls the execution of user programs and operations of I/O devices .
- Kernel – the one program running at all times (all else being application programs).

# Computer-System Architecture

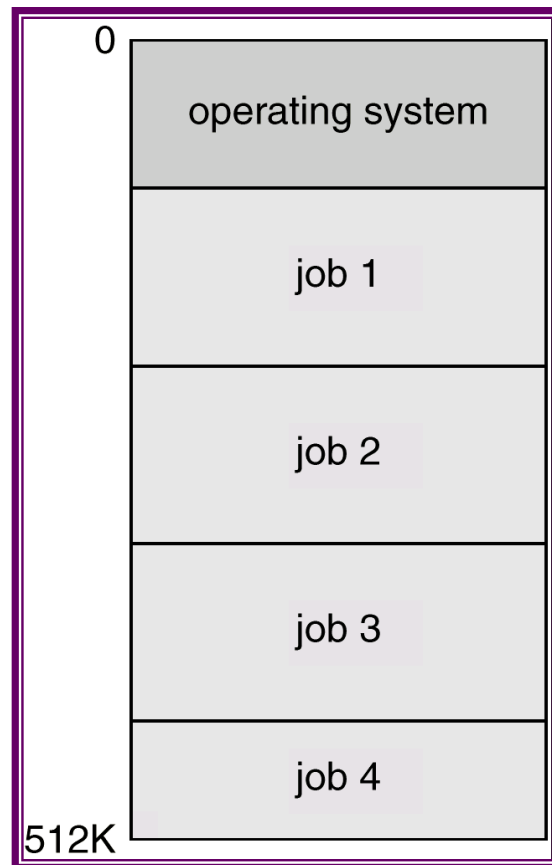


# Memory Layout for a Simple Batch System



# Multiprogrammed Batch Systems

Several jobs are kept in main memory at the same time, and the CPU is multiplexed among them.





# Computer-System Operation

- I/O devices and the CPU can execute concurrently.
- Each device controller is in charge of a particular device type.
- Each device controller has a local buffer.
- CPU moves data from/to main memory to/from local buffers
- I/O is from the device to local buffer of controller.
- Device controller informs CPU that it has finished its operation by causing an *interrupt*.

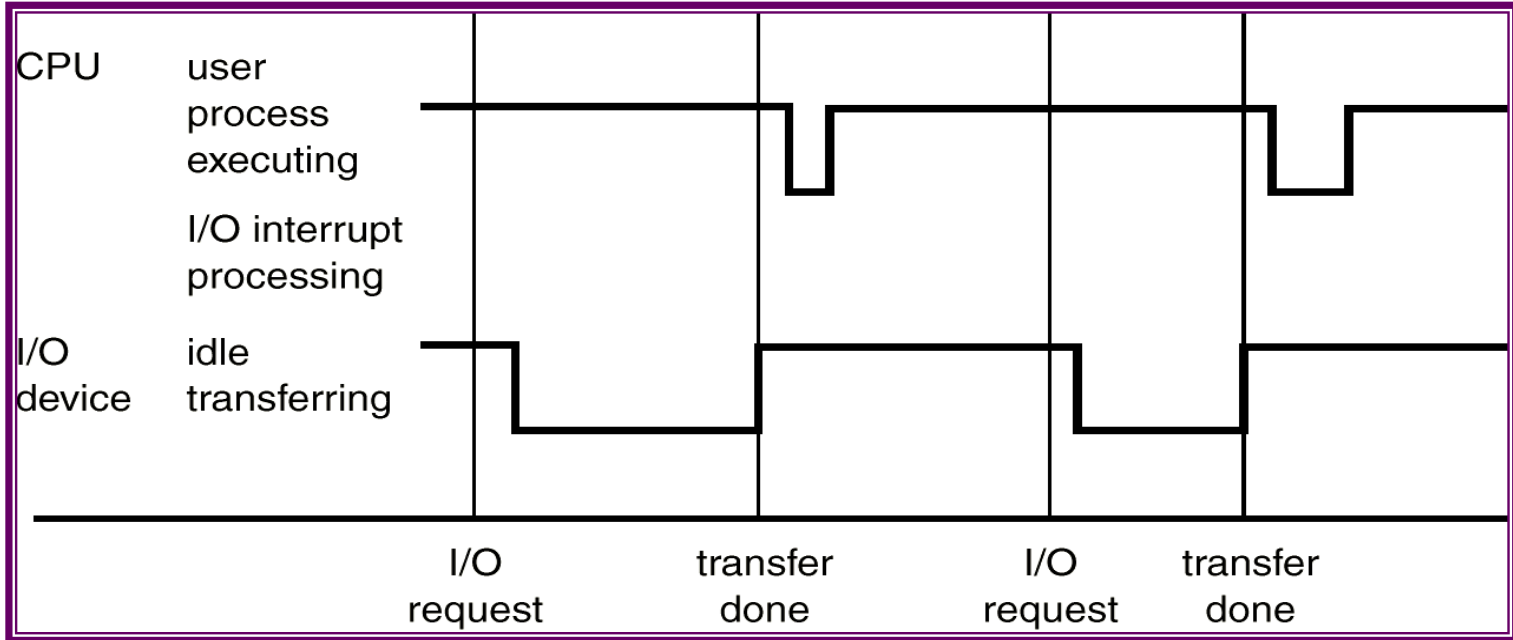
# Common Functions of Interrupts

- Interrupt transfers control to the interrupt service routine generally, through the *interrupt vector*, which contains the addresses of all the service routines.
- Interrupt architecture must save the address of the interrupted instruction.
- Incoming interrupts are *disabled* while another interrupt is being processed to prevent a *lost interrupt*.
- A *trap* is a software-generated interrupt caused either by an error or a user request.
- An operating system is *interrupt driven*.

# Interrupt Handling

- The operating system preserves the state of the CPU by storing registers and the program counter.
- Determines which type of interrupt has occurred:
  - ☞ *polling*
  - ☞ *vectored* interrupt system
- Separate segments of code determine what action should be taken for each type of interrupt

# Interrupt Time Line For a Single Process Doing Output



# OS Features Needed for Multiprogramming

- I/O routine supplied by the system.
- Memory management – the system must allocate the memory to several jobs.
- CPU scheduling – the system must choose among several jobs ready to run.
- Allocation of devices.

# Desktop Systems

- *Personal computers* – computer system dedicated to a single user.
- I/O devices – keyboards, mice, display screens, small printers.
- User convenience and responsiveness.
- Can adopt technology developed for larger operating system' often individuals have sole use of computer and do not need advanced CPU utilization or protection features.
- May run several different types of operating systems (Windows, MacOS, UNIX, Linux)

# Parallel Systems

- Multiprocessor systems with more than one CPU in close communication.
- *Tightly coupled system* – processors share memory and a clock; communication usually takes place through the shared memory.
- Advantages of parallel system:
  - ☞ Increased *throughput*
  - ☞ Economical
  - ☞ Increased reliability
    - 📄 graceful degradation
    - 📄 fail-soft systems

# Parallel Systems (Cont.)

- *Symmetric multiprocessing (SMP)*
  - ☞ Each processor runs an identical copy of the operating system.
  - ☞ Many processes can run at once without performance deterioration.
  - ☞ Most modern operating systems support SMP
- *Asymmetric multiprocessing*
  - ☞ Each processor is assigned a specific task; master processor schedules and allocates work to slave processors.
  - ☞ More common in extremely large systems



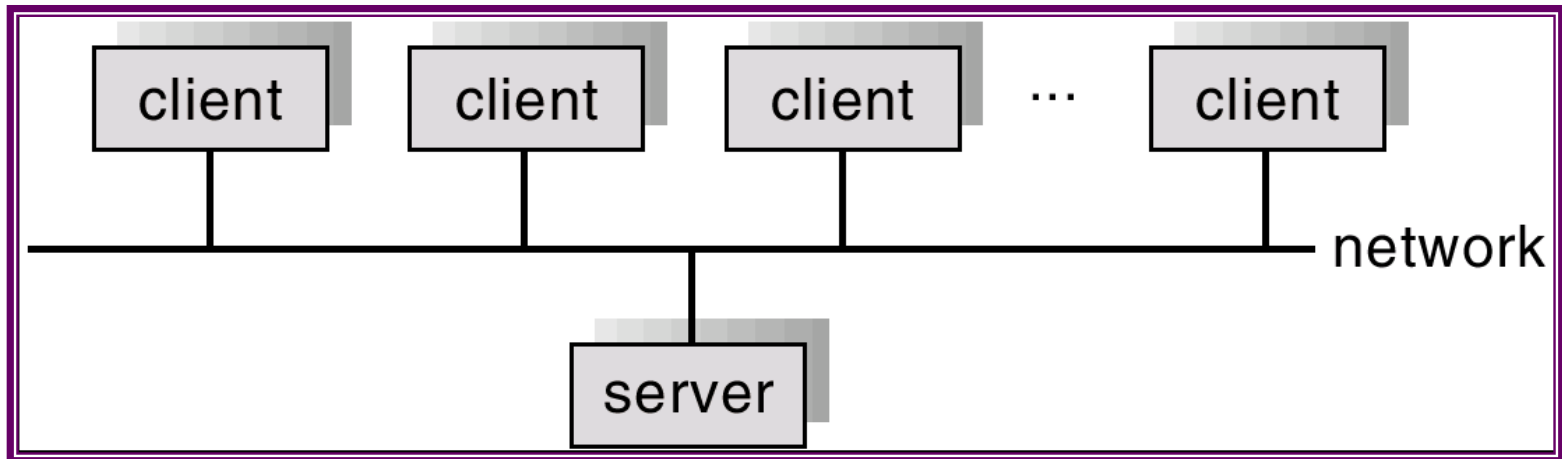
# Distributed Systems

- Distribute the computation among several physical processors.
- *Loosely coupled system* – each processor has its own local memory; processors communicate with one another through various communications lines, such as high-speed buses or telephone lines.
- Advantages of distributed systems.
  - ☞ Resources Sharing
  - ☞ Computation speed up – load sharing
  - ☞ Reliability
  - ☞ Communications

# Distributed Systems (cont)

- Requires networking infrastructure.
- Local area networks (LAN) or Wide area networks (WAN)
- May be either client-server or peer-to-peer systems.

# General Structure of Client-Server



# Real-Time Systems

- Often used as a control device in a dedicated application such as controlling scientific experiments, medical imaging systems, industrial control systems, and some display systems.
- Well-defined fixed-time constraints.
- Real-Time systems may be either *hard* or *soft* real-time.

# Real-Time Systems (Cont.)

## ■ Hard real-time:

- ☞ Secondary storage limited or absent, data stored in short term memory, or read-only memory (ROM)
- ☞ Conflicts with time-sharing systems, not supported by general-purpose operating systems.

## ■ Soft real-time

- ☞ Limited utility in industrial control of robotics
- ☞ Useful in applications (multimedia, virtual reality) requiring advanced operating-system features.

# Operating System Services

- Program execution – system capability to load a program into memory and to run it.
- I/O operations – since user programs cannot execute I/O operations directly, the operating system must provide some means to perform I/O.
- File-system manipulation – program capability to read, write, create, and delete files.
- Communications – exchange of information between processes executing either on the same computer or on different systems tied together by a network. Implemented via *shared memory* or *message passing*.
- Error detection – ensure correct computing by detecting errors in the CPU and memory hardware, in I/O devices, or in user programs.

# Additional Operating System Functions

Additional functions exist not for helping the user, but rather for ensuring efficient system operations.

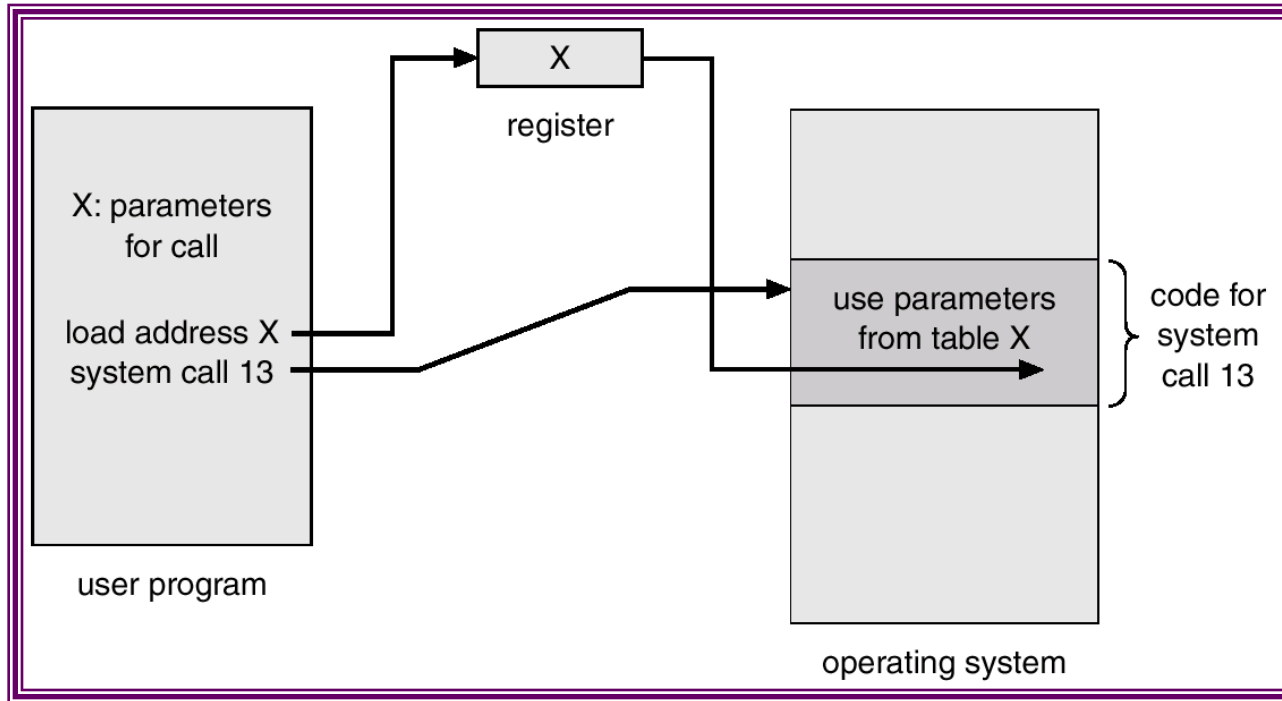
- Resource allocation – allocating resources to multiple users or multiple jobs running at the same time.
- Accounting – keep track of and record which users use how much and what kinds of computer resources for account billing or for accumulating usage statistics.
- Protection – ensuring that all access to system resources is controlled.

# System Calls

- System calls provide the interface between a running program and the operating system.
  - ☞ Generally available as assembly-language instructions.
  - ☞ Languages defined to replace assembly language for systems programming allow system calls to be made directly (e.g., C, C++)
- Three general methods are used to pass parameters between a running program and the operating system.
  - ☞ Pass parameters in *registers*.
  - ☞ Store the parameters in a table in memory, and the table address is passed as a parameter in a register.
  - ☞ *Push* (store) the parameters onto the *stack* by the program, and *pop* off the stack by operating system.



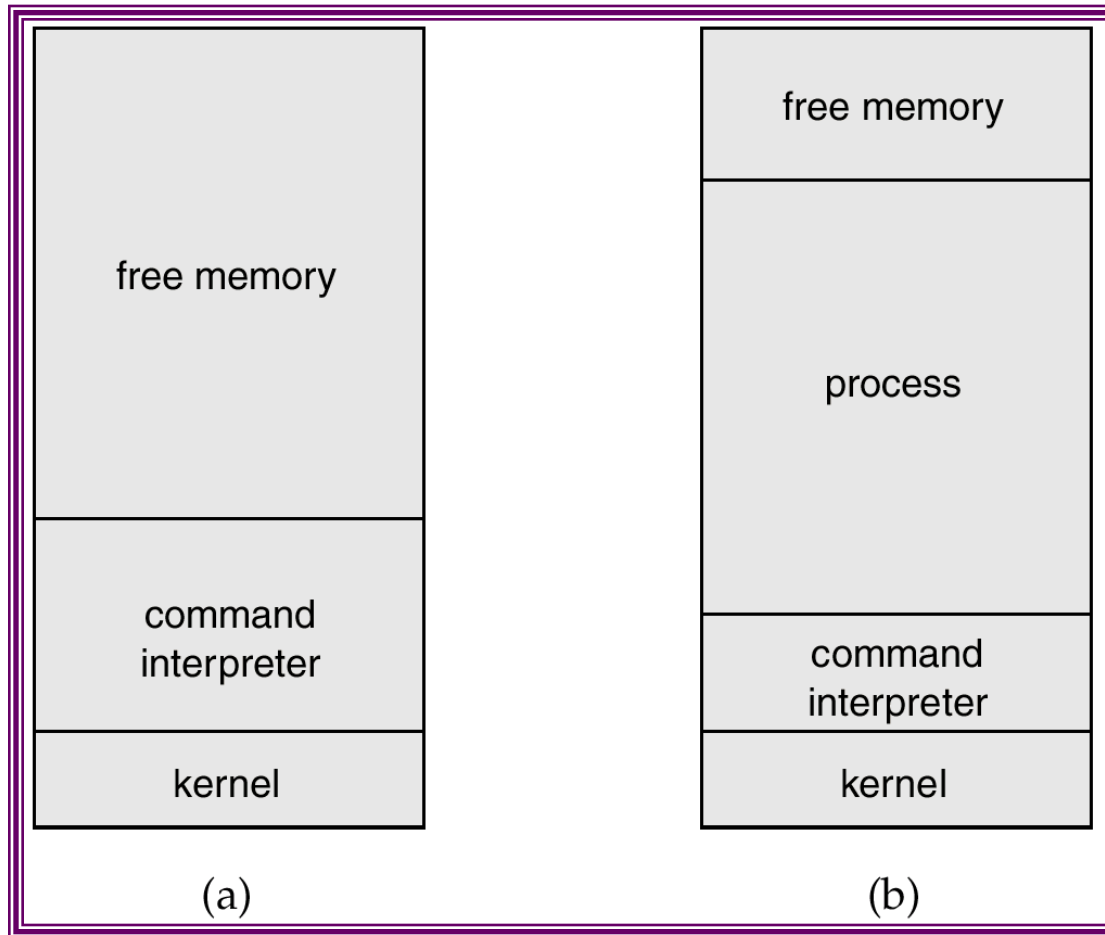
# Passing of Parameters As A Table



# Types of System Calls

- Process control
- File management
- Device management
- Information maintenance
- Communications

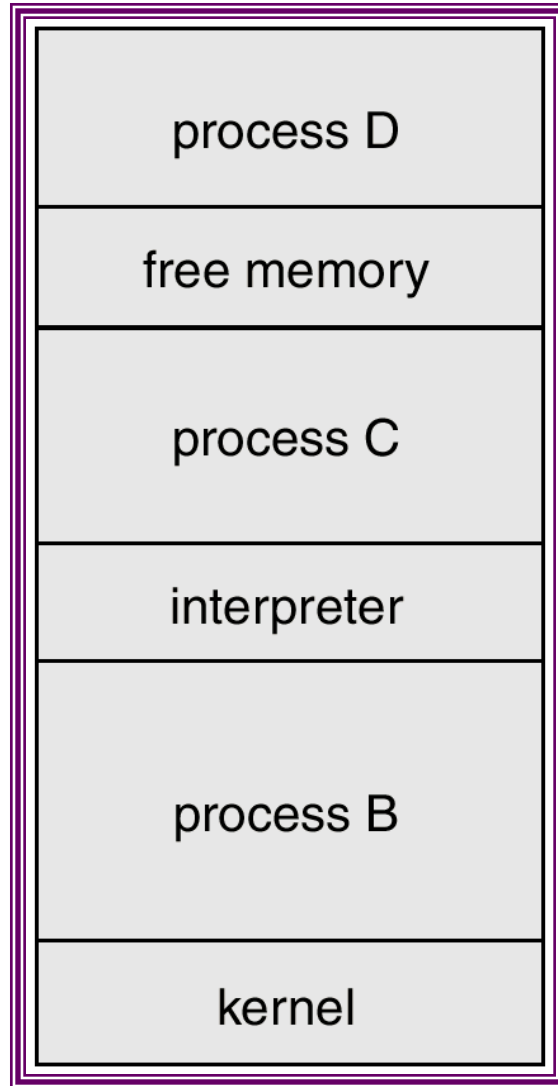
# MS-DOS Execution



At System Start-up

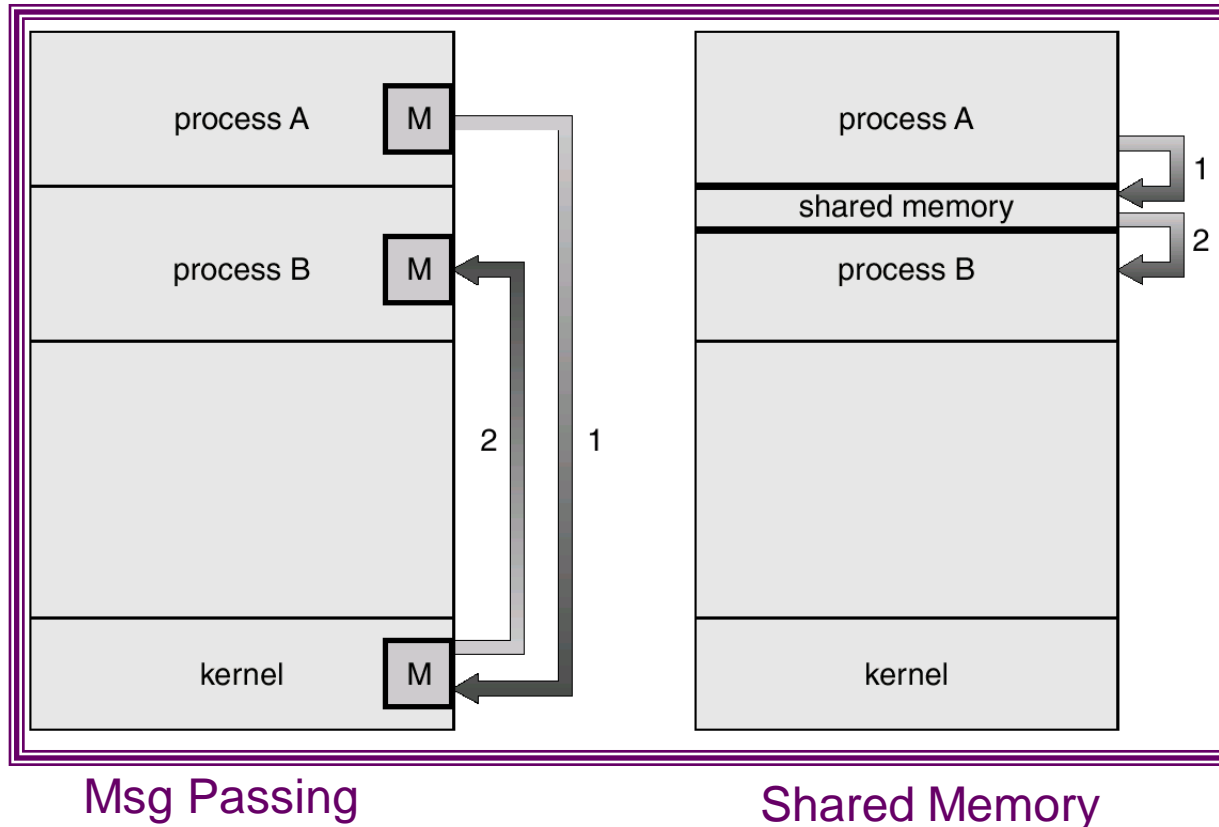
Running a Program

# UNIX Running Multiple Programs



# Communication Models

- Communication may take place using either message passing or shared memory.



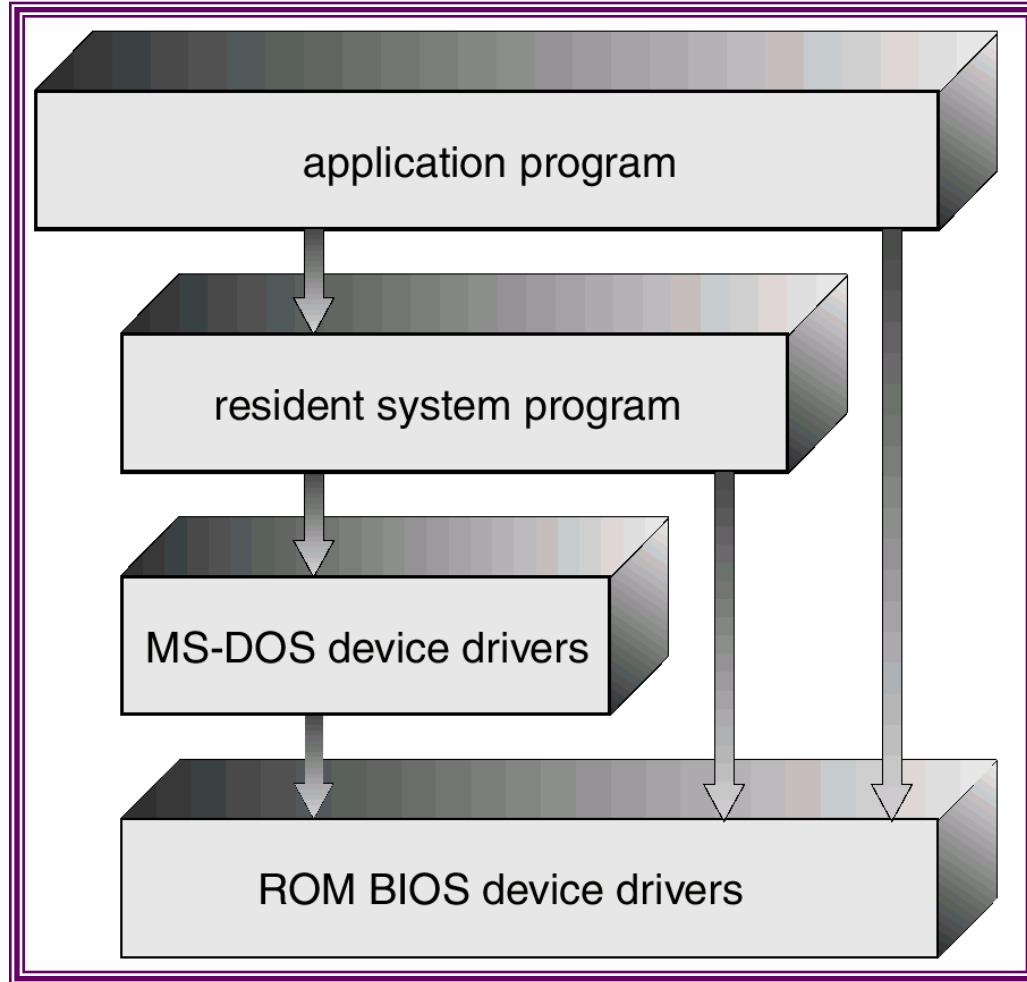
# System Programs

- 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.

# MS-DOS System Structure

- MS-DOS – written to provide the most functionality in the least space
  - ☞ not divided into modules
  - ☞ Although MS-DOS has some structure, its interfaces and levels of functionality are not well separated

# MS-DOS Layer Structure

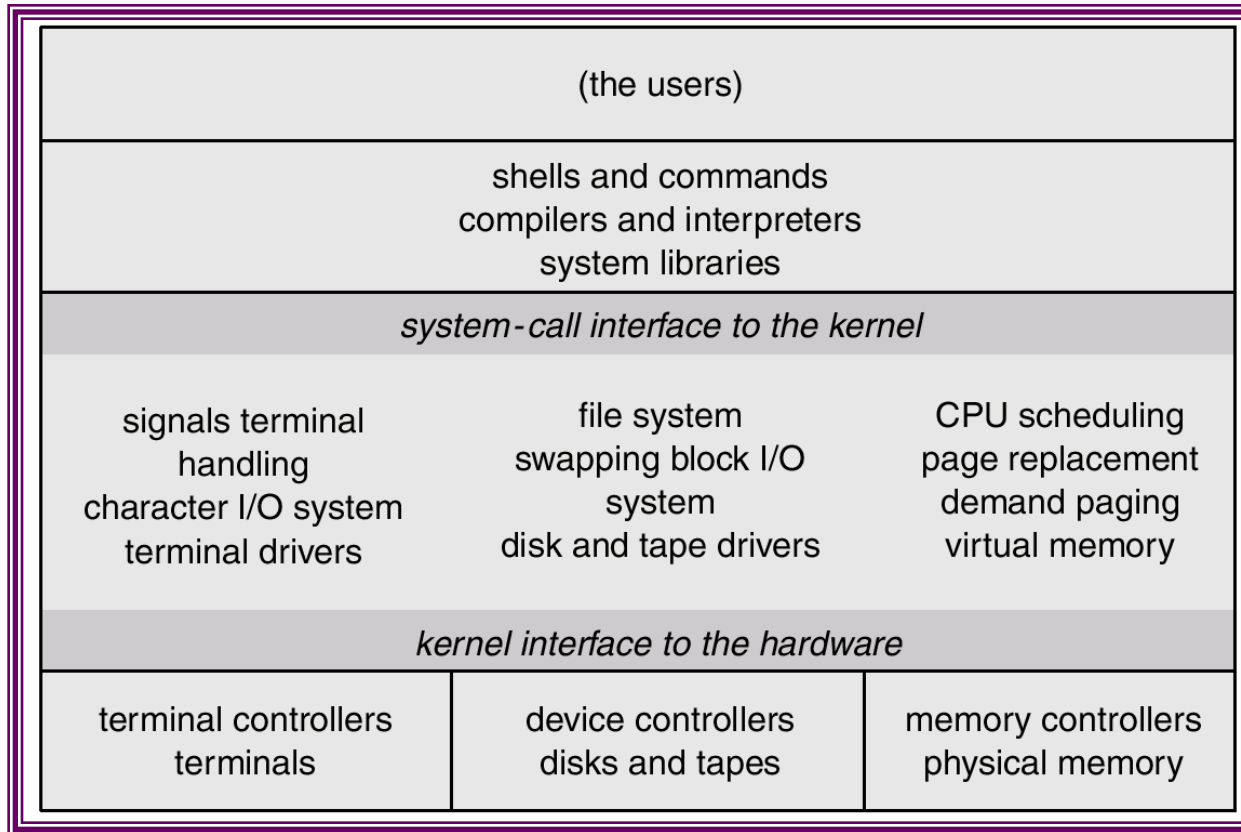




# UNIX System Structure

- UNIX – limited by hardware functionality, the original UNIX operating system had limited structuring. The UNIX OS consists of two separable parts.
  - ☞ Systems programs
  - ☞ The kernel
    - 📄 Consists of everything below the system-call interface and above the physical hardware
    - 📄 Provides the file system, CPU scheduling, memory management, and other operating-system functions; a large number of functions for one level.

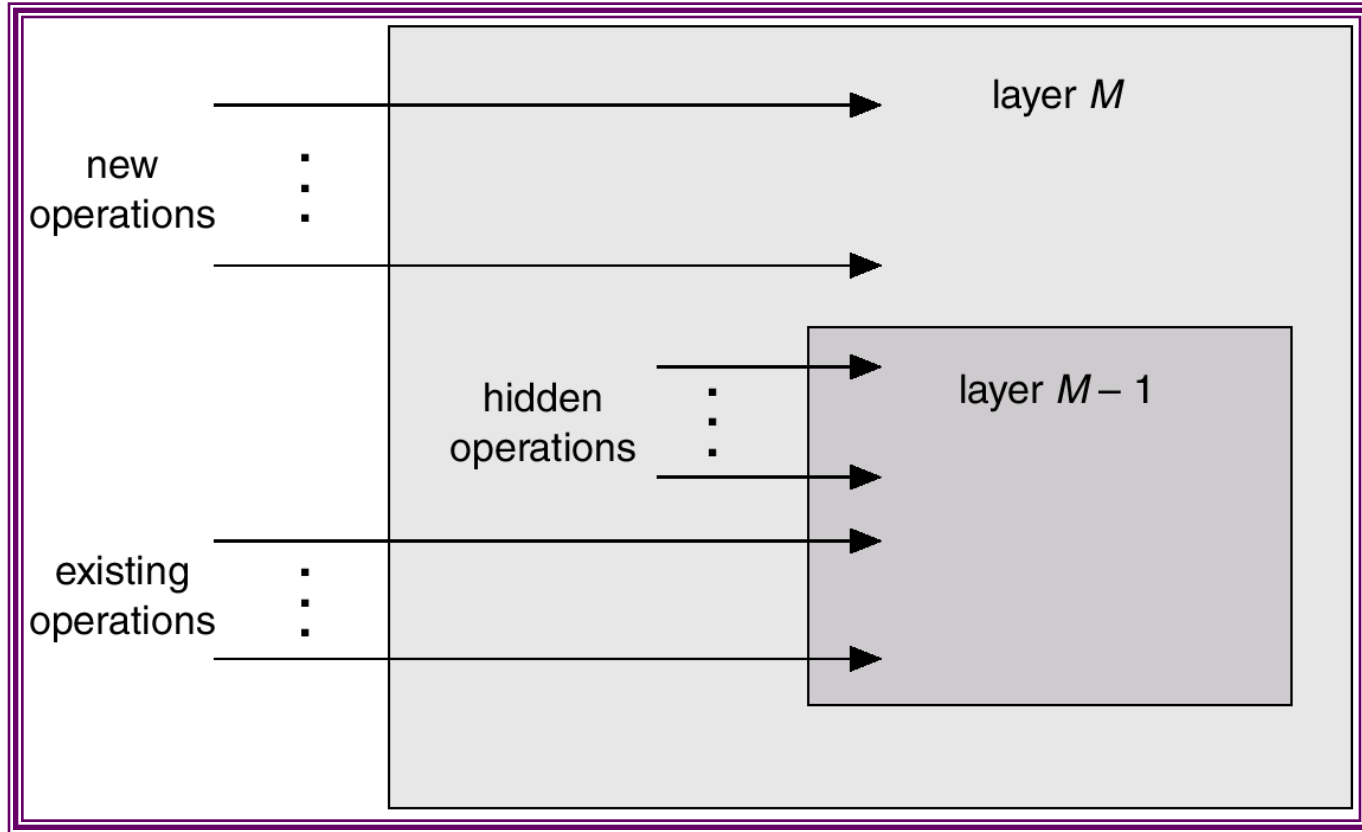
# UNIX System Structure



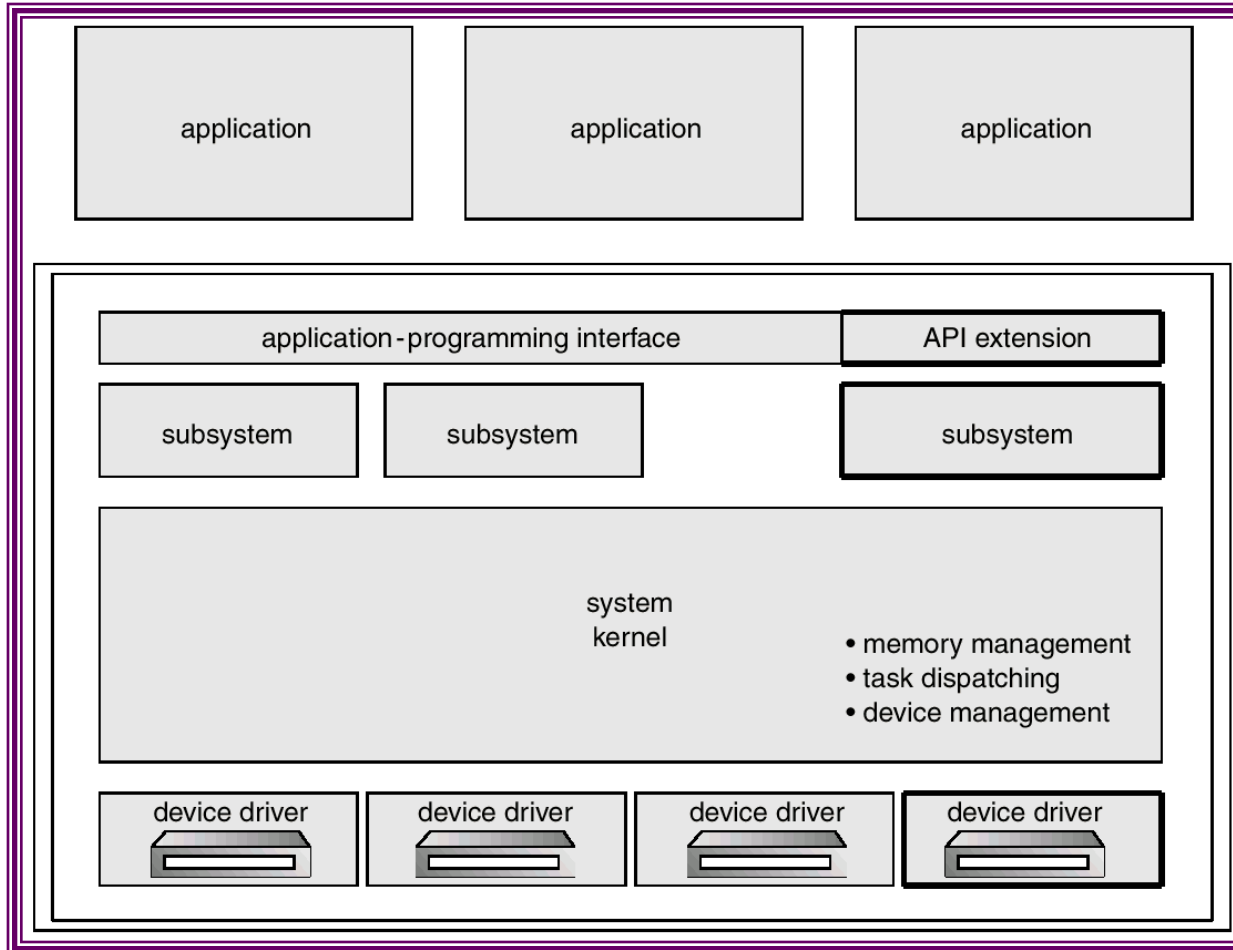
# Layered Approach

- The operating system is divided into a number of layers (levels), each built on top of lower layers. The bottom layer (layer 0), is the hardware; the highest (layer N) is the user interface.
- With modularity, layers are selected such that each uses functions (operations) and services of only lower-level layers.

# An Operating System Layer



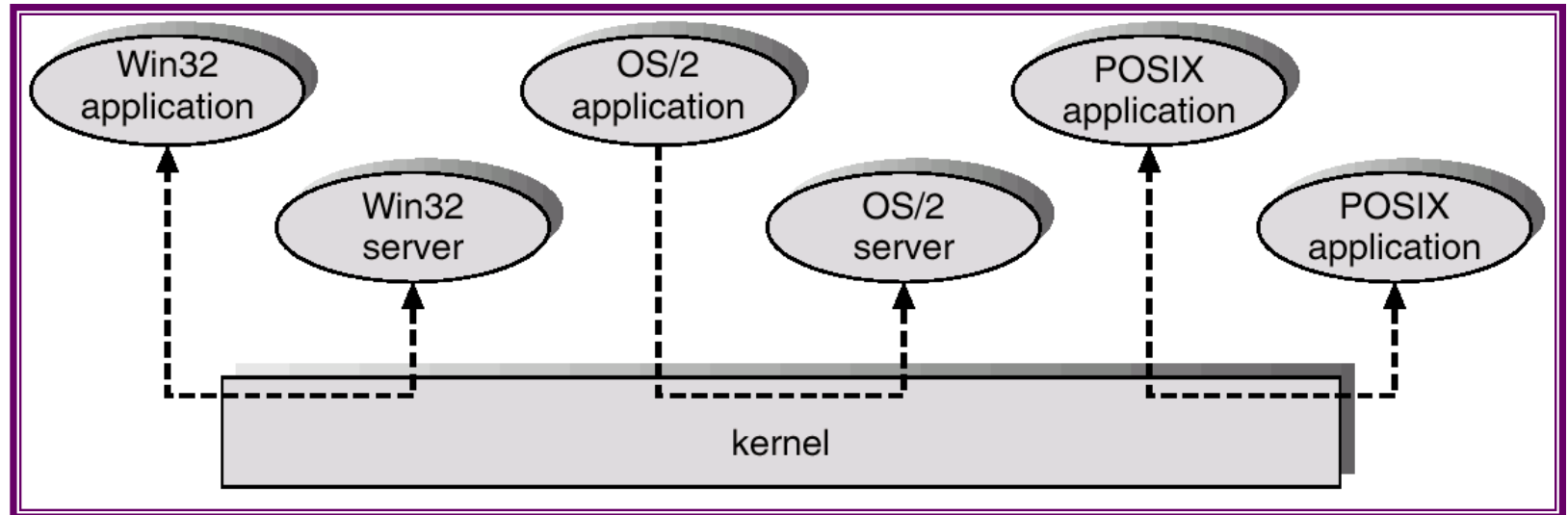
# OS/2 Layer Structure



# Microkernel System Structure

- Moves as much from the kernel into “*user*” space.
- Communication takes place between user modules using message passing.
- Benefits:
  - easier to extend a microkernel
  - easier to port the operating system to new architectures
  - more reliable (less code is running in kernel mode)
  - more secure

# Windows NT Client-Server Structure



# Virtual Machines

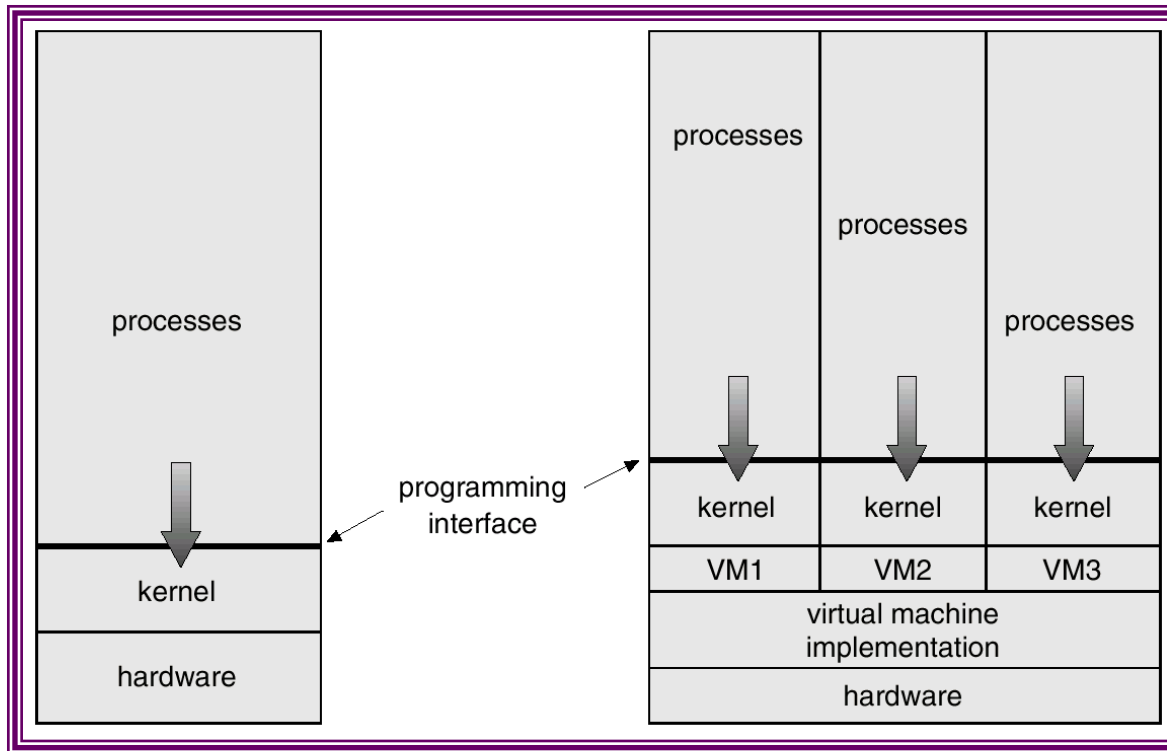
- A *virtual machine* takes the layered approach to its logical conclusion. It treats hardware and the operating system kernel as though they were all hardware.
- A virtual machine provides an interface *identical* to the underlying bare hardware.
- The operating system creates the illusion of multiple processes, each executing on its own processor with its own (virtual) memory.



# Virtual Machines (Cont.)

- The resources of the physical computer are shared to create the virtual machines.
  - ☞ CPU scheduling can create the appearance that users have their own processor.
  - ☞ Spooling and a file system can provide virtual card readers and virtual line printers.
  - ☞ A normal user time-sharing terminal serves as the virtual machine operator's console.

# System Models



Non-virtual Machine

Virtual Machine

# Advantages/Disadvantages of Virtual Machines

- The virtual-machine concept provides complete protection of system resources since each virtual machine is isolated from all other virtual machines. This isolation, however, permits no direct sharing of resources.
- A virtual-machine system is a perfect vehicle for operating-systems research and development. System development is done on the virtual machine, instead of on a physical machine and so does not disrupt normal system operation.
- The virtual machine concept is difficult to implement due to the effort required to provide an *exact* duplicate to the underlying machine.

# System Design Goals

- User goals – operating system should be convenient to use, easy to learn, reliable, safe, and fast.
- System goals – operating system should be easy to design, implement, and maintain, as well as flexible, reliable, error-free, and efficient.

# System Implementation

- Traditionally written in assembly language, operating systems can now be written in higher-level languages.
- Code written in a high-level language:
  - ☞ can be written faster.
  - ☞ is more compact.
  - ☞ is easier to understand and debug.
- An operating system is far easier to *port* (move to some other hardware) if it is written in a high-level language.

# Unit 2: Processes and CPU Scheduling

- Process Concept
- Process Scheduling
- Operations on Processes
- Cooperating Processes
- Interprocess Communication
- Scheduling Criteria
- Scheduling algorithms
- Multiple-Processor Scheduling
- Real-Time Scheduling,
- Thread Scheduling

# Process Concept

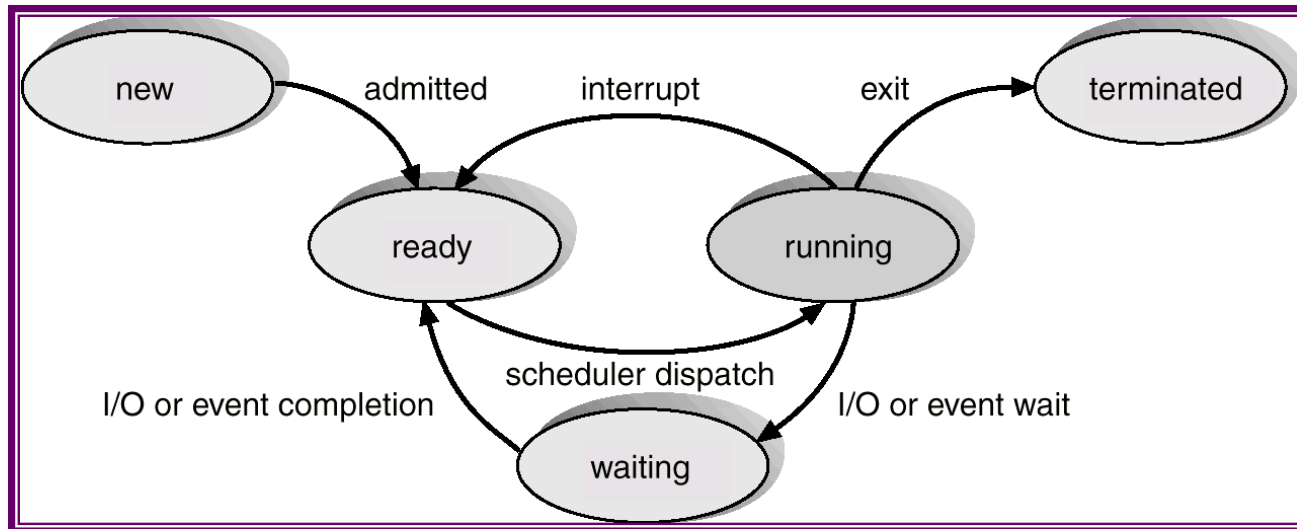
- An operating system executes a variety of programs:
  - ☞ Batch system – jobs
  - ☞ Time-shared systems – user programs or tasks
- Textbook uses the terms *job* and *process* almost interchangeably.
- Process – a program in execution; process execution must progress in sequential fashion.
- A process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

# 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 process.
  - **terminated**: The process has finished execution.



# Diagram of Process State

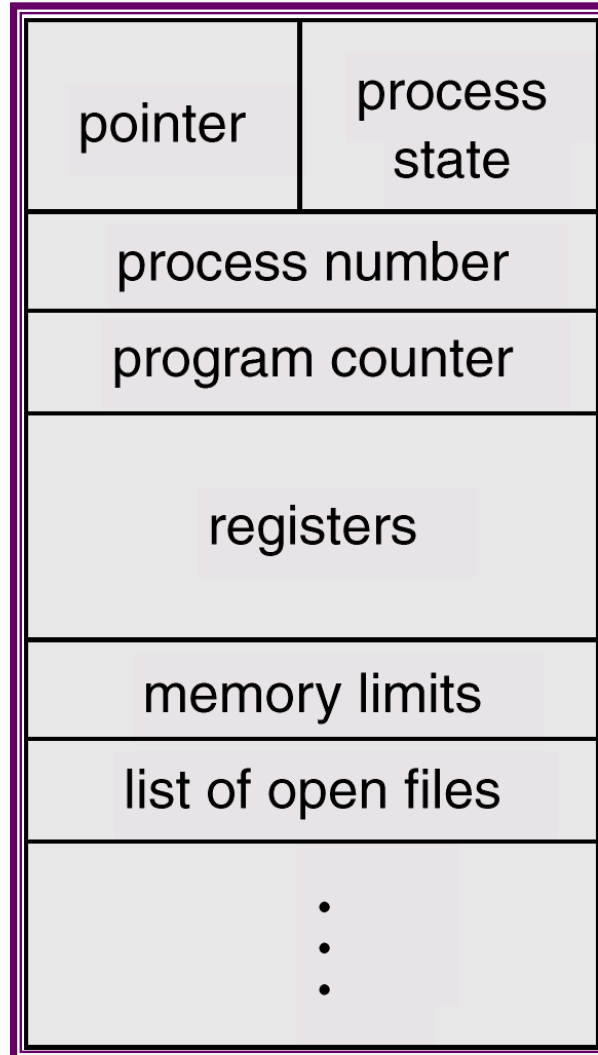


# Process Control Block (PCB)

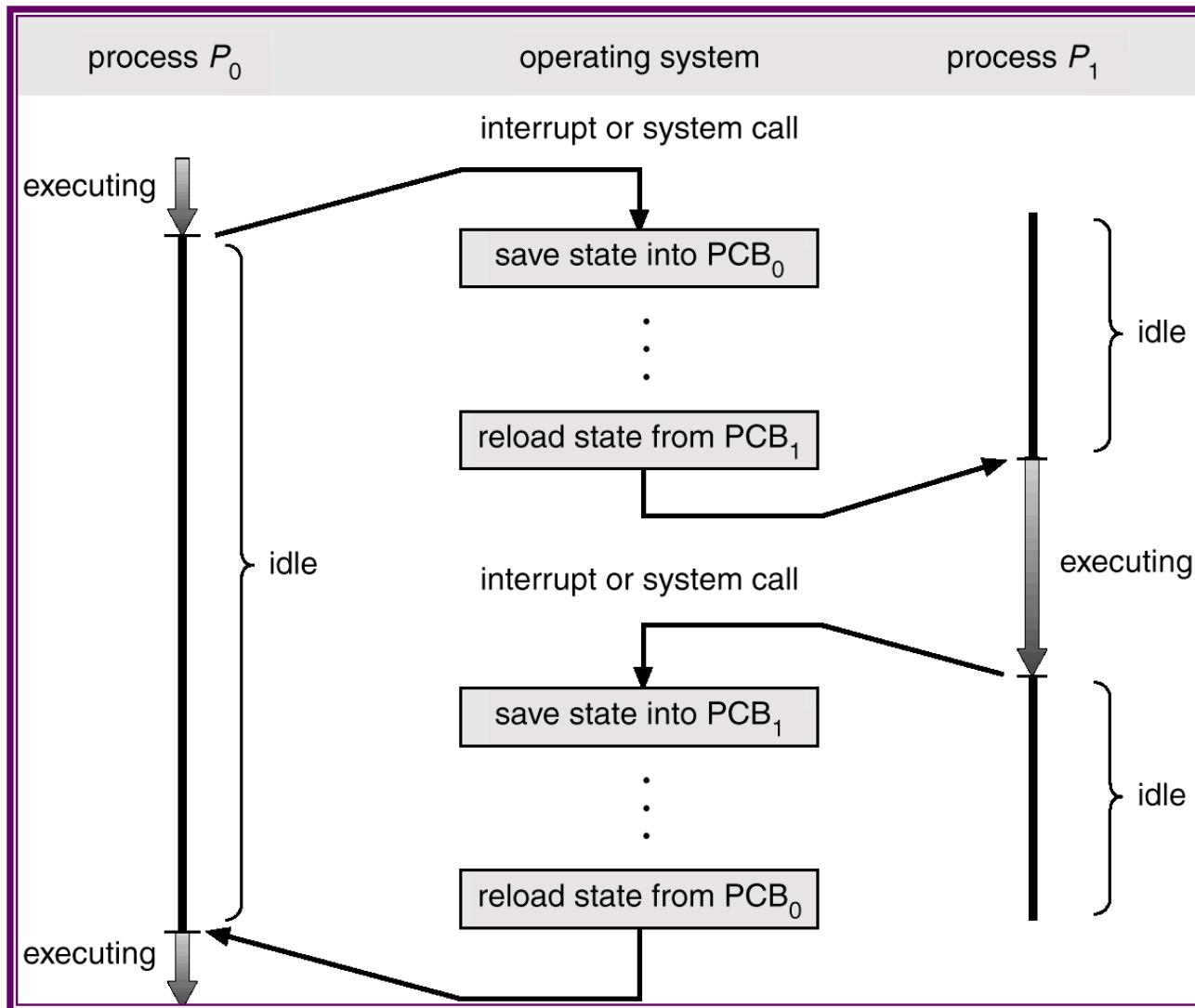
Information associated with each process.

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

# Process Control Block (PCB)



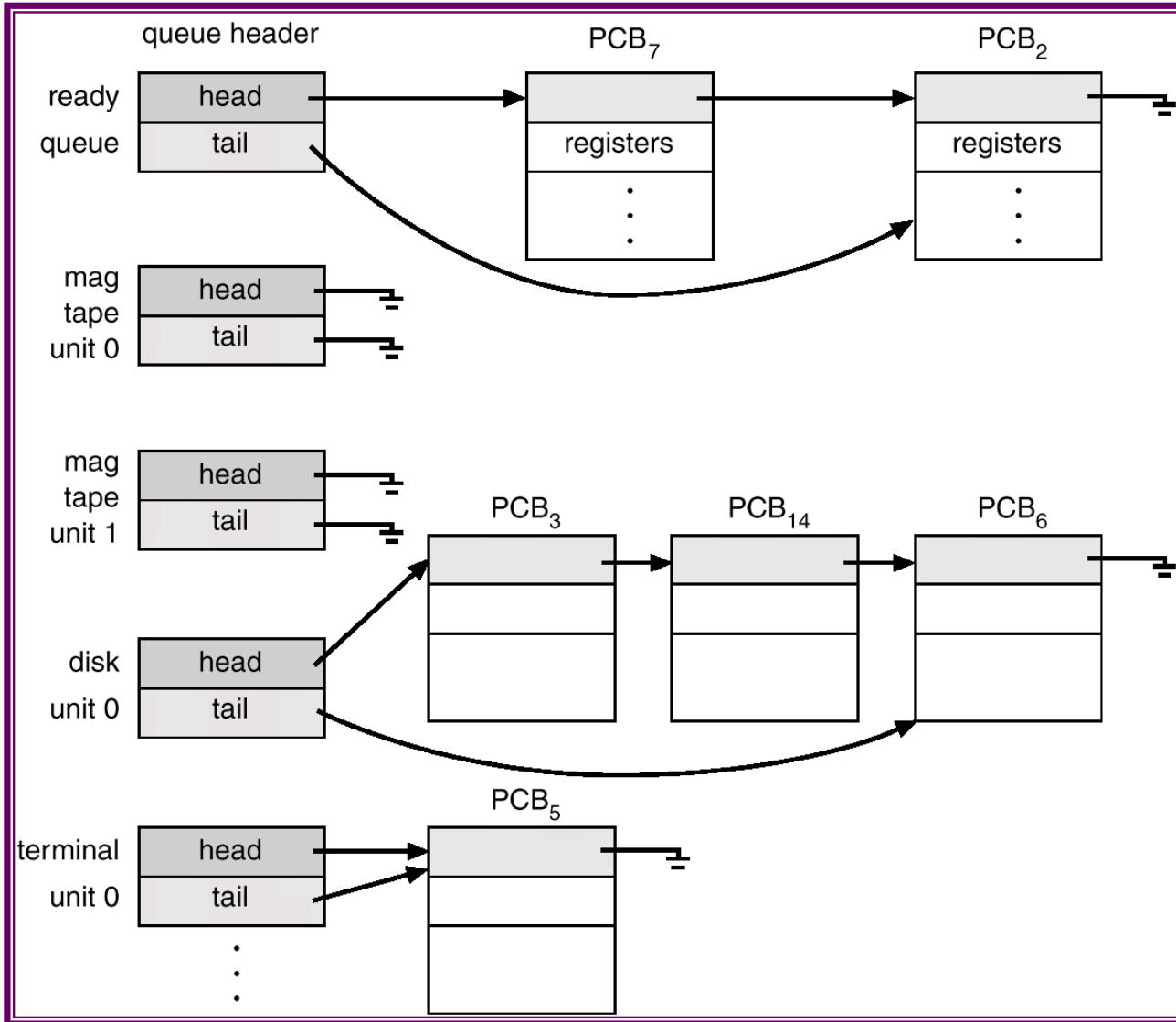
# CPU Switch From Process to Process



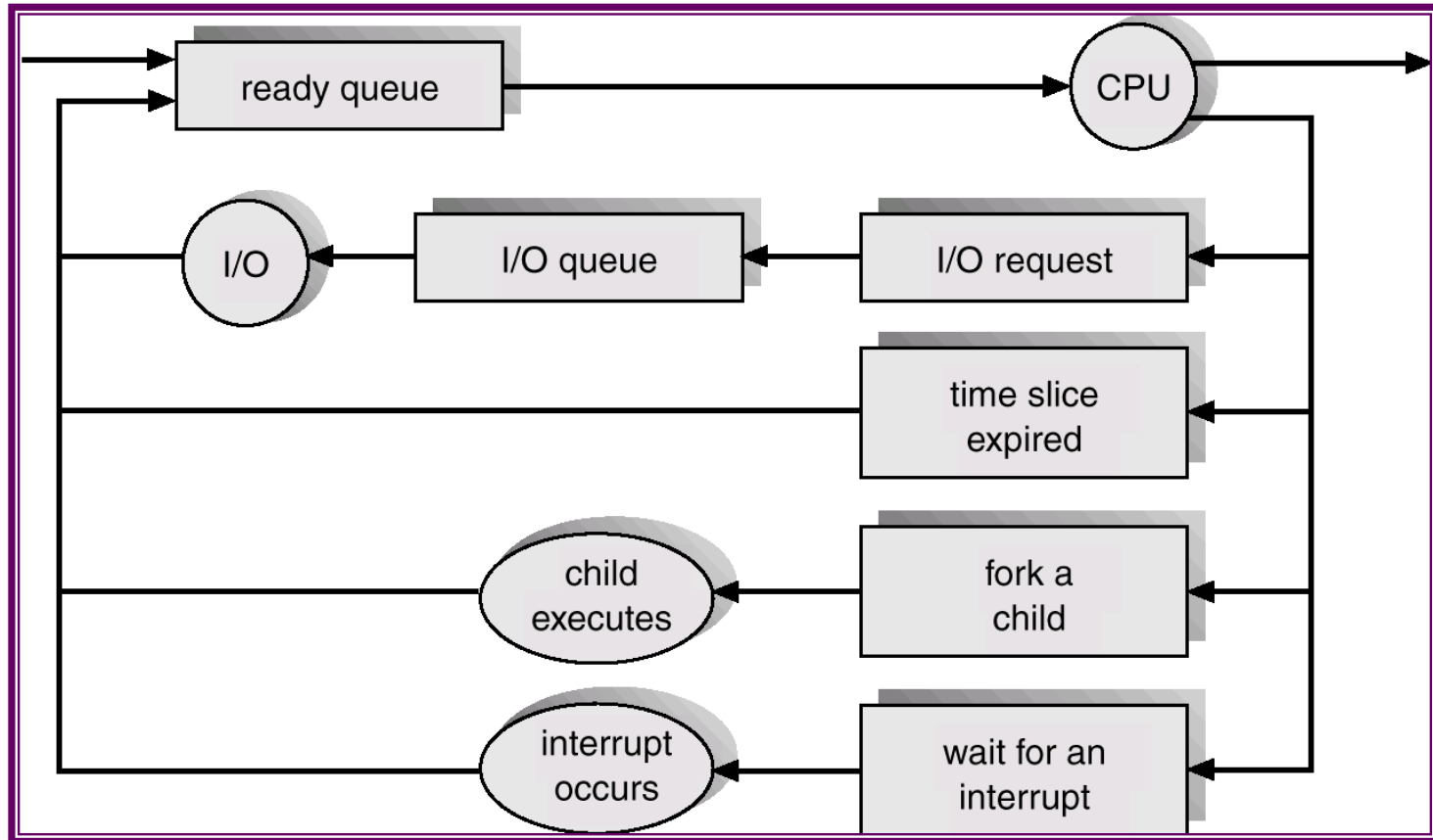
# Process Scheduling Queues

- 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.
- Process migration between the various queues.

# Ready Queue And Various I/O Device Queues



# Representation of Process Scheduling

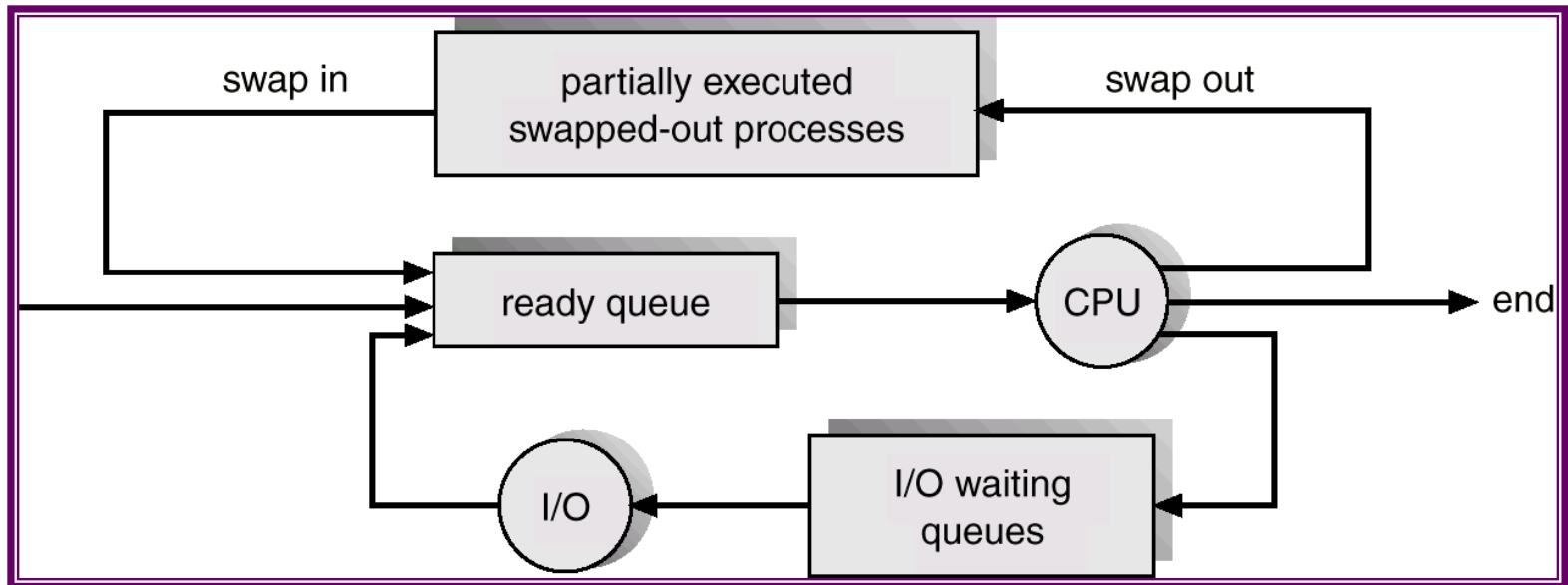


# Schedulers

- 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.



# Addition of Medium Term Scheduling



# Schedulers (Cont.)

- 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.

# Context Switch

- When CPU switches to another process, the system must save the state of the old process and load the saved state for the new process.
- Context-switch time is overhead; the system does no useful work while switching.
- Time dependent on hardware support.

# Process Creation

- Parent process create children processes, which, in turn create other processes, forming a tree of processes.
- Resource sharing
  - ☞ Parent and children share all resources.
  - ☞ Children share subset of parent's resources.
  - ☞ Parent and child share no resources.
- Execution
  - ☞ Parent and children execute concurrently.
  - ☞ Parent waits until children terminate.

# Process Creation (Cont.)

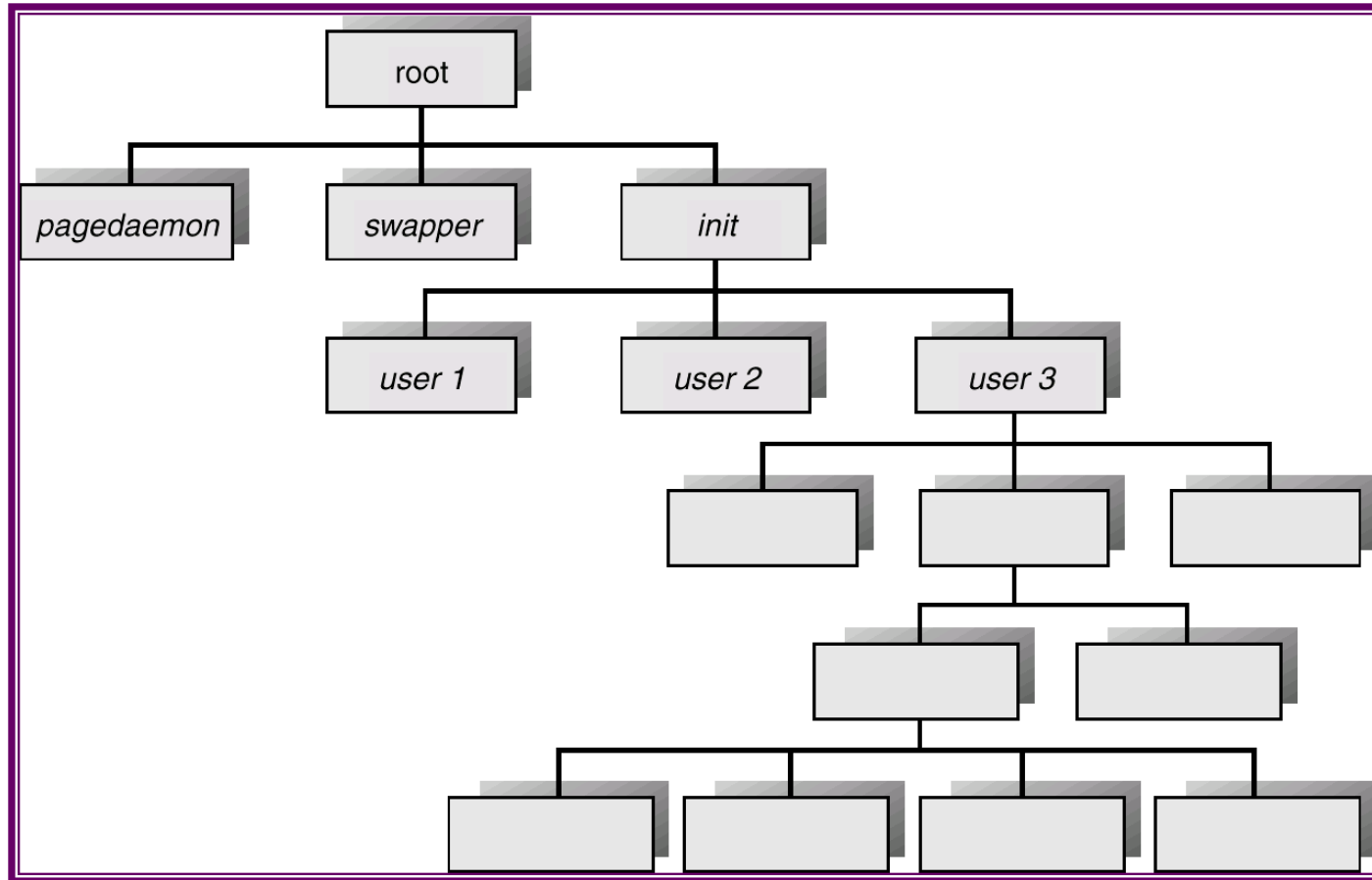
## ■ Address space

- ☞ Child duplicate of parent.
- ☞ Child has a program loaded into it.

## ■ UNIX examples

- ☞ **fork** system call creates new process
- ☞ **exec** system call used after a **fork** to replace the process' memory space with a new program.

# Processes Tree on a UNIX System



# Process Termination

- Process executes last statement and asks the operating system to decide it (**exit**).
  - ☞ Output data from child to parent (via **wait**).
  - ☞ Process' resources are deallocated by operating system.
- Parent may terminate execution of children processes (**abort**).
  - ☞ Child has exceeded allocated resources.
  - ☞ Task assigned to child is no longer required.
  - ☞ Parent is exiting.
    - 📄 Operating system does not allow child to continue if its parent terminates.
    - 📄 Cascading termination.

# Cooperating Processes

- *Independent* process cannot affect or be affected by the execution of another process.
- *Cooperating* process can affect or be affected by the execution of another process
- Advantages of process cooperation
  - ☞ Information sharing
  - ☞ Computation speed-up
  - ☞ Modularity
  - ☞ Convenience



# Producer-Consumer Problem

- Paradigm for cooperating processes, *producer* process produces information that is consumed by a *consumer* process.
  - ☞ *unbounded-buffer* places no practical limit on the size of the buffer.
  - ☞ *bounded-buffer* assumes that there is a fixed buffer size.

# Bounded-Buffer – Shared-Memory Solution

- Shared data

```
#define BUFFER_SIZE 10
typedef struct {
    . . .
} item;
item buffer[BUFFER_SIZE];
int in = 0;
int out = 0;
```

- Solution is correct, but can only use `BUFFER_SIZE-1` elements

# Bounded-Buffer – Producer Process

```
item nextProduced;
```

```
while (1) {  
    while (((in + 1) % BUFFER_SIZE) == out)  
        ; /* do nothing */  
    buffer[in] = nextProduced;  
    in = (in + 1) % BUFFER_SIZE;  
}
```

# Bounded-Buffer – Consumer Process

```
item nextConsumed;
```

```
while (1) {  
    while (in == out)  
        ; /* do nothing */  
    nextConsumed = buffer[out];  
    out = (out + 1) % BUFFER_SIZE;  
}
```

# Interprocess Communication (IPC)

- Mechanism for processes to communicate and to synchronize their actions.
- Message system – processes communicate with each other without resorting to shared variables.
- IPC facility provides two operations:
  - ☞ **send**(*message*) – message size fixed or variable
  - ☞ **receive**(*message*)
- If  $P$  and  $Q$  wish to communicate, they need to:
  - ☞ establish a *communication link* between them
  - ☞ exchange messages via send/receive
- Implementation of communication link
  - ☞ physical (e.g., shared memory, hardware bus)
  - ☞ logical (e.g., logical properties)

# Implementation Questions

- How are links established?
- Can a link be associated with more than two processes?
- How many links can there be between every pair of communicating processes?
- What is the capacity of a link?
- Is the size of a message that the link can accommodate fixed or variable?
- Is a link unidirectional or bi-directional?

# Direct Communication

- Processes must name each other explicitly:
  - ☞ **send** ( $P$ , *message*) – send a message to process  $P$
  - ☞ **receive**( $Q$ , *message*) – receive a message from process  $Q$
- Properties of communication link
  - ☞ Links are established automatically.
  - ☞ A link is associated with exactly one pair of communicating processes.
  - ☞ Between each pair there exists exactly one link.
  - ☞ The link may be unidirectional, but is usually bi-directional.

# Indirect Communication

- Messages are directed and received from mailboxes (also referred to as ports).
  - ☞ Each mailbox has a unique id.
  - ☞ Processes can communicate only if they share a mailbox.
- Properties of communication link
  - ☞ Link established only if processes share a common mailbox
  - ☞ A link may be associated with many processes.
  - ☞ Each pair of processes may share several communication links.
  - ☞ Link may be unidirectional or bi-directional.



# Indirect Communication

## ■ Operations

- ☞ create a new mailbox
- ☞ send and receive messages through mailbox
- ☞ destroy a mailbox

## ■ Primitives are defined as:

**send**(*A, message*) – send a message to mailbox A

**receive**(*A, message*) – receive a message from mailbox A

# Indirect Communication

## ■ Mailbox sharing

- ☞  $P_1$ ,  $P_2$ , and  $P_3$  share mailbox A.
- ☞  $P_1$  sends;  $P_2$  and  $P_3$  receive.
- ☞ Who gets the message?

## ■ Solutions

- ☞ Allow a link to be associated with at most two processes.
- ☞ Allow only one process at a time to execute a receive operation.
- ☞ Allow the system to select arbitrarily the receiver. Sender is notified who the receiver was.

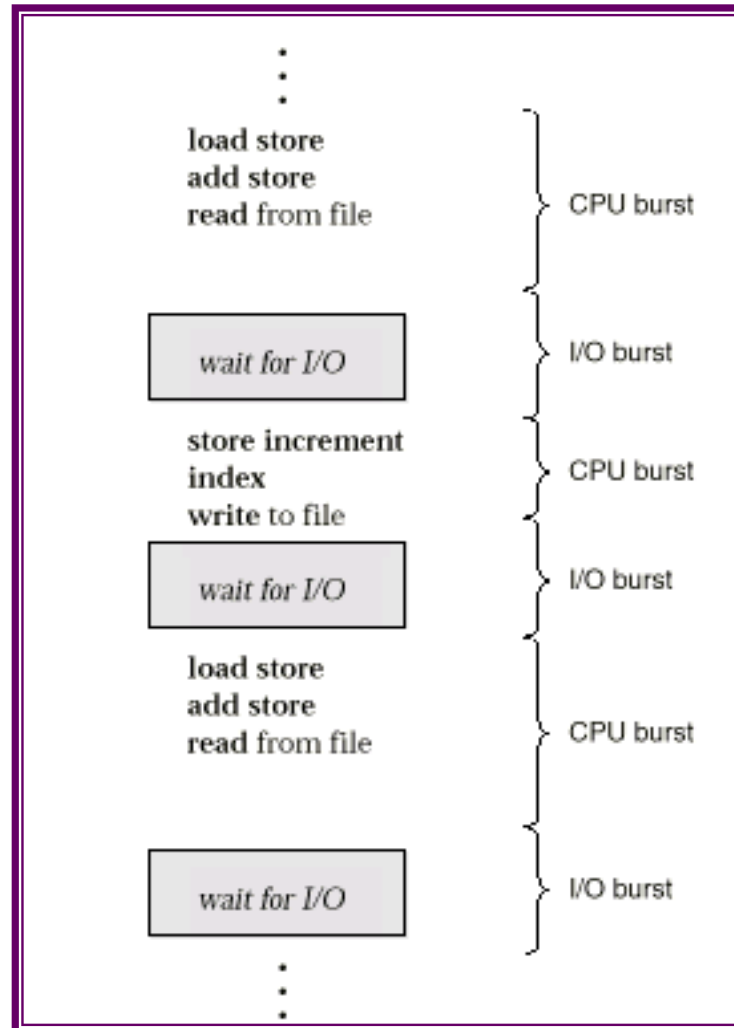
# Synchronization

- Message passing may be either blocking or non-blocking.
- **Blocking** is considered **synchronous**
- **Non-blocking** is considered **asynchronous**
- **send** and **receive** primitives may be either blocking or non-blocking.

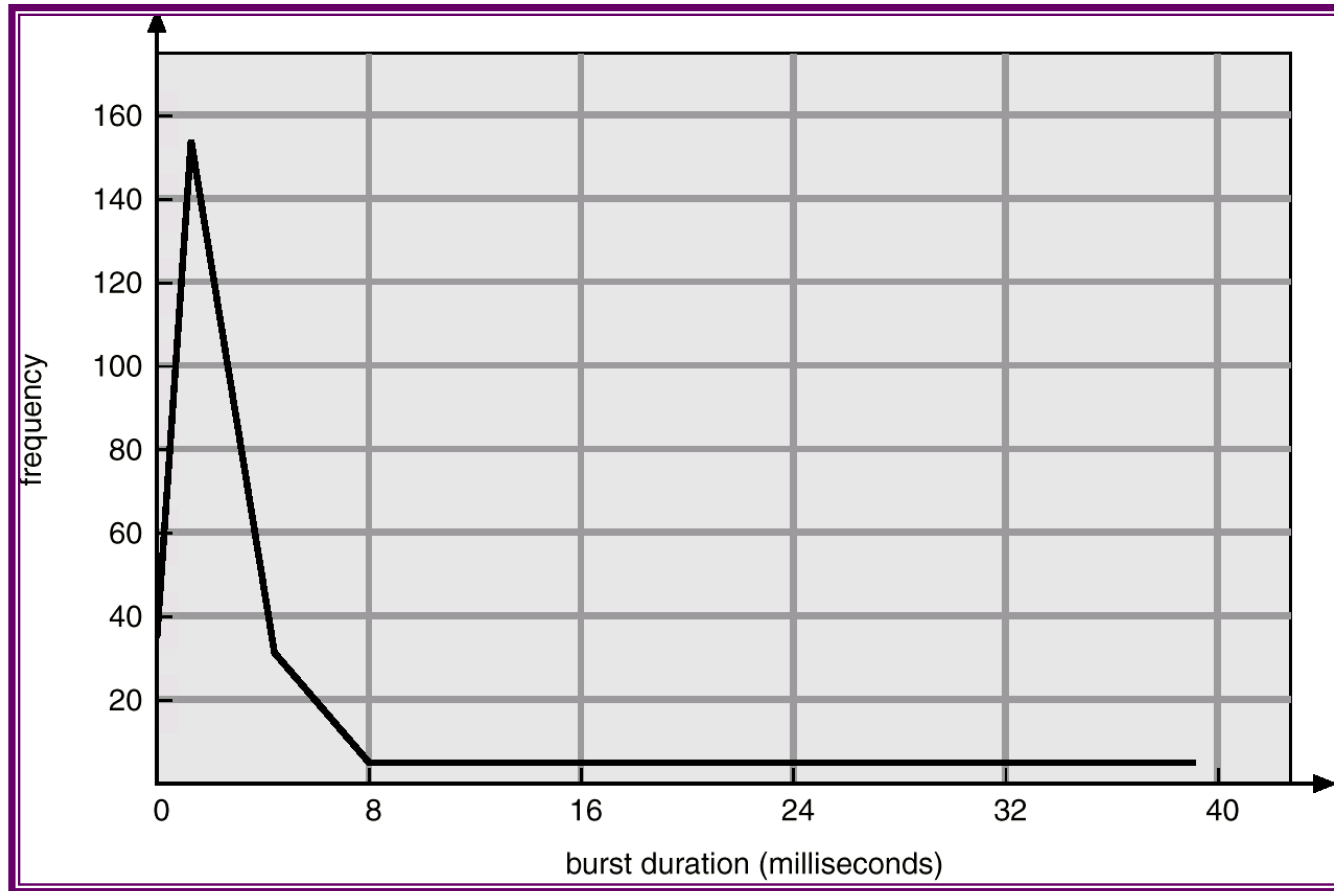
# Basic Concepts

- Maximum CPU utilization obtained with multiprogramming
- CPU–I/O Burst Cycle – Process execution consists of a *cycle* of CPU execution and I/O wait.
- CPU burst distribution

# Alternating Sequence of CPU And I/O Bursts



# Histogram of CPU-burst Times



# CPU Scheduler

- Selects from among the processes in memory that are ready to execute, and allocates the CPU to one of them.
- CPU scheduling decisions may take place when a process:
  1. Switches from running to waiting state.
  2. Switches from running to ready state.
  3. Switches from waiting to ready.
  4. Terminates.
- Scheduling under 1 and 4 is *nonpreemptive*.
- All other scheduling is *preemptive*.

# Dispatcher

- Dispatcher module gives control of the CPU to the process selected by the short-term scheduler; this involves:
  - ☞ switching context
  - ☞ switching to user mode
  - ☞ jumping to the proper location in the user program to restart that program
- *Dispatch latency* – time it takes for the dispatcher to stop one process and start another running.



# Scheduling Criteria

- CPU utilization – keep the CPU as busy as possible
- Throughput – # of processes that complete their execution per time unit
- Turnaround time – amount of time to execute a particular process
- Waiting time – amount of time a process has been waiting in the ready queue
- 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)

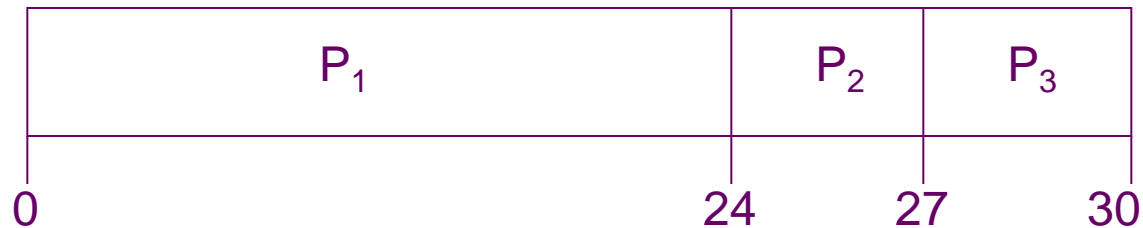
# Optimization Criteria

- Max CPU utilization
- Max throughput
- Min turnaround time
- Min waiting time
- 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:

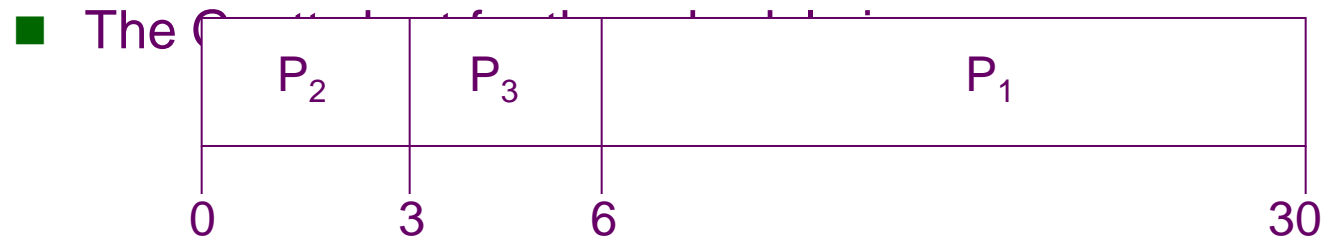


- Waiting time for  $P_1 = 0$ ;  $P_2 = 24$ ;  $P_3 = 27$
- Average waiting time:  $(0 + 24 + 27)/3 = 17$

# FCFS Scheduling (Cont.)

Suppose that the processes arrive in the order

$P_2, P_3, P_1$ .



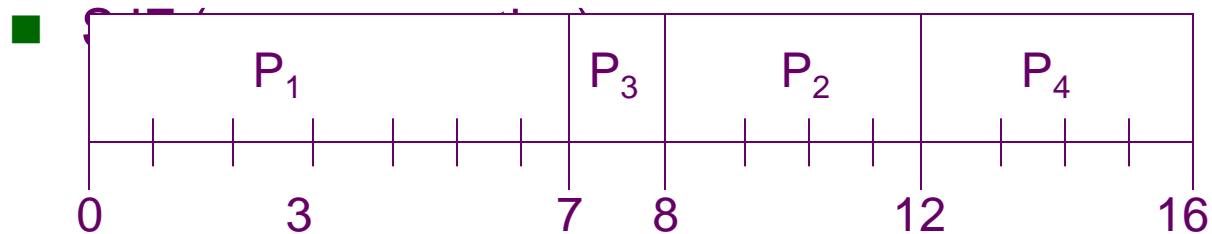
- Waiting time for  $P_1 = 6$ ;  $P_2 = 0$ ;  $P_3 = 3$
- Average waiting time:  $(6 + 0 + 3)/3 = 3$
- Much better than previous case.
- *Convoy effect* short process behind long process

# Shortest-Job-First (SJF) Scheduling

- Associate with each process the length of its next CPU burst. Use these lengths to schedule the process with the shortest time.
- Two schemes:
  - ☞ nonpreemptive – once CPU given to the process it cannot be preempted until completes its CPU burst.
  - ☞ preemptive – if a new process arrives with CPU burst length less than remaining time of current executing process, preempt. This scheme is known as the Shortest-Remaining-Time-First (SRTF).
- SJF is optimal – gives minimum average waiting time for a given set of processes.

# Example of Non-Preemptive SJF

<u>Process</u>	<u>Arrival Time</u>	<u>Burst Time</u>
$P_1$	0.0	7
$P_2$	2.0	4
$P_3$	4.0	1
$P_4$	5.0	4

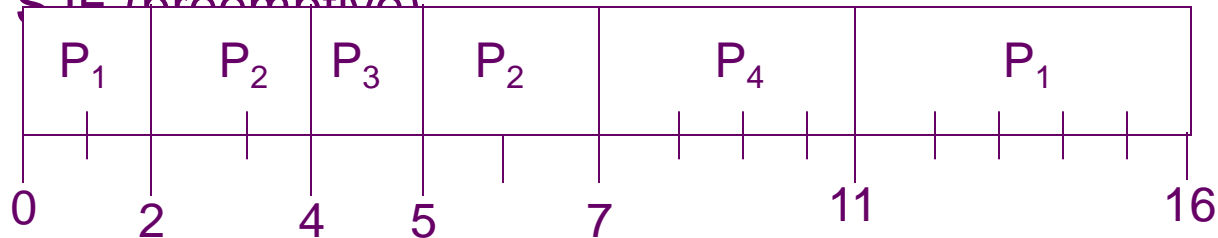


■ Average waiting time =  $(0 + 6 + 3 + 7)/4 - 4$

# Example of Preemptive SJF

<u>Process</u>	<u>Arrival Time</u>	<u>Burst Time</u>
$P_1$	0.0	7
$P_2$	2.0	4
$P_3$	4.0	1
$P_4$	5.0	4

## ■ SJF (preemptive)



■ Average waiting time =  $(9 + 1 + 0 + 2)/4 - 3$

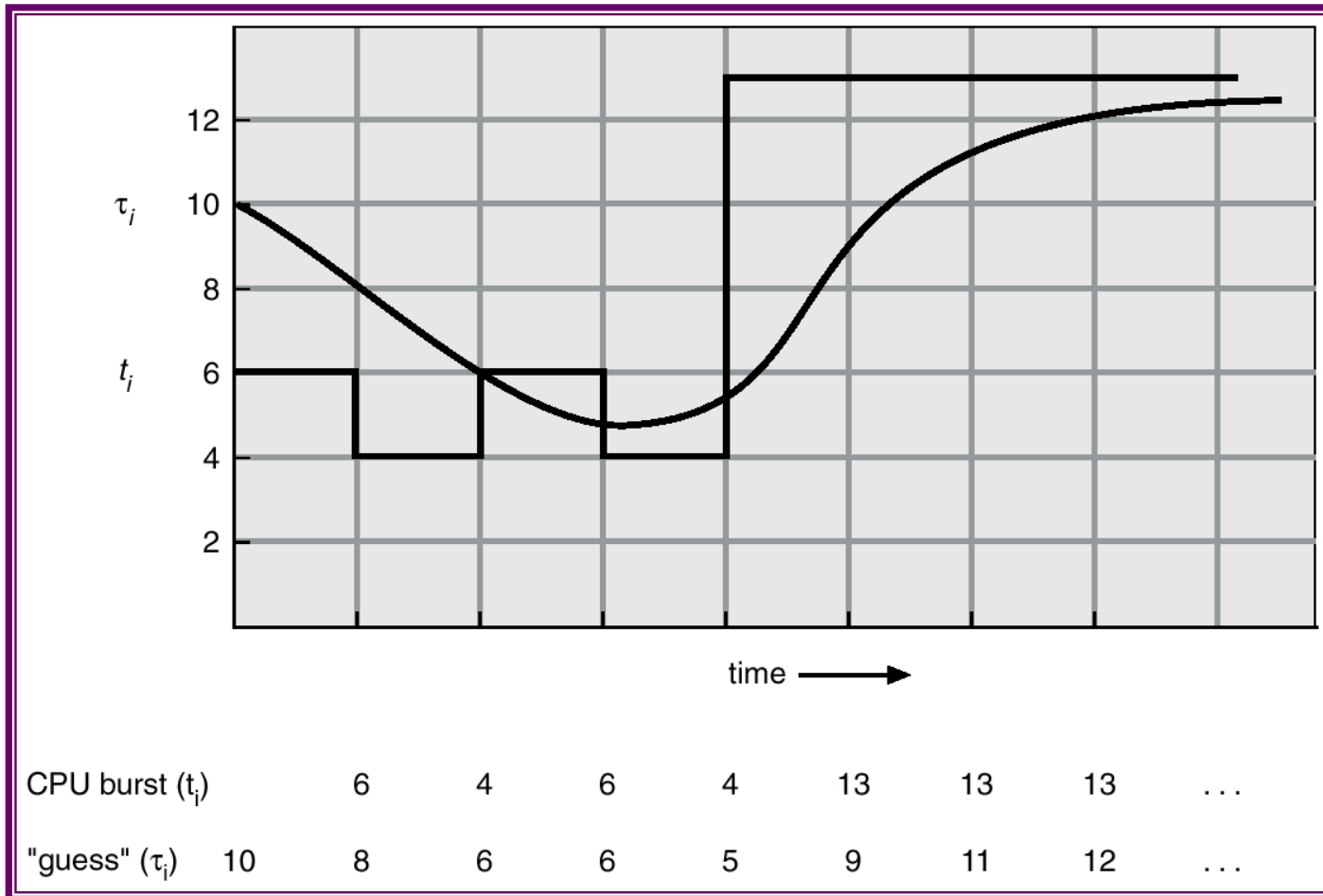
# Determining Length of Next CPU Burst

- Can only estimate the length.
- Can be done by using the length of previous CPU bursts, using exponential averaging.
  1.  $t_n$  = actual length of  $n^{\text{th}}$  CPU burst
  2.  $\tau_{n+1}$  = predicted value for the next CPU burst
  3.  $\alpha, 0 \leq \alpha \leq 1$
  4. Define :

$$\tau_{n+1} = \alpha t_n + (1 - \alpha)\tau_n.$$



# Prediction of the Length of the Next CPU Burst



# Examples of Exponential Averaging

- $\alpha = 0$

- ☞  $\tau_{n+1} = \tau_n$

- ☞ Recent history does not count.

- $\alpha = 1$

- ☞  $\tau_{n+1} = t_n$

- ☞ Only the actual last CPU burst counts.

- If we expand the formula, we get:

$$\begin{aligned}\tau_{n+1} = & \alpha t_n + (1 - \alpha) \alpha t_{n-1} + \dots \\ & + (1 - \alpha)^j \alpha t_{n-1} + \dots \\ & + (1 - \alpha)^{n=1} t_n \tau_0\end{aligned}$$

- Since both  $\alpha$  and  $(1 - \alpha)$  are less than or equal to 1, each successive term has less weight than its predecessor.

# Priority Scheduling

- A priority number (integer) is associated with each process
- The CPU is allocated to the process with the highest priority (smallest integer  $\equiv$  highest priority).
  - ☞ Preemptive
  - ☞ nonpreemptive
- SJF is a priority scheduling where priority is the predicted next CPU burst time.
- Problem  $\equiv$  Starvation – low priority processes may never execute.
- Solution  $\equiv$  Aging – as time progresses increase the priority of the process.

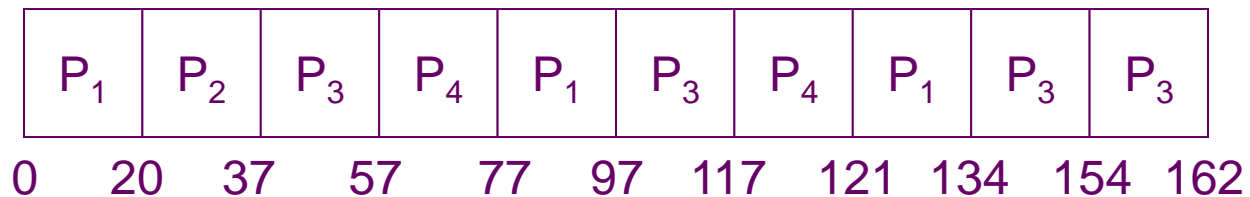
# Round Robin (RR)

- Each process gets a small unit of CPU time (*time quantum*), usually 10-100 milliseconds. After this time has elapsed, the process is preempted and added to the end of the ready queue.
- If there are  $n$  processes in the ready queue and the time quantum is  $q$ , then each process gets  $1/n$  of the CPU time in chunks of at most  $q$  time units at once. No process waits more than  $(n-1)q$  time units.
- Performance
  - ☞  $q$  large  $\Rightarrow$  FIFO
  - ☞  $q$  small  $\Rightarrow q$  must be large with respect to context switch, otherwise overhead is too high.

# Example of RR with Time Quantum = 20

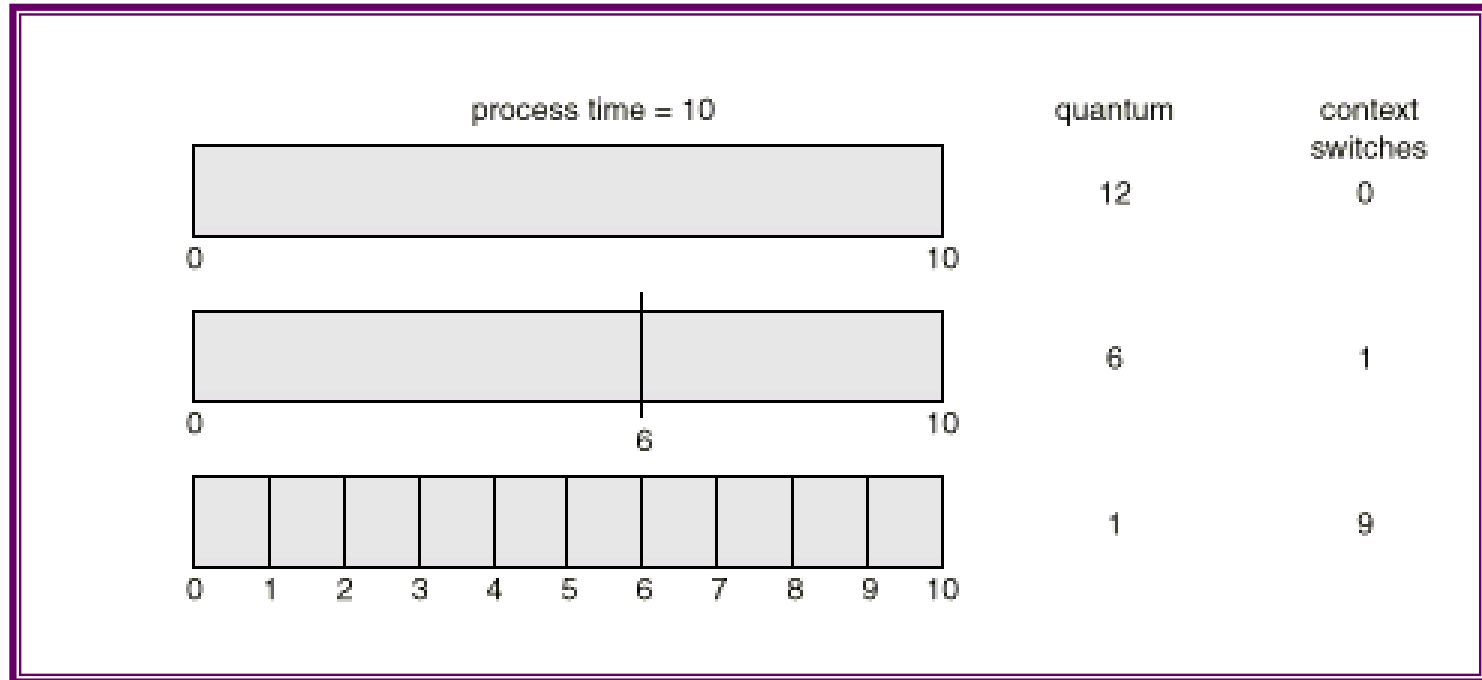
<u>Process</u>	<u>Burst Time</u>
$P_1$	53
$P_2$	17
$P_3$	68
$P_4$	24

- The Gantt chart is:

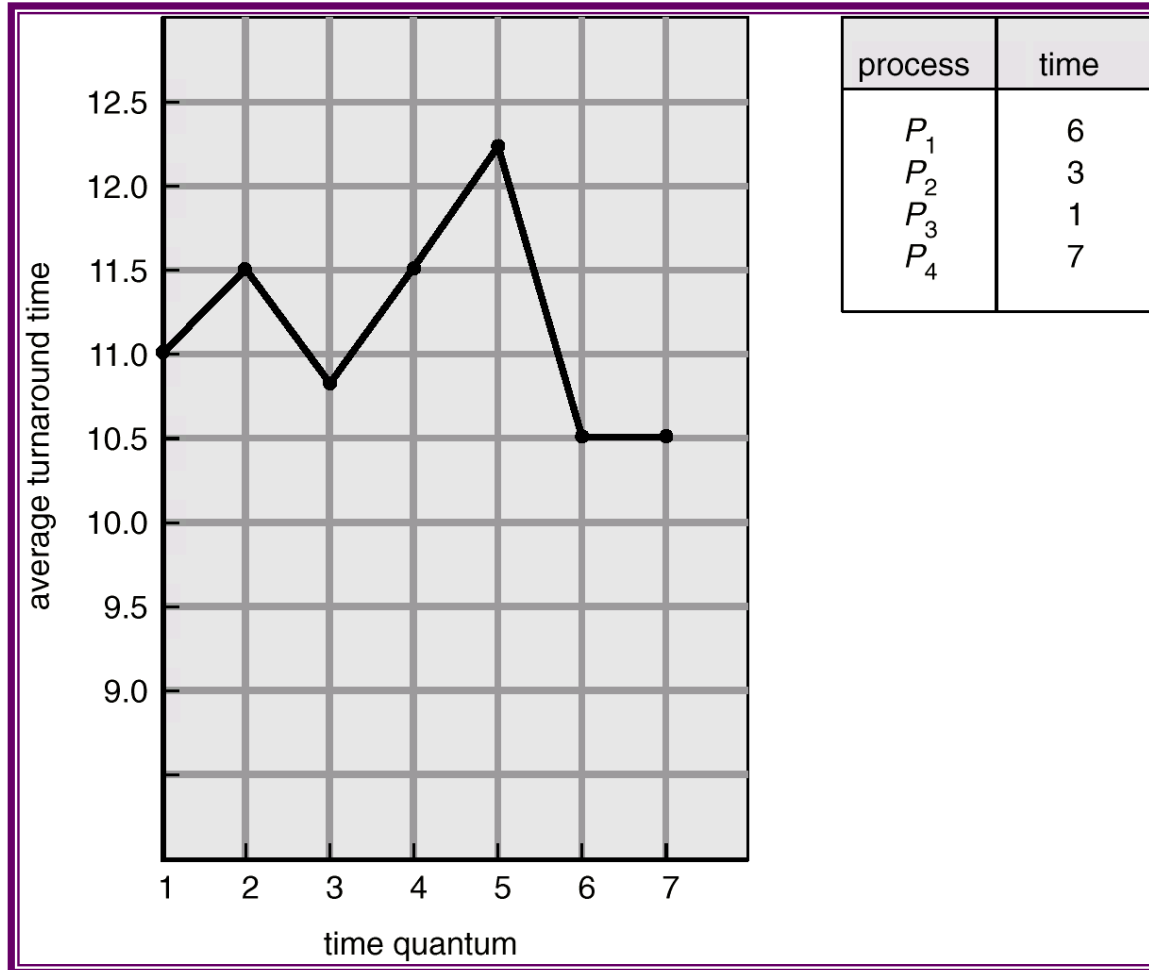


- Typically, higher average turnaround than SJF, but better *response*.

# Time Quantum and Context Switch Time



# Turnaround Time Varies With The Time Quantum

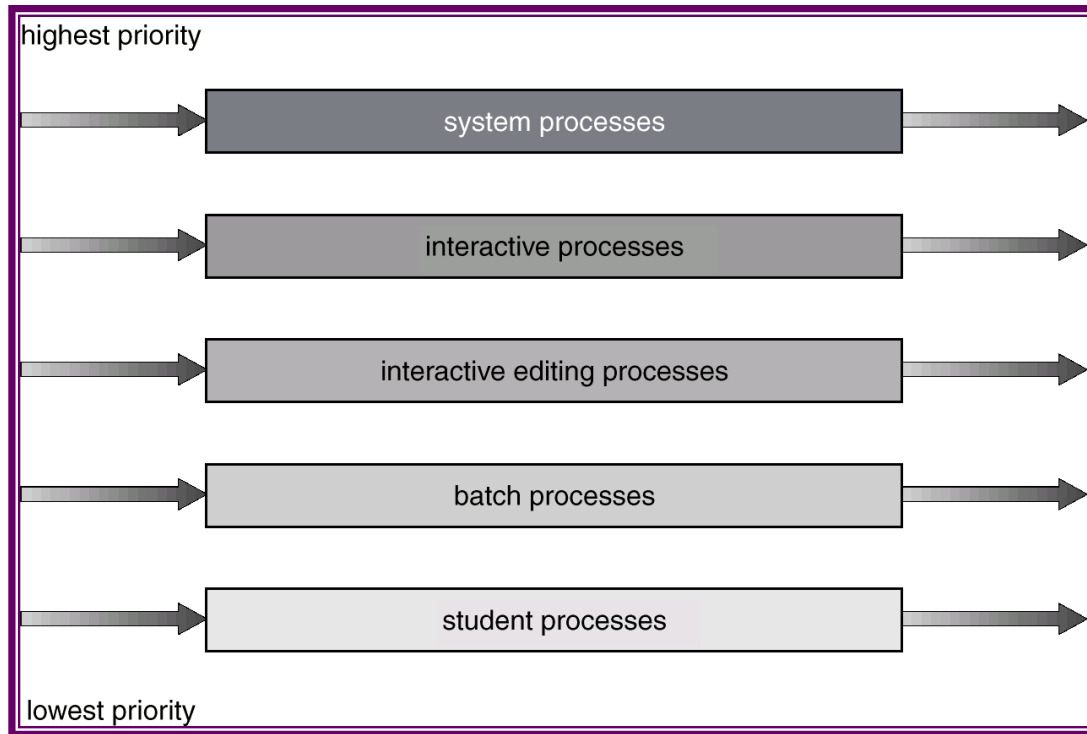


# Multilevel Queue

- Ready queue is partitioned into separate queues:
  - foreground (interactive)
  - background (batch)
- Each queue has its own scheduling algorithm,
  - foreground – RR
  - background – FCFS
- Scheduling must be done between the queues.
  - ☞ Fixed priority scheduling; (i.e., serve all from foreground then from background). Possibility of starvation.
  - ☞ Time slice – each queue gets a certain amount of CPU time which it can schedule amongst its processes; i.e., 80% to foreground in RR
  - ☞ 20% to background in FCFS



# Multilevel Queue Scheduling



# Multilevel Feedback Queue

- A process can move between the various queues; aging can be implemented this way.
- Multilevel-feedback-queue scheduler defined by the following parameters:
  - ☞ number of queues
  - ☞ scheduling algorithms for each queue
  - ☞ method used to determine when to upgrade a process
  - ☞ method used to determine when to demote a process
  - ☞ method used to determine which queue a process will enter when that process needs service

# Example of Multilevel Feedback Queue

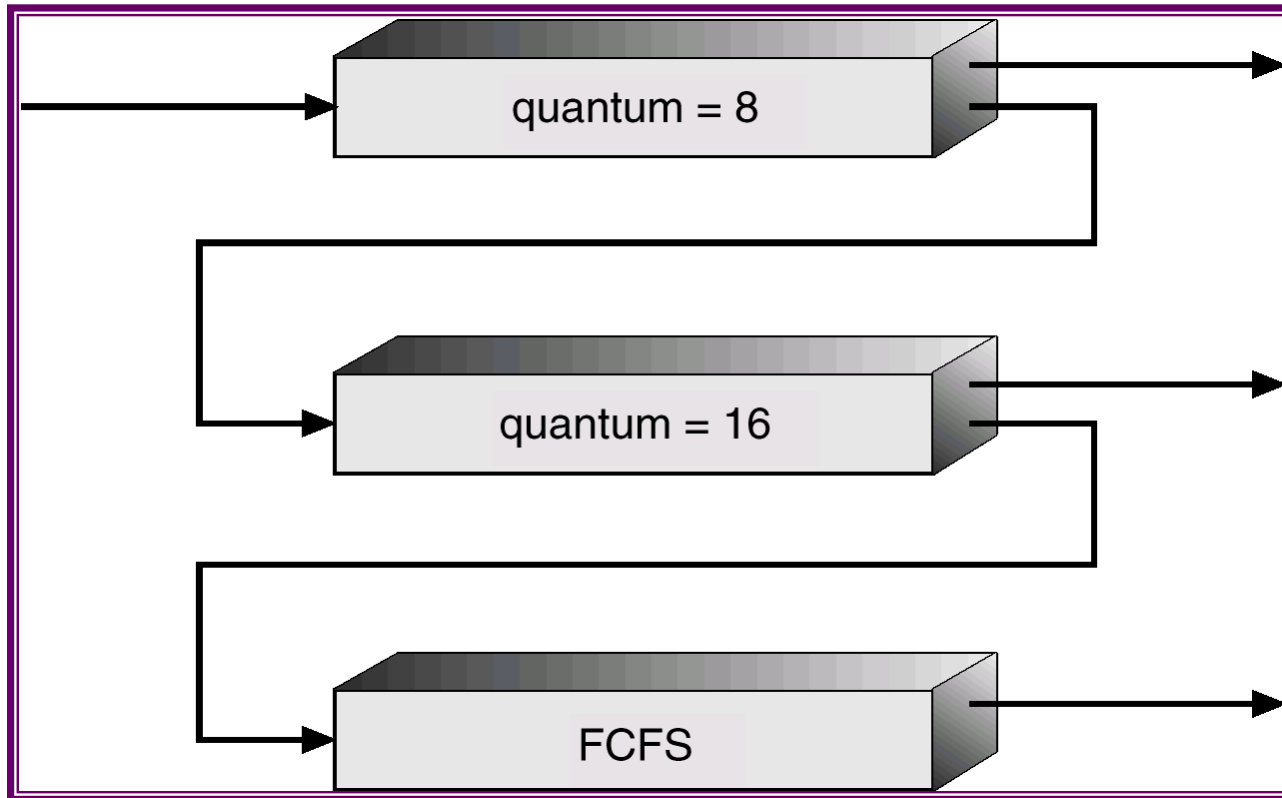
## ■ Three queues:

- ☞  $Q_0$  – time quantum 8 milliseconds
- ☞  $Q_1$  – time quantum 16 milliseconds
- ☞  $Q_2$  – FCFS

## ■ Scheduling

- ☞ A new job enters queue  $Q_0$  which is served FCFS. When it gains CPU, job receives 8 milliseconds. If it does not finish in 8 milliseconds, job is moved to queue  $Q_1$ .
- ☞ At  $Q_1$  job is again served FCFS and receives 16 additional milliseconds. If it still does not complete, it is preempted and moved to queue  $Q_2$ .

# Multilevel Feedback Queues



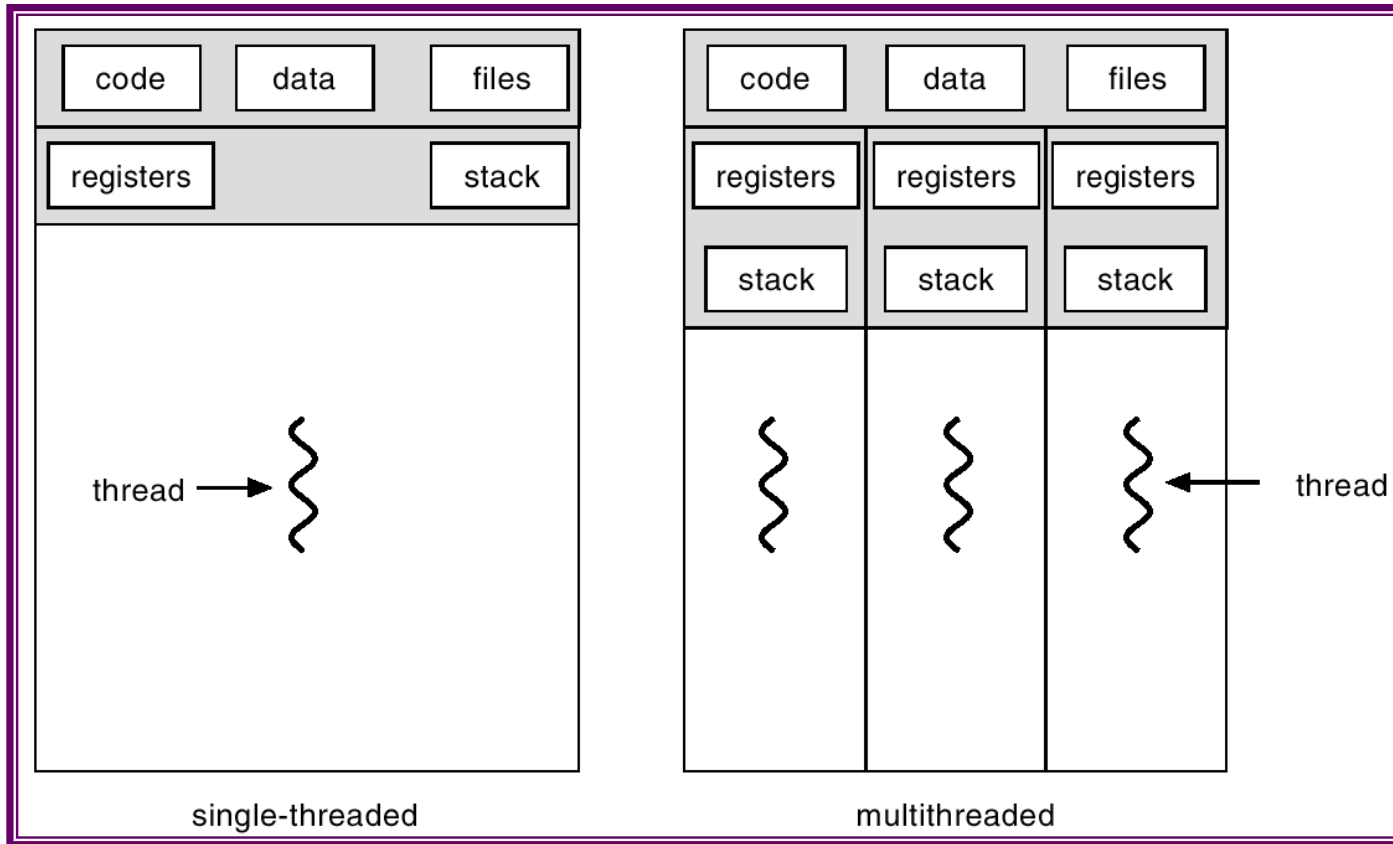
# Multiple-Processor Scheduling

- CPU scheduling more complex when multiple CPUs are available.
- *Homogeneous processors* within a multiprocessor.
- *Load sharing*
- *Asymmetric multiprocessing* – only one processor accesses the system data structures, alleviating the need for data sharing.

# Real-Time Scheduling

- *Hard real-time* systems – required to complete a critical task within a guaranteed amount of time.
- *Soft real-time* computing – requires that critical processes receive priority over less fortunate ones.

# Single and Multithreaded Processes



# Benefits

- Responsiveness
- Resource Sharing
- Economy
- Utilization of MP Architectures



# User Threads

- Thread management done by user-level threads library
- Examples
  - POSIX *Pthreads*
  - Mach *C-threads*
  - Solaris *threads*

# Kernel Threads

- Supported by the Kernel
- Examples
  - Windows 95/98/NT/2000
  - Solaris
  - Tru64 UNIX
  - BeOS
  - Linux

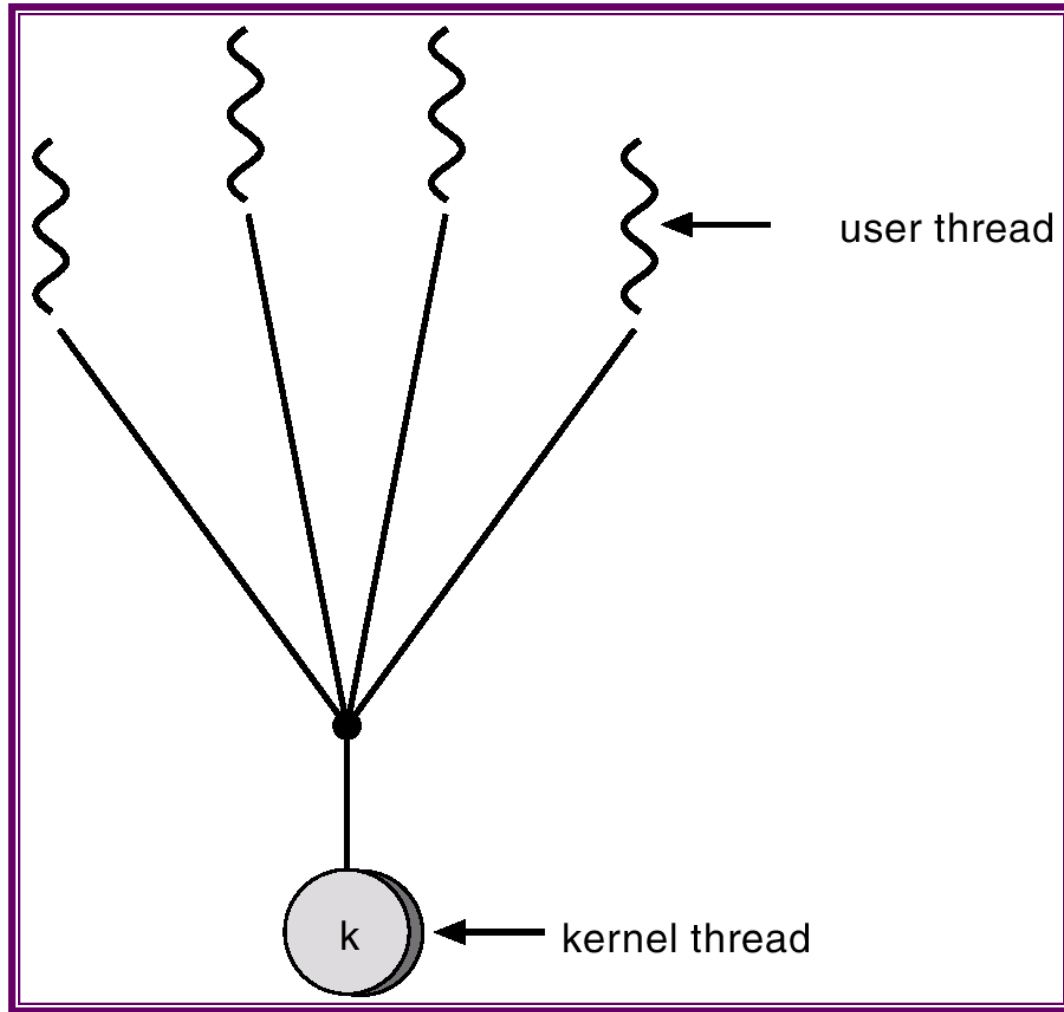
# Multithreading Models

- Many-to-One
- One-to-One
- Many-to-Many

# Many-to-One

- Many user-level threads mapped to single kernel thread.
- Used on systems that do not support kernel threads.

# Many-to-One Model



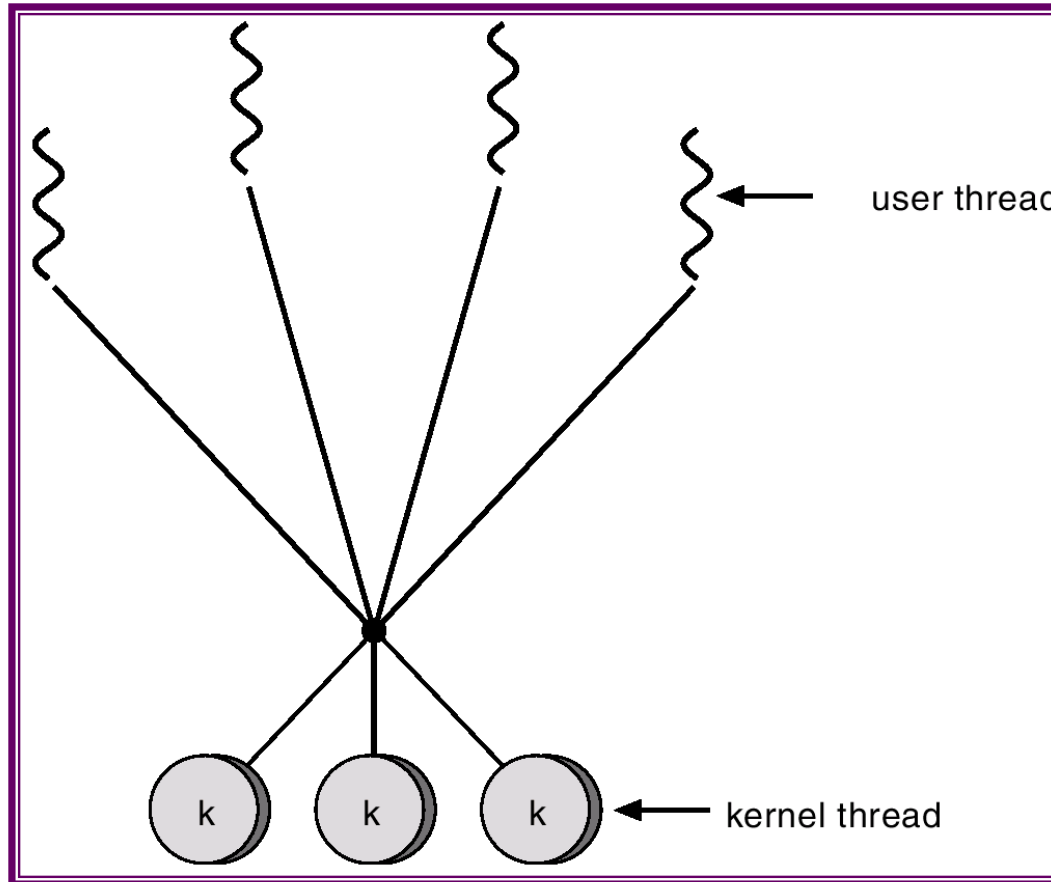
# One-to-One

- Each user-level thread maps to kernel thread.
- Examples
  - Windows 95/98/NT/2000
  - OS/2

# Many-to-Many Model

- Allows many user level threads to be mapped to many kernel threads.
- Allows the operating system to create a sufficient number of kernel threads.
- Solaris 2
- Windows NT/2000 with the *ThreadFiber* package

# Many-to-Many Model





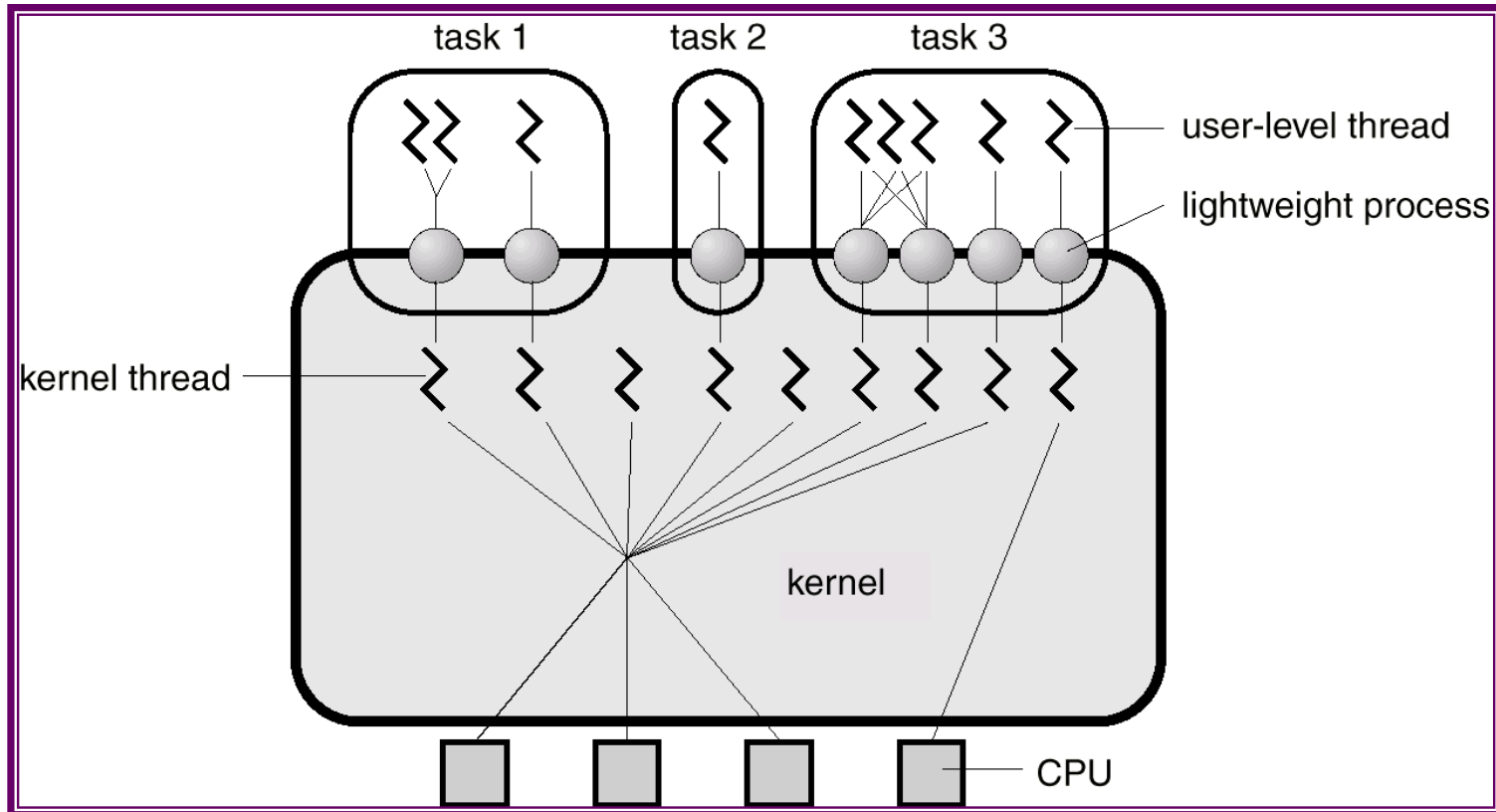
# Threading Issues

- Semantics of fork() and exec() system calls.
- Thread cancellation.
- Signal handling
- Thread pools
- Thread specific data

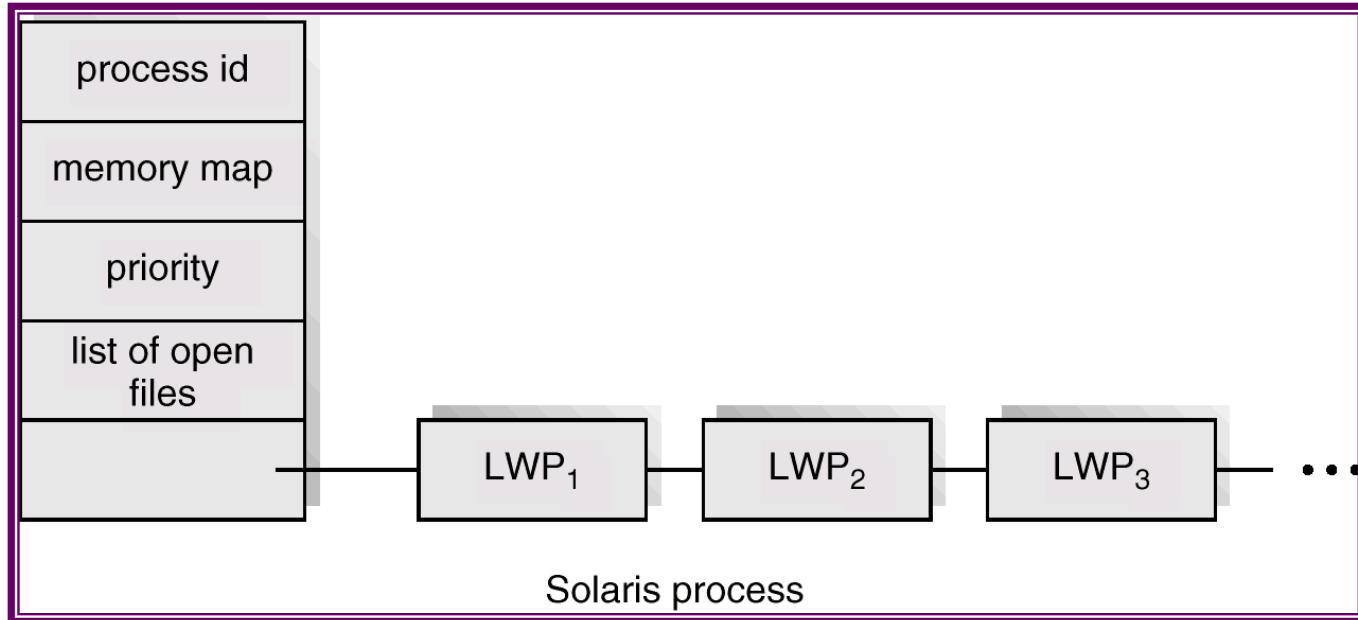
# Pthreads

- a POSIX standard (IEEE 1003.1c) API for thread creation and synchronization.
- API specifies behavior of the thread library, implementation is up to development of the library.
- Common in UNIX operating systems.

# Solaris 2 Threads



# Solaris Process



# Windows 2000 Threads

- Implements the one-to-one mapping.
- Each thread contains
  - a thread id
  - register set
  - separate user and kernel stacks
  - private data storage area

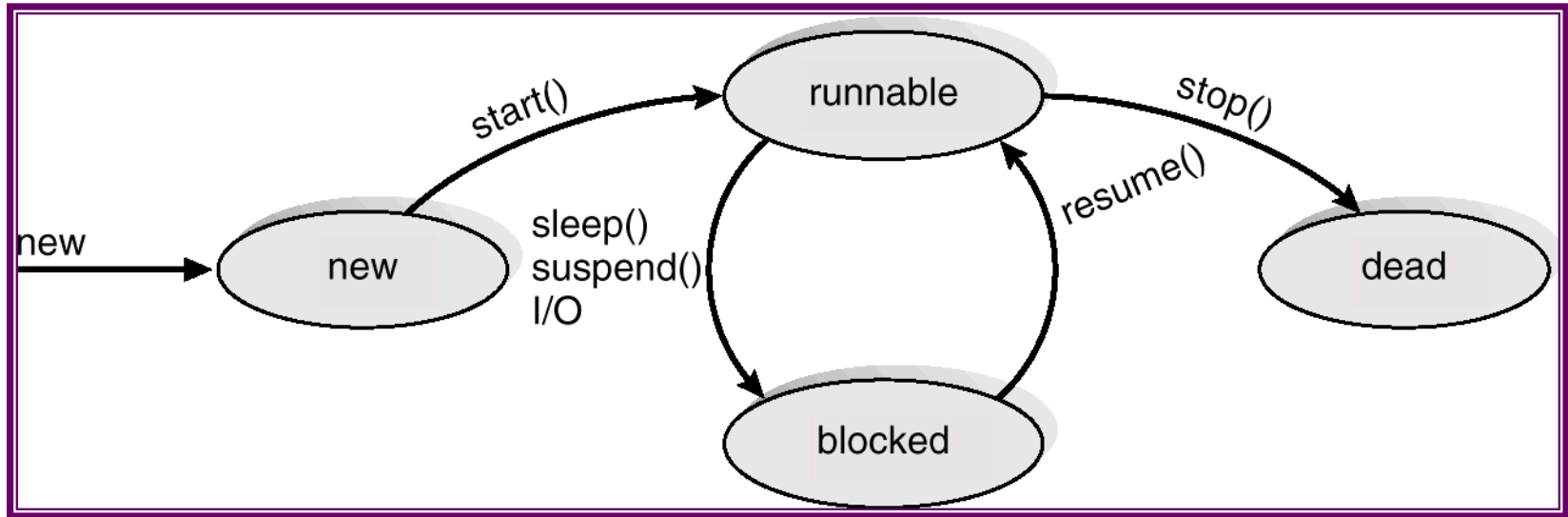
# Linux Threads

- Linux refers to them as *tasks* rather than *threads*.
- Thread creation is done through `clone()` system call.
- `Clone()` allows a child task to share the address space of the parent task (process)

# Java Threads

- Java threads may be created by:
  - ☞ Extending Thread class
  - ☞ Implementing the Runnable interface
- Java threads are managed by the JVM.

# Java Thread States







# UNIT 2: Process Synchronization

- Background
- The Critical-Section Problem
- Synchronization Hardware
- Semaphores
- Classical Problems of Synchronization
- Critical Regions
- Monitors
- Synchronization in Solaris 2 & Windows 2000

# Background

- Concurrent access to shared data may result in data inconsistency.
- Maintaining data consistency requires mechanisms to ensure the orderly execution of cooperating processes.
- Shared-memory solution to bounded-buffer problem (Chapter 4) allows at most  $n - 1$  items in buffer at the same time. A solution, where all  $N$  buffers are used is not simple.
  - ☞ Suppose that we modify the producer-consumer code by adding a variable *counter*, initialized to 0 and incremented each time a new item is added to the buffer

# Bounded-Buffer

- Shared data

```
#define BUFFER_SIZE 10
typedef struct {
    ...
} item;
item buffer[BUFFER_SIZE];
int in = 0;
int out = 0;
int counter = 0;
```

# Bounded-Buffer

- Producer process

```
item nextProduced;
```

```
while (1) {  
    while (counter == BUFFER_SIZE)  
        ; /* do nothing */  
    buffer[in] = nextProduced;  
    in = (in + 1) % BUFFER_SIZE;  
    counter++;  
}
```

# Bounded-Buffer

- Consumer process

```
item nextConsumed;
```

```
while (1) {  
    while (counter == 0)  
        ; /* do nothing */  
    nextConsumed = buffer[out];  
    out = (out + 1) % BUFFER_SIZE;  
    counter--;  
}
```

# Bounded Buffer

- The statements

```
counter++;  
counter--;
```

must be performed *atomically*.

- Atomic operation means an operation that completes in its entirety without interruption.

# Bounded Buffer

- The statement “**count++**” may be implemented in machine language as:

**register1 = counter**

**register1 = register1 + 1**

**counter = register1**

- The statement “**count—**” may be implemented as:

**register2 = counter**

**register2 = register2 – 1**

**counter = register2**



# Bounded Buffer

- If both the producer and consumer attempt to update the buffer concurrently, the assembly language statements may get interleaved.
- Interleaving depends upon how the producer and consumer processes are scheduled.

# Bounded Buffer

- Assume **counter** is initially 5. One interleaving of statements is:

producer: **register1 = counter** (*register1 = 5*)

producer: **register1 = register1 + 1** (*register1 = 6*)

consumer: **register2 = counter** (*register2 = 5*)

consumer: **register2 = register2 - 1** (*register2 = 4*)

producer: **counter = register1** (*counter = 6*)

consumer: **counter = register2** (*counter = 4*)

- The value of **count** may be either 4 or 6, where the correct result should be 5.

# Race Condition

- **Race condition:** The situation where several processes access – and manipulate shared data concurrently. The final value of the shared data depends upon which process finishes last.
- To prevent race conditions, concurrent processes must be **synchronized**.

# The Critical-Section Problem

- $n$  processes all competing to use some shared data
- Each process has a code segment, called *critical section*, in which the shared data is accessed.
- Problem – ensure that when one process is executing in its critical section, no other process is allowed to execute in its critical section.

# 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.
  - Assume that each process executes at a nonzero speed
  - No assumption concerning relative speed of the  $n$  processes.

# Initial Attempts to Solve Problem

- Only 2 processes,  $P_0$  and  $P_1$
- General structure of process  $P_i$  (other process  $P_j$ )

```
do {   
    entry section  
      
    exit section  
    reminder section  
} while (1);
```
- Processes may share some common variables to synchronize their actions.

# Algorithm 1

- Shared variables:

- ☞ **int turn;**  
initially **turn = 0**

- ☞ **turn = i** ⇒  $P_i$  can enter its critical section

- Process  $P_i$

```
do {  
    while (turn != i) ;  
    critical section  
    turn = j;  
    reminder section  
} while (1);
```

- Satisfies mutual exclusion, but not progress

# Algorithm 2

- Shared variables

  - ☞ **boolean flag[2];**

    - initially **flag [0] = flag [1] = false.**

  - ☞ **flag [i] = true**  $\Rightarrow P_i$  ready to enter its critical section

- Process  $P_i$

      - do {**

        - flag[i] := true;**

        - while (flag[j]) ;**

          - critical section

        - flag [i] = false;**

          - remainder section

      - } while (1);**

- Satisfies mutual exclusion, but not progress requirement.



# Algorithm 3

- Combined shared variables of algorithms 1 and 2.
- Process  $P_i$ 
  - do** {
    - flag [i] := true;**
    - turn = j;**
    - while (flag [j] and turn = j) ;**
      - critical section
    - flag [i] = false;**
    - remainder section
  - } while (1);**
- Meets all three requirements; solves the critical-section problem for two processes.

# Bakery Algorithm

## Critical section for n processes

- Before entering its critical section, process receives a number. Holder of the smallest number enters the critical section.
- If processes  $P_i$  and  $P_j$  receive the same number, if  $i < j$ , then  $P_i$  is served first; else  $P_j$  is served first.
- The numbering scheme always generates numbers in increasing order of enumeration; i.e., 1,2,3,3,3,3,4,5...

# Bakery Algorithm

- Notation  $\leq$  lexicographical order (ticket #, process id #)

- ☞  $(a,b) < c,d$  if  $a < c$  or if  $a = c$  and  $b < d$

- ☞  $\max(a_0, \dots, a_{n-1})$  is a number,  $k$ , such that  $k \geq a_i$  for  $i = 0, \dots, n-1$

- Shared data

- boolean choosing[n];**

- int number[n];**

Data structures are initialized to **false** and **0** respectively

# Bakery Algorithm

```
do {
    choosing[i] = true;
    number[i] = max(number[0], number[1], ..., number [n - 1])+1;
    choosing[i] = false;
    for (j = 0; j < n; j++) {
        while (choosing[j]) ;
        while ((number[j] != 0) && (number[j,j] < number[i,i])) ;
    }
    critical section
    number[i] = 0;
    remainder section
} while (1);
```

# Synchronization Hardware

- Test and modify the content of a word atomically

.

```
boolean TestAndSet(boolean &target) {  
    boolean rv = target;  
    target = true;  
  
    return rv;  
}
```

# Mutual Exclusion with Test-and-Set

- Shared data:

```
boolean lock = false;
```

- Process  $P_i$

```
do {  
    while (TestAndSet(lock)) ;  
    critical section  
    lock = false;  
    remainder section  
}
```

# Synchronization Hardware

- Atomically swap two variables.

```
void Swap(boolean &a, boolean &b) {  
    boolean temp = a;  
    a = b;  
    b = temp;  
}
```

# Mutual Exclusion with Swap

- Shared data (initialized to **false**):  
**boolean lock;**  
**boolean waiting[n];**
- Process  $P_i$   
do {  
    **key = true;**  
    **while (key == true)**  
        **Swap(lock,key);**  
        critical section  
    **lock = false;**  
    remainder section  
}



# Semaphores

- Synchronization tool that does not require busy waiting.
- Semaphore  $S$  – integer variable
- can only be accessed via two indivisible (atomic) operations

*wait* ( $S$ ):

```
while  $S \leq 0$  do no-op;  
S--;
```

*signal* ( $S$ ):

```
S++;
```

# Critical Section of $n$ Processes

- Shared data:  
**semaphore mutex;** //initially *mutex* = 1

- Process  $P_i$ :

```
do {  
    wait(mutex);  
    critical section  
    signal(mutex);  
    remainder section  
} while (1);
```

# Semaphore Implementation

- Define a semaphore as a record

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

- Assume two simple operations:

- ☞ **block** suspends the process that invokes it.
- ☞ **wakeup(P)** resumes the execution of a blocked process **P**.

# Implementation

- Semaphore operations now defined as

*wait(S):*

```
S.value--;  
if (S.value < 0) {  
    add this process to S.L;  
    block;  
}
```

*signal(S):*

```
S.value++;  
if (S.value <= 0) {  
    remove a process P from S.L;  
    wakeup(P);  
}
```

# Semaphore as a General Synchronization Tool

- Execute  $B$  in  $P_j$  only after  $A$  executed in  $P_i$
- Use semaphore  $flag$  initialized to 0
- Code:

$P_i$	$P_j$
$\vdots$	$\vdots$
$A$	$wait(flag)$
$signal(flag)$	$B$

# Deadlock and Starvation

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

$P_0$	$P_1$
<i>wait(S);</i>	<i>wait(Q);</i>
<i>wait(Q);</i>	<i>wait(S);</i>
$\vdots$	$\vdots$
<i>signal(S);</i>	<i>signal(Q);</i>
<i>signal(Q)</i>	<i>signal(S);</i>

- **Starvation** – indefinite blocking. A process may never be removed from the semaphore queue in which it is suspended.

# Two Types of Semaphores

- *Counting* semaphore – integer value can range over an unrestricted domain.
- *Binary* semaphore – integer value can range only between 0 and 1; can be simpler to implement.
- Can implement a counting semaphore  $S$  as a binary semaphore.

# Implementing $S$ as a Binary Semaphore

- Data structures:

**binary-semaphore S1, S2;**

**int C;**

- Initialization:

**S1 = 1**

**S2 = 0**

**C = initial value of semaphore S**



# Implementing S

- *wait* operation

```
wait(S1);  
C--;  
if (C < 0) {  
    signal(S1);  
    wait(S2);  
}  
signal(S1);
```

- *signal* operation

```
wait(S1);  
C ++;  
if (C <= 0)  
    signal(S2);  
else  
    signal(S1);
```

# Classical Problems of Synchronization

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

# Bounded-Buffer Problem

- Shared data

**semaphore full, empty, mutex;**

Initially:

**full = 0, empty = n, mutex = 1**

# Bounded-Buffer Problem Producer Process

```
do {  
    ...  
    produce an item in nextp  
    ...  
    wait(empty);  
    wait(mutex);  
    ...  
    add nextp to buffer  
    ...  
    signal(mutex);  
    signal(full);  
} while (1);
```

# Bounded-Buffer Problem Consumer Process

```
do {  
    wait(full)  
    wait(mutex);  
    ...  
    remove an item from buffer to nextc  
    ...  
    signal(mutex);  
    signal(empty);  
    ...  
    consume the item in nextc  
    ...  
} while (1);
```

# Readers-Writers Problem

- Shared data

**semaphore mutex, wrt;**

Initially

**mutex = 1, wrt = 1, readcount = 0**

# Readers-Writers Problem Writer Process

**wait(wrt);**

...

writing is performed

...

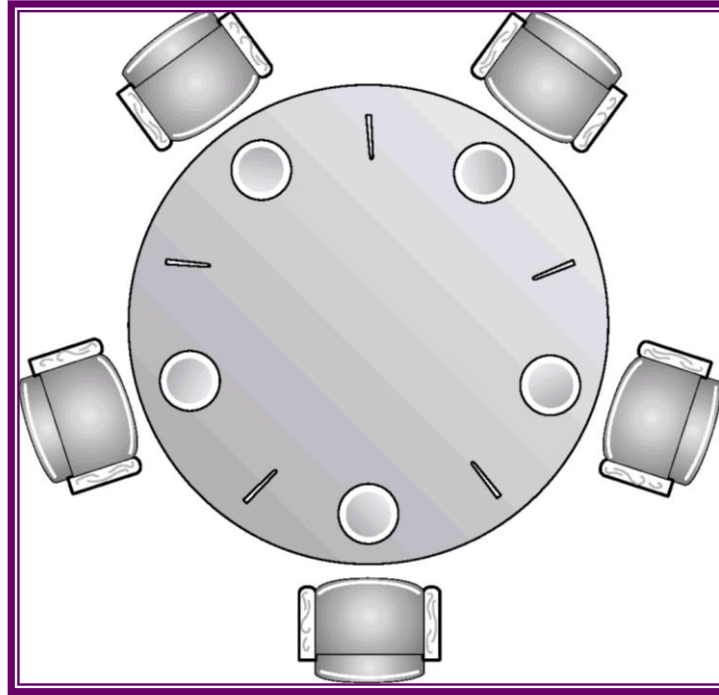
**signal(wrt);**

# Readers-Writers Problem Reader Process

```
wait(mutex);  
readcount++;  
if (readcount == 1)  
    wait(rt);  
signal(mutex);  
    ...  
    reading is performed  
    ...  
wait(mutex);  
readcount--;  
if (readcount == 0)  
    signal(wrt);  
signal(mutex);
```



# Dining-Philosophers Problem



- Shared data

**semaphore chopstick[5];**

Initially all values are 1

# Dining-Philosophers Problem

- Philosopher  $i$ :

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

# Critical Regions

- High-level synchronization construct
- A shared variable  $v$  of type  $T$ , is declared as:  
 **$v$ : shared  $T$**
- Variable  $v$  accessed only inside statement  
**region  $v$  when  $B$  do  $S$**

where  $B$  is a boolean expression.

- While statement  $S$  is being executed, no other process can access variable  $v$ .

# Critical Regions

- Regions referring to the same shared variable exclude each other in time.
- When a process tries to execute the region statement, the Boolean expression  $B$  is evaluated. If  $B$  is true, statement  $S$  is executed. If it is false, the process is delayed until  $B$  becomes true and no other process is in the region associated with  $v$ .

# Example – Bounded Buffer

- Shared data:

```
struct buffer {  
    int pool[n];  
    int count, in, out;  
}
```

# Bounded Buffer Producer Process

- Producer process inserts **nextp** into the shared buffer

```
region buffer when( count < n) {  
    pool[in] = nextp;  
    in:= (in+1) % n;  
    count++;  
}
```

# Bounded Buffer Consumer Process

- Consumer process removes an item from the shared buffer and puts it in **nextc**

```
region buffer when (count > 0) {  
    nextc = pool[out];  
    out = (out+1) % n;  
    count--;  
}
```

# Implementation region $x$ when $B$ do $S$

- Associate with the shared variable  $x$ , the following variables:  

```
semaphore mutex, first-delay, second-delay;  
int first-count, second-count;
```
- Mutually exclusive access to the critical section is provided by **mutex**.
- If a process cannot enter the critical section because the Boolean expression **B** is false, it initially waits on the **first-delay** semaphore; moved to the **second-delay** semaphore before it is allowed to reevaluate  $B$ .



# Implementation

- Keep track of the number of processes waiting on **first-delay** and **second-delay**, with **first-count** and **second-count** respectively.
- The algorithm assumes a FIFO ordering in the queuing of processes for a semaphore.
- For an arbitrary queuing discipline, a more complicated implementation is required.

# Monitors

- High-level synchronization construct that allows the safe sharing of an abstract data type among concurrent processes.

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

# Monitors

- To allow a process to wait within the monitor, a **condition** variable must be declared, as

**condition x, y;**

- Condition variable can only be used with the operations **wait** and **signal**.

☞ The operation

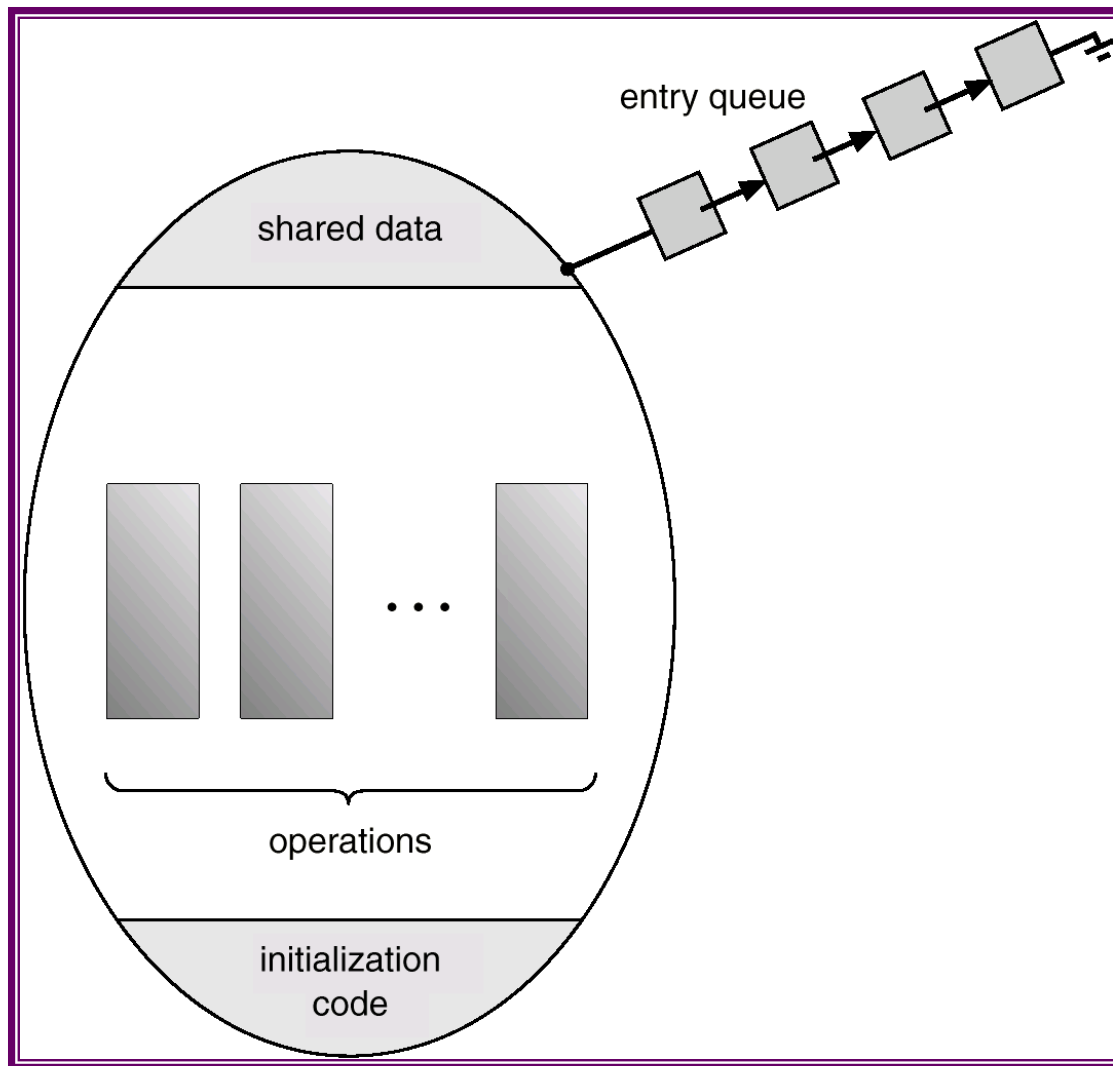
**x.wait();**

means that the process invoking this operation is suspended until another process invokes

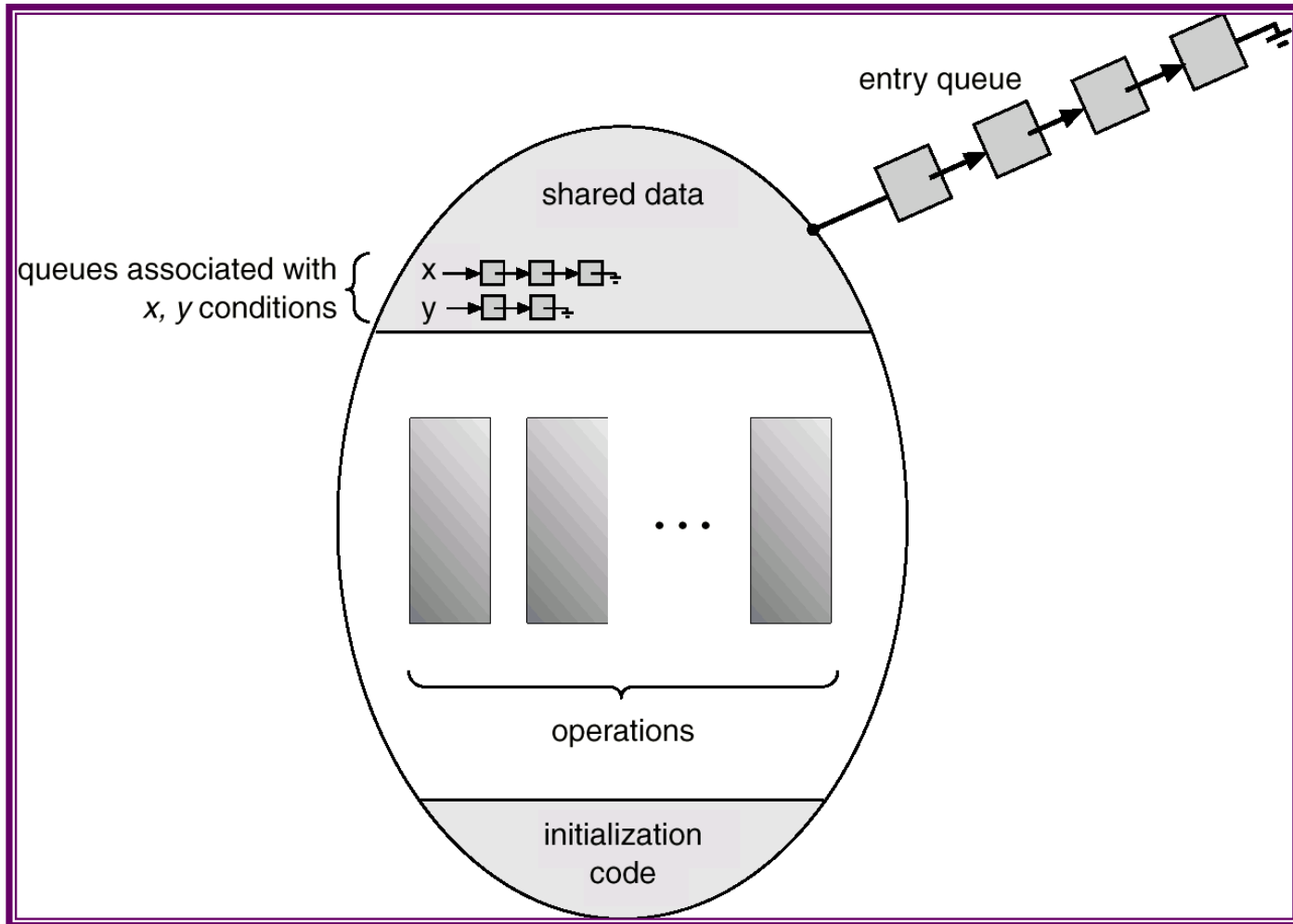
**x.signal();**

☞ The **x.signal** operation resumes exactly one suspended process. If no process is suspended, then the **signal** operation has no effect.

# Schematic View of a Monitor



# Monitor With Condition Variables



# Dining Philosophers Example

```
monitor dp
{
    enum {thinking, hungry, eating} state[5];
    condition self[5];
    void pickup(int i)           // following slides
    void putdown(int i)         // following slides
    void test(int i)            // following slides
    void init() {
        for (int i = 0; i < 5; i++)
            state[i] = thinking;
    }
}
```

# Dining Philosophers

```
void pickup(int i) {  
    state[i] = hungry;  
    test(i);  
    if (state[i] != eating)  
        self[i].wait();  
}
```

```
void putdown(int i) {  
    state[i] = thinking;  
    // test left and right neighbors  
    test((i+4) % 5);  
    test((i+1) % 5);  
}
```

# Dining Philosophers

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



# Monitor Implementation Using Semaphores

- Variables

```
semaphore mutex; // (initially = 1)
semaphore next;  // (initially = 0)
int next-count = 0;
```

- Each external procedure  $F$  will be replaced by

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

- Mutual exclusion within a monitor is ensured.

# Monitor Implementation

- For each condition variable  $x$ , we have:  
**semaphore x-sem; // (initially = 0)**  
**int x-count = 0;**

- The operation  $x.wait$  can be implemented as:

```
x-count++;  
if (next-count > 0)  
    signal(next);  
else  
    signal(mutex);  
wait(x-sem);  
x-count--;
```

# Monitor Implementation

- The operation **x.signal** can be implemented as:

```
if (x-count > 0) {  
    next-count++;  
    signal(x-sem);  
    wait(next);  
    next-count--;  
}
```

# Monitor Implementation

- **Conditional-wait** construct: **x.wait(c);**
  - ☞ **c** – integer expression evaluated when the **wait** operation is executed.
  - ☞ value of **c** (a *priority number*) stored with the name of the process that is suspended.
  - ☞ when **x.signal** is executed, process with smallest associated priority number is resumed next.
- Check two conditions to establish correctness of system:
  - ☞ User processes must always make their calls on the monitor in a correct sequence.
  - ☞ Must ensure that an uncooperative process does not ignore the mutual-exclusion gateway provided by the monitor, and try to access the shared resource directly, without using the access protocols.

# Solaris 2 Synchronization

- Implements a variety of locks to support multitasking, multithreading (including real-time threads), and multiprocessing.
- Uses *adaptive mutexes* for efficiency when protecting data from short code segments.
- Uses *condition variables* and *readers-writers* locks when longer sections of code need access to data.
- Uses *turnstile*s to order the list of threads waiting to acquire either an adaptive mutex or reader-writer lock.

# Windows 2000 Synchronization

- Uses interrupt masks to protect access to global resources on uniprocessor systems.
- Uses *spinlocks* on multiprocessor systems.
- Also provides *dispatcher objects* which may act as wither mutexes and semaphores.
- Dispatcher objects may also provide *events*. An event acts much like a condition variable.



# UNIT 3: Memory Management

- Background
- Swapping
- Contiguous Allocation
- Paging
- Segmentation
- Segmentation with Paging
- Demand Paging
- Process Creation
- Page Replacement
- Allocation of Frames
- Thrashing
- Operating System Examples



# Background

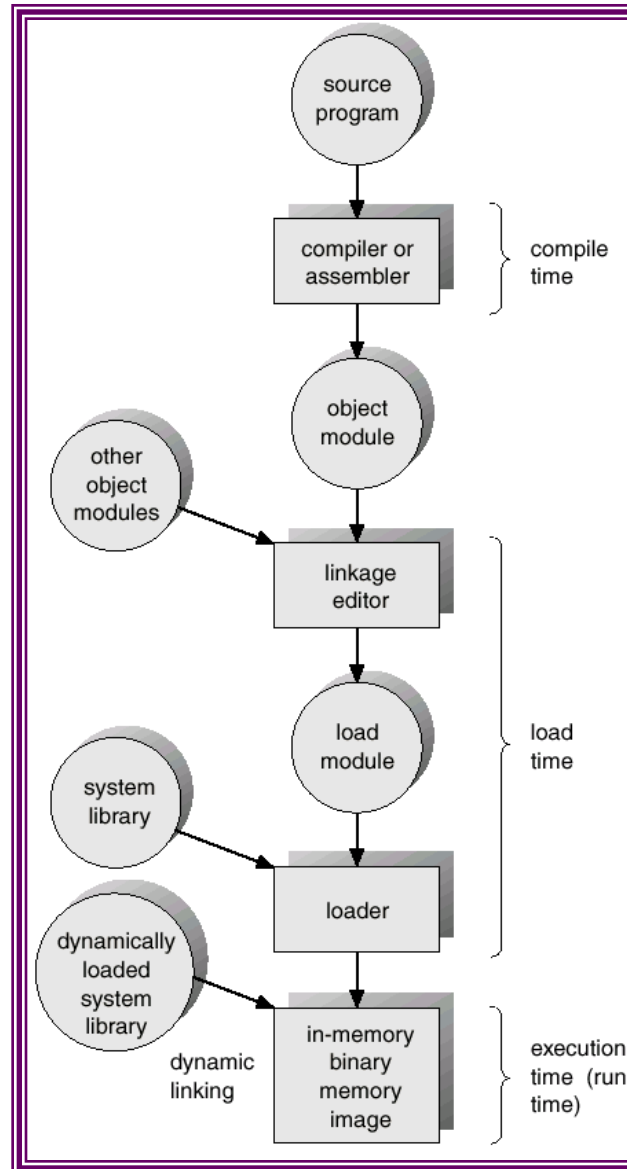
- Program must be brought into memory and placed within a process for it to be run.
- *Input queue* – collection of processes on the disk that are waiting to be brought into memory to run the program.
- User programs go through several steps before being run.

# Binding of Instructions and Data to Memory

Address binding of instructions and data to memory addresses can happen at three different stages.

- **Compile time:** If memory location known a priori, absolute code can be generated; must recompile code if starting location changes.
- **Load time:** Must generate *relocatable* code if memory location is not known at compile time.
- **Execution time:** Binding delayed until run time if the process can be moved during its execution from one memory segment to another. Need hardware support for address maps (e.g., *base* and *limit registers*).

# Multistep Processing of a User Program



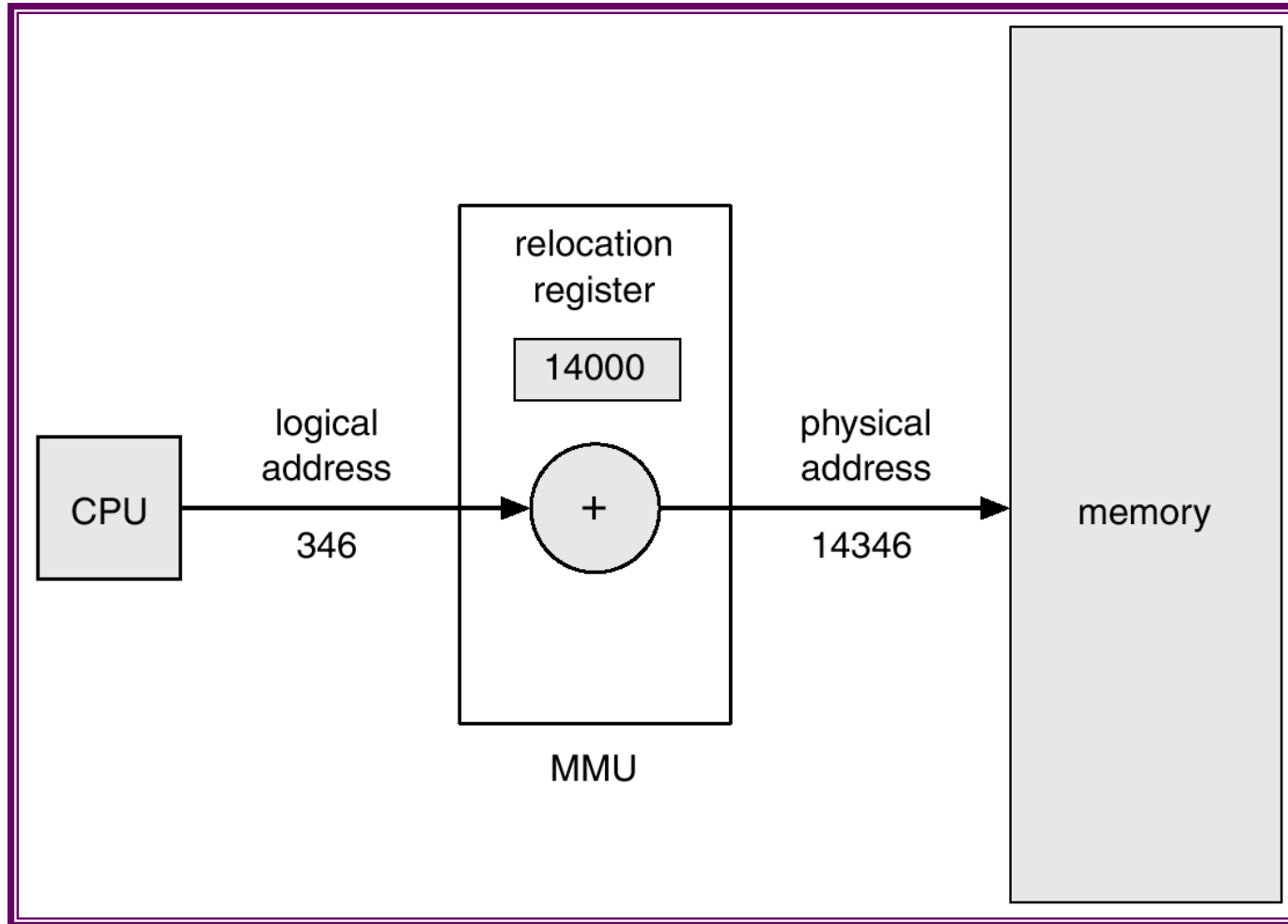
# Logical vs. Physical Address Space

- 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.

# Memory-Management Unit (MMU)

- Hardware device that maps virtual to physical address.
- In MMU scheme, the value in the relocation register is added to every address generated by a user process at the time it is sent to memory.
- The user program deals with *logical* addresses; it never sees the *real* physical addresses.

# Dynamic relocation using a relocation register



# Dynamic Loading

- Routine is not loaded until it is called
- Better memory-space utilization; unused routine is never loaded.
- Useful when large amounts of code are needed to handle infrequently occurring cases.
- No special support from the operating system is required implemented through program design.

# Dynamic Linking

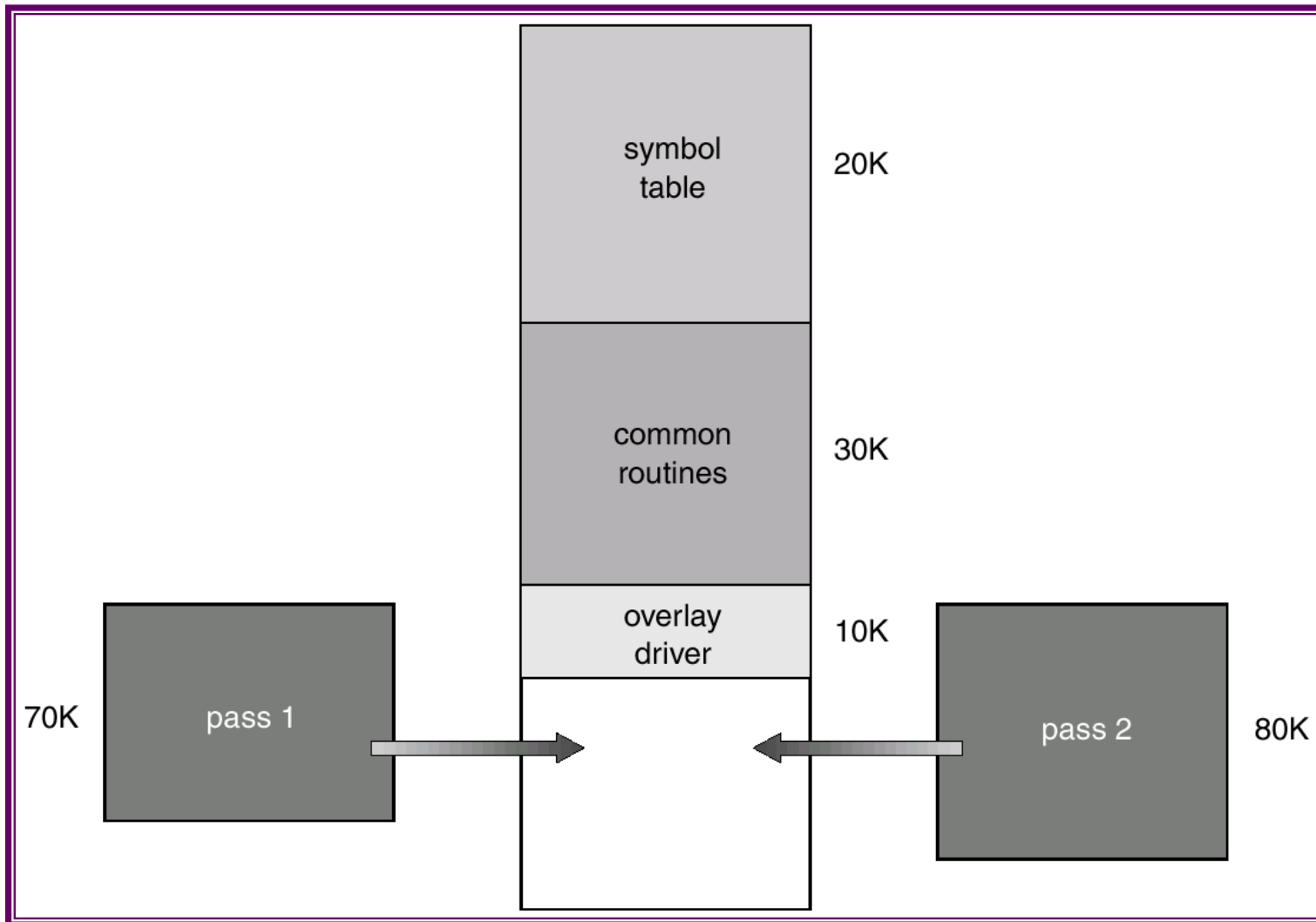
- Linking postponed until execution time.
- Small piece of code, *stub*, used to locate the appropriate memory-resident library routine.
- Stub replaces itself with the address of the routine, and executes the routine.
- Operating system needed to check if routine is in processes' memory address.
- Dynamic linking is particularly useful for libraries.



# Overlays

- Keep in memory only those instructions and data that are needed at any given time.
- Needed when process is larger than amount of memory allocated to it.
- Implemented by user, no special support needed from operating system, programming design of overlay structure is complex

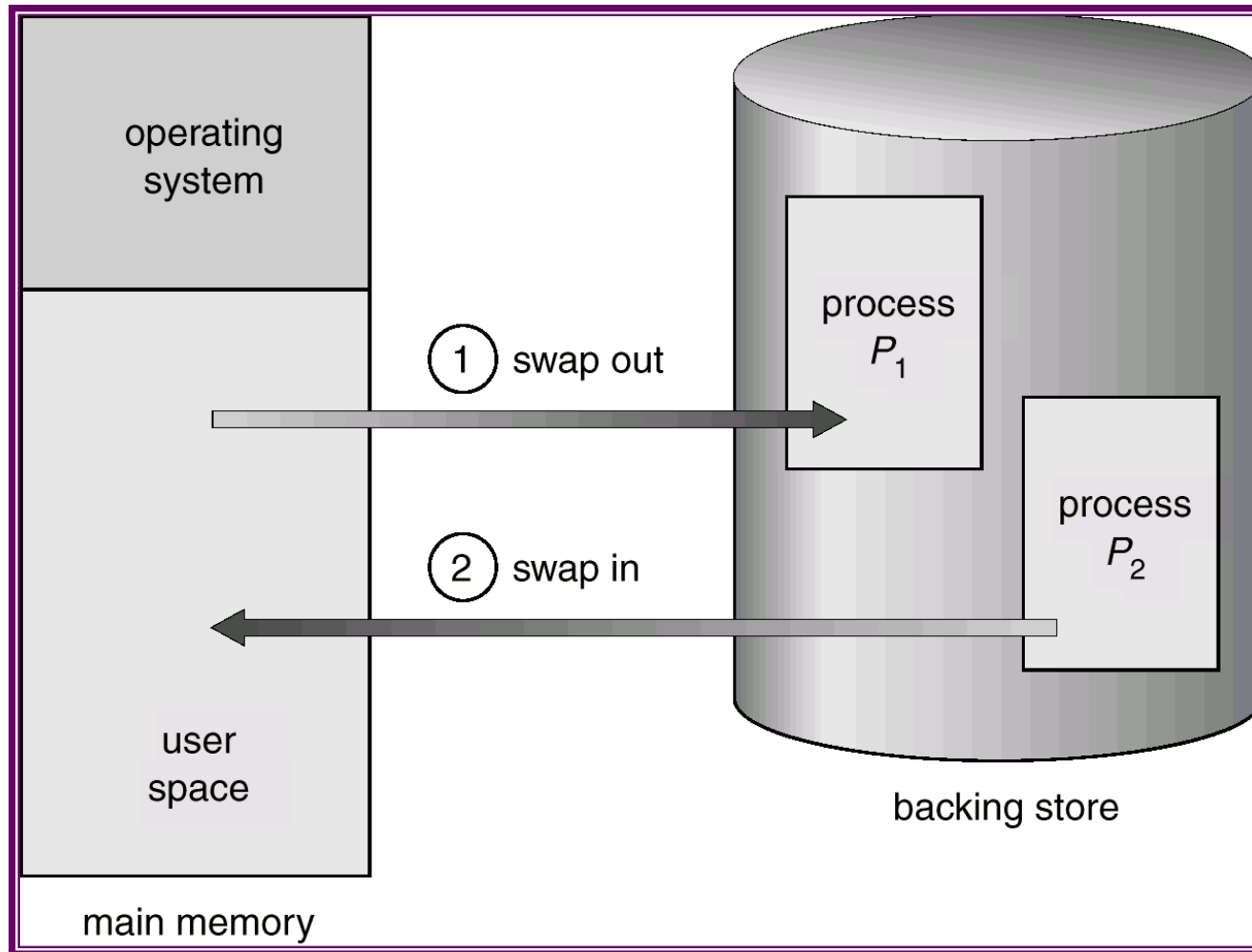
# Overlays for a Two-Pass Assembler



# Swapping

- A process can be *swapped* temporarily out of memory to a *backing store*, and then brought back into memory for continued execution.
- Backing store – fast disk large enough to accommodate copies of all memory images for all users; must provide direct access to these memory images.
- *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.
- Major part of swap time is transfer time; total transfer time is directly proportional to the *amount* of memory swapped.
- Modified versions of swapping are found on many systems, i.e., UNIX, Linux, and Windows.

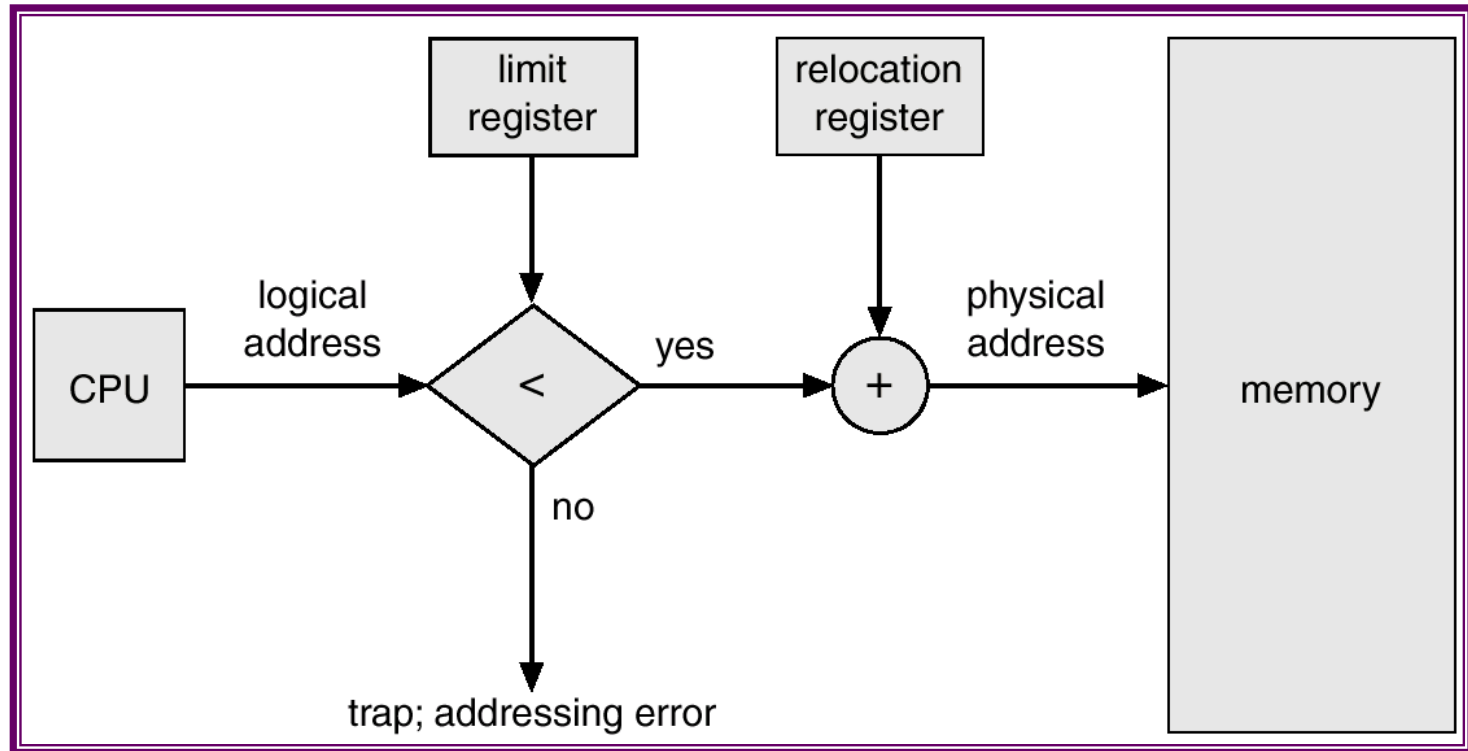
# Schematic View of Swapping



# Contiguous Allocation

- Main memory usually into two partitions:
  - ☞ Resident operating system, usually held in low memory with interrupt vector.
  - ☞ User processes then held in high memory.
- Single-partition allocation
  - ☞ Relocation-register scheme used to protect user processes from each other, and from changing operating-system code and data.
  - ☞ Relocation register contains value of smallest physical address; limit register contains range of logical addresses – each logical address must be less than the limit register.

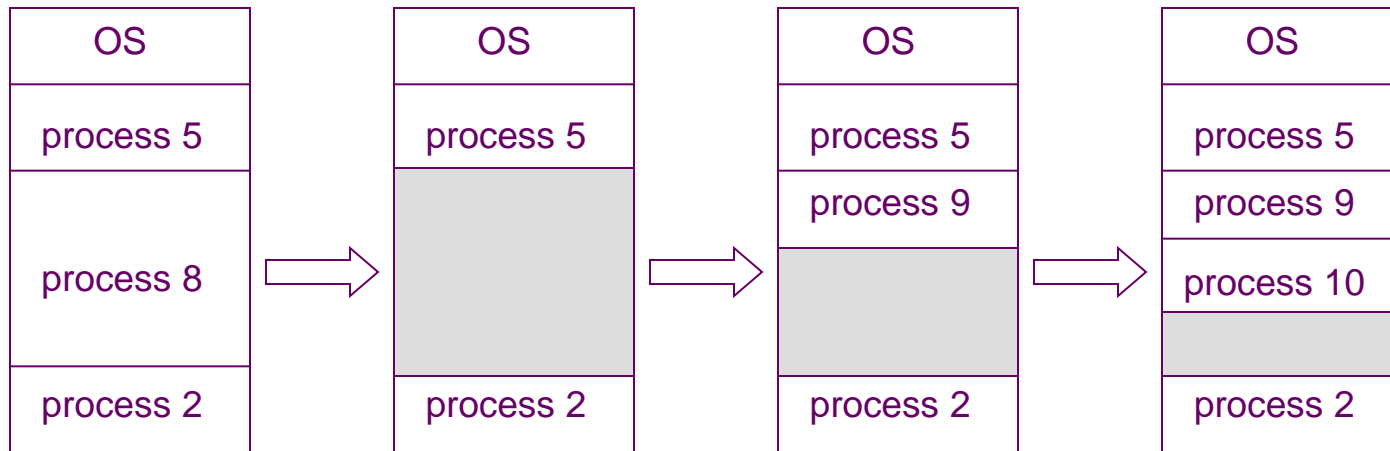
# Hardware Support for Relocation and Limit Registers



# Contiguous Allocation (Cont.)

## ■ Multiple-partition allocation

- ☞ *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.
- ☞ Operating system maintains information about:  
a) allocated partitions    b) free partitions (hole)



# Dynamic Storage-Allocation Problem

How to satisfy a request of size  $n$  from a list of free holes.

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

First-fit and best-fit better than worst-fit in terms of speed and storage utilization.



# Fragmentation

- **External Fragmentation** – total memory space exists to satisfy a request, but it is not contiguous.
- **Internal Fragmentation** – allocated memory may be slightly larger than requested memory; this size difference is memory internal to a partition, but not being used.
- Reduce external fragmentation by compaction
  - ☞ Shuffle memory contents to place all free memory together in one large block.
  - ☞ Compaction is possible *only* if relocation is dynamic, and is done at execution time.
  - ☞ I/O problem
    - 📄 Latch job in memory while it is involved in I/O.
    - 📄 Do I/O only into OS buffers.

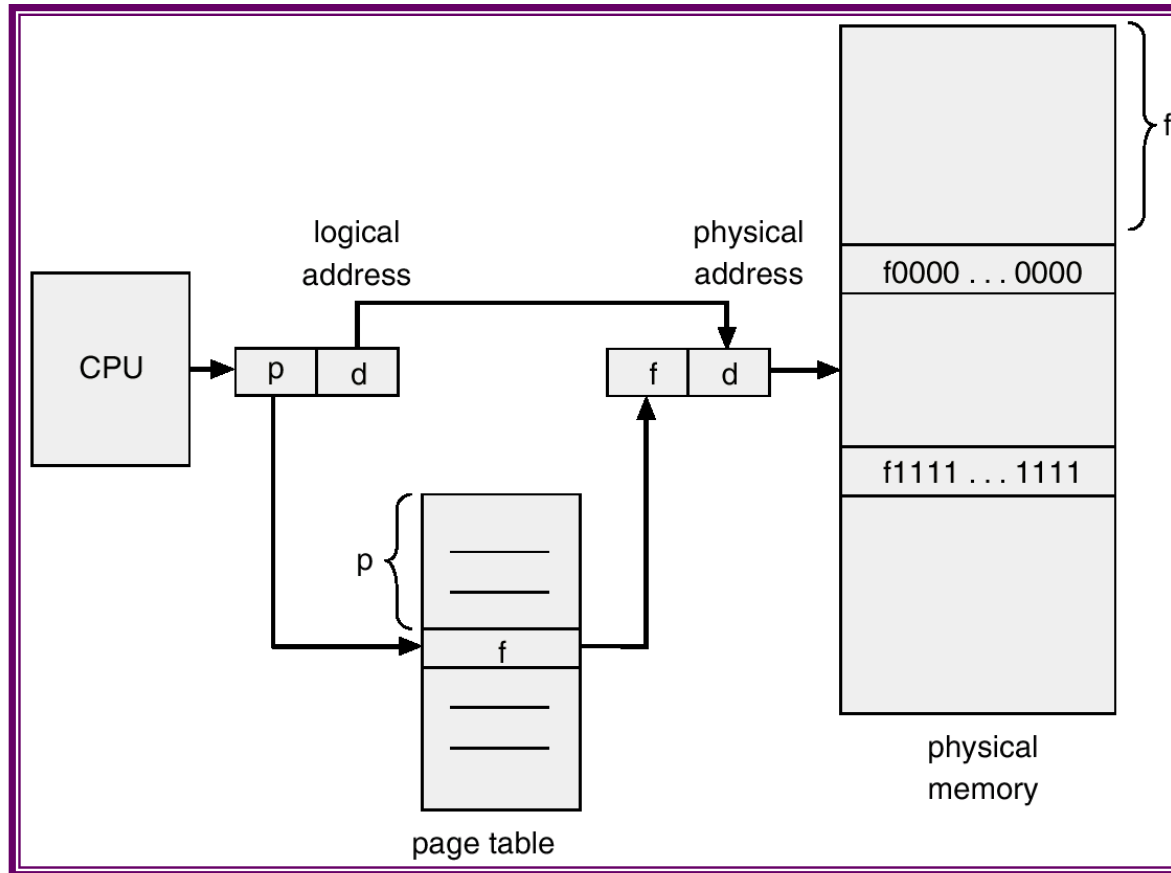
# Paging

- Logical address space of a process can be noncontiguous; process is allocated physical memory whenever the latter is available.
- Divide physical memory into fixed-sized blocks called **frames** (size is power of 2, between 512 bytes and 8192 bytes).
- Divide logical memory into blocks of same size called **pages**.
- Keep track of all free frames.
- To run a program of size  $n$  pages, need to find  $n$  free frames and load program.
- Set up a page table to translate logical to physical addresses.
- Internal fragmentation.

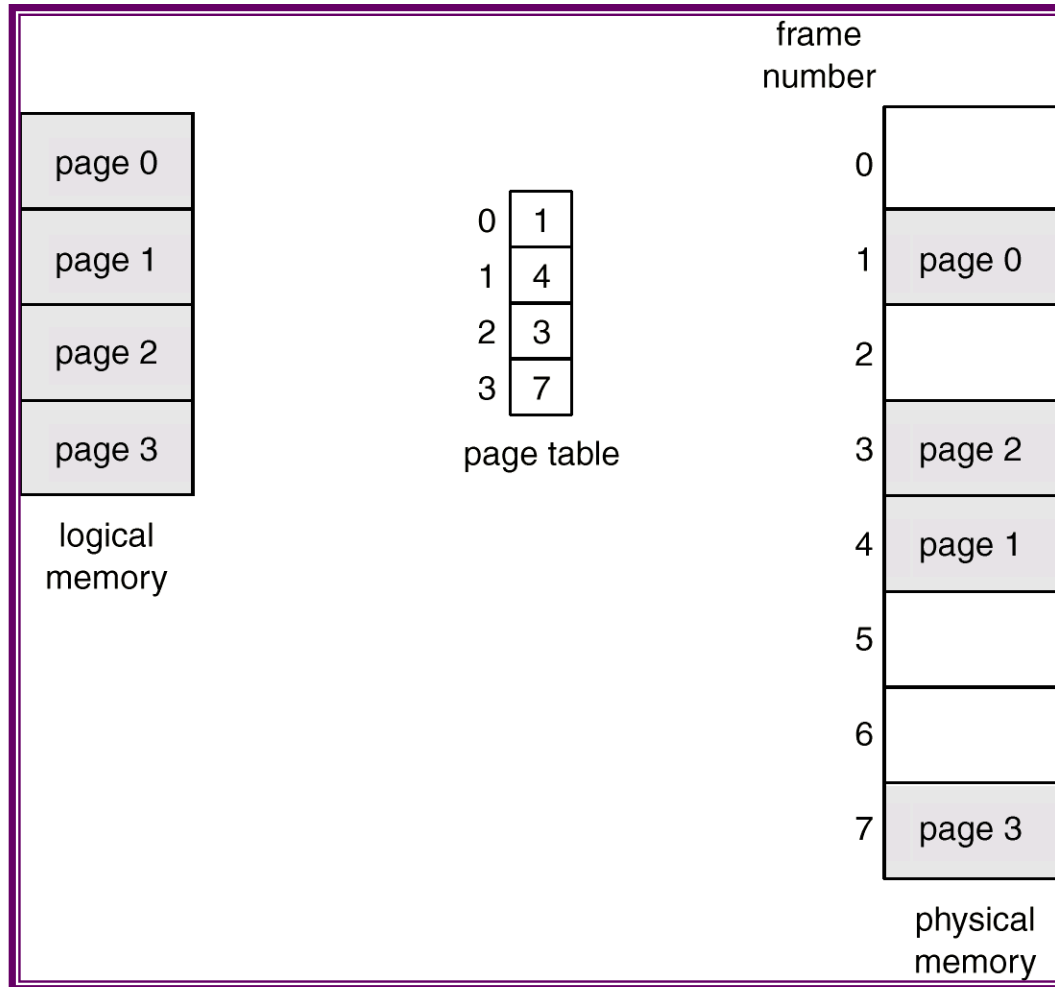
# Address Translation Scheme

- 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.

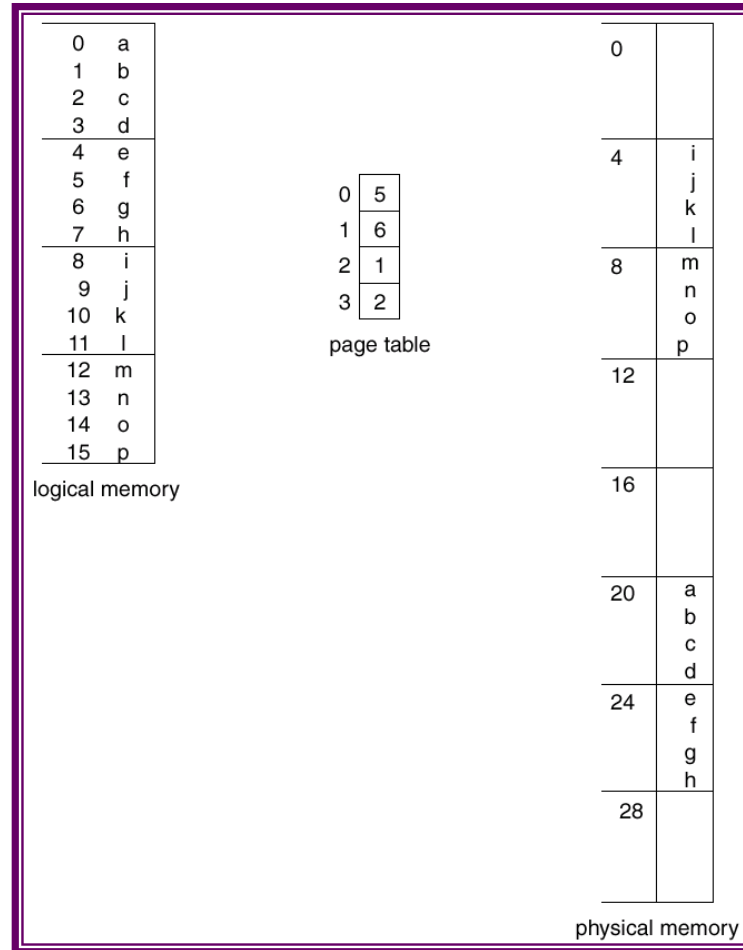
# Address Translation Architecture



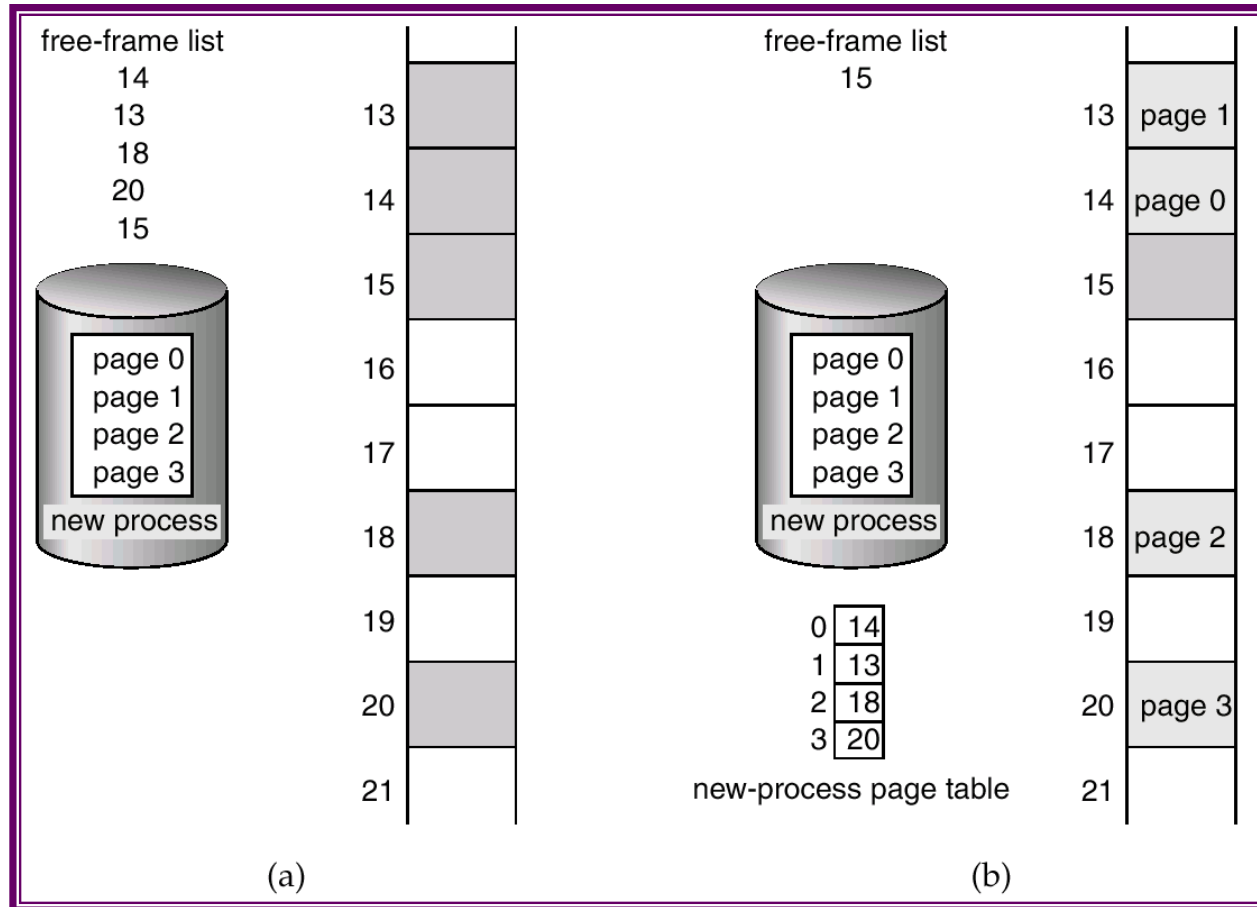
# Paging Example



# Paging Example



# Free Frames



Before allocation

After allocation

# Implementation of Page Table

- Page table is kept in main memory.
- *Page-table base register (PTBR)* points to the page table.
- *Page-table length register (PRLR)* indicates size of the page table.
- In this scheme every data/instruction access requires two memory accesses. One for the page table and one for the data/instruction.
- The two memory access problem can be solved by the use of a special fast-lookup hardware cache called *associative memory* or *translation look-aside buffers (TLBs)*



# Associative Memory

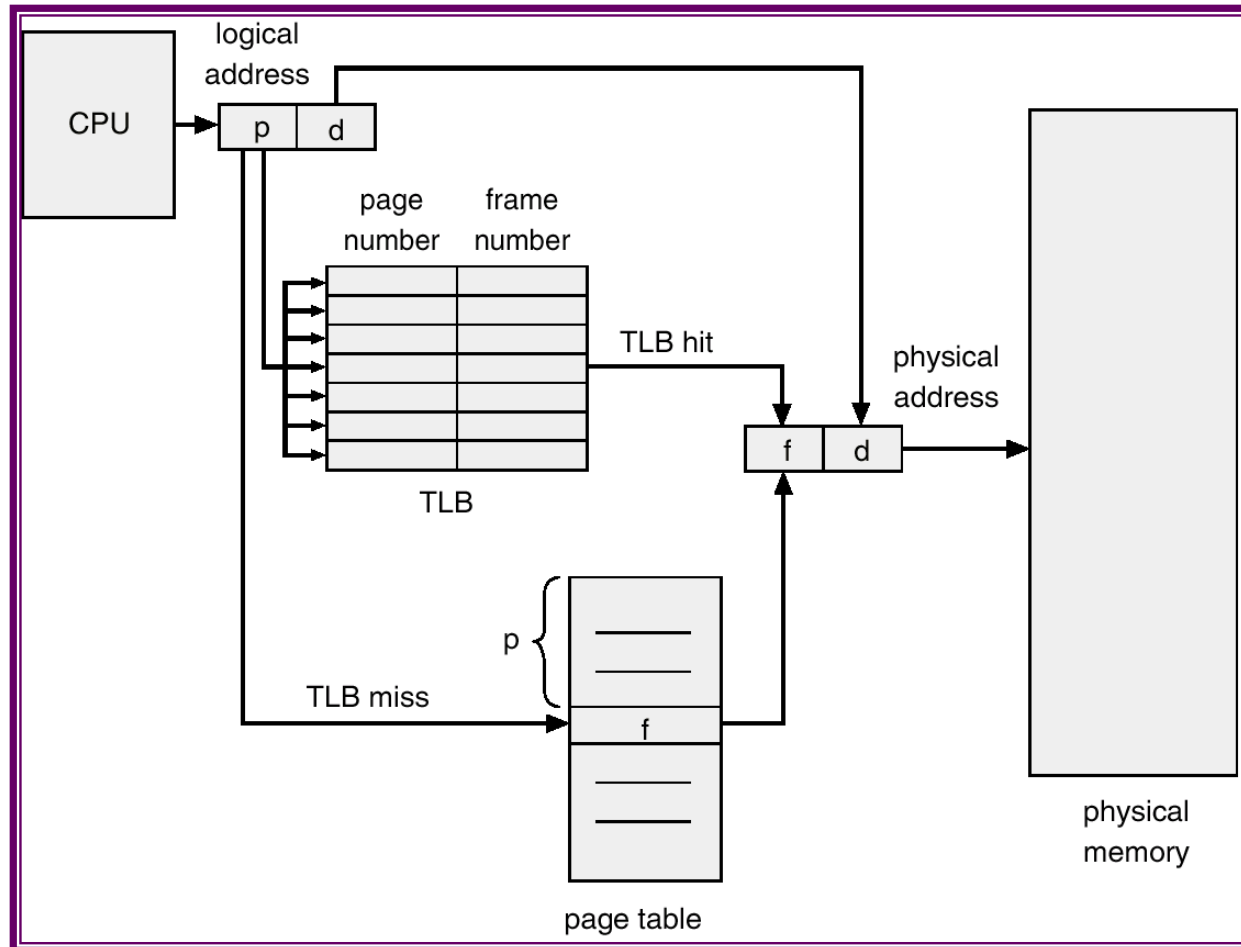
■ Associative memory – parallel search

Page #	Frame #

Address translation ( $A'$ ,  $A''$ )

- ☞ If  $A'$  is in associative register, get frame # out.
- ☞ Otherwise get frame # from page table in memory

# Paging Hardware With TLB



# Effective Access Time

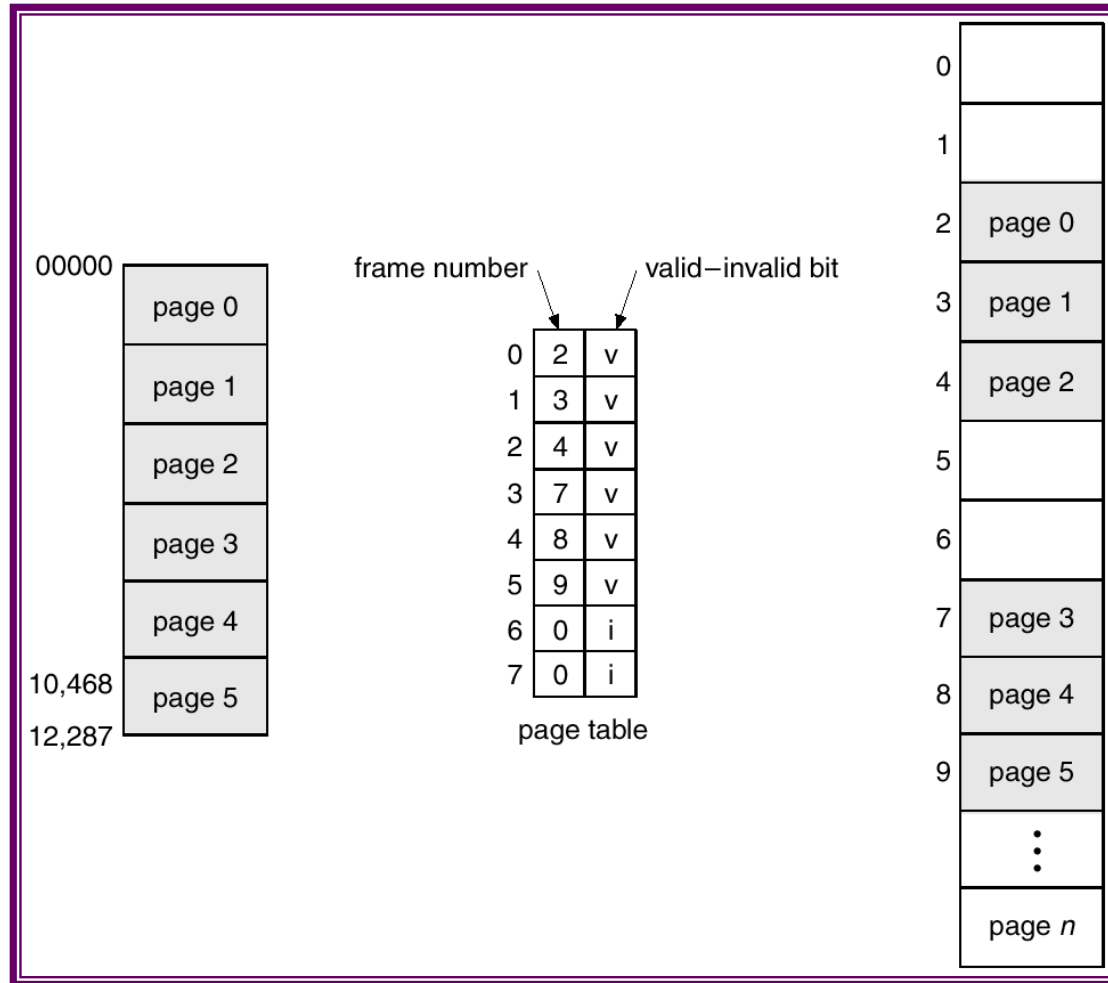
- Associative Lookup =  $\varepsilon$  time unit
- Assume memory cycle time is 1 microsecond
- Hit ratio – percentage of times that a page number is found in the associative registers; ration related to number of associative registers.
- Hit ratio =  $\alpha$
- Effective Access Time (EAT)

$$\begin{aligned} \text{EAT} &= (1 + \varepsilon) \alpha + (2 + \varepsilon)(1 - \alpha) \\ &= 2 + \varepsilon - \alpha \end{aligned}$$

# Memory Protection

- Memory protection implemented by associating protection bit with each frame.
- *Valid-invalid* bit attached to each entry in the page table:
  - ☞ “valid” indicates that the associated page is in the process’ logical address space, and is thus a legal page.
  - ☞ “invalid” indicates that the page is not in the process’ logical address space.

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



# Page Table Structure

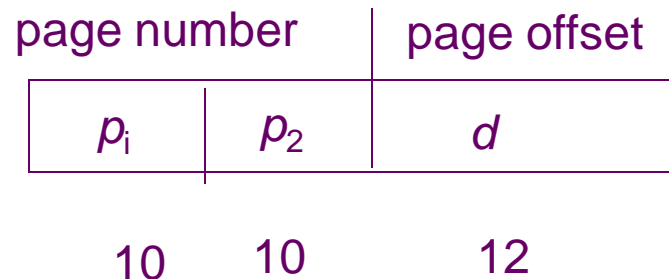
- Hierarchical Paging
- Hashed Page Tables
- Inverted Page Tables

# Hierarchical Page Tables

- Break up the logical address space into multiple page tables.
- A simple technique is a two-level page table.

# Two-Level Paging Example

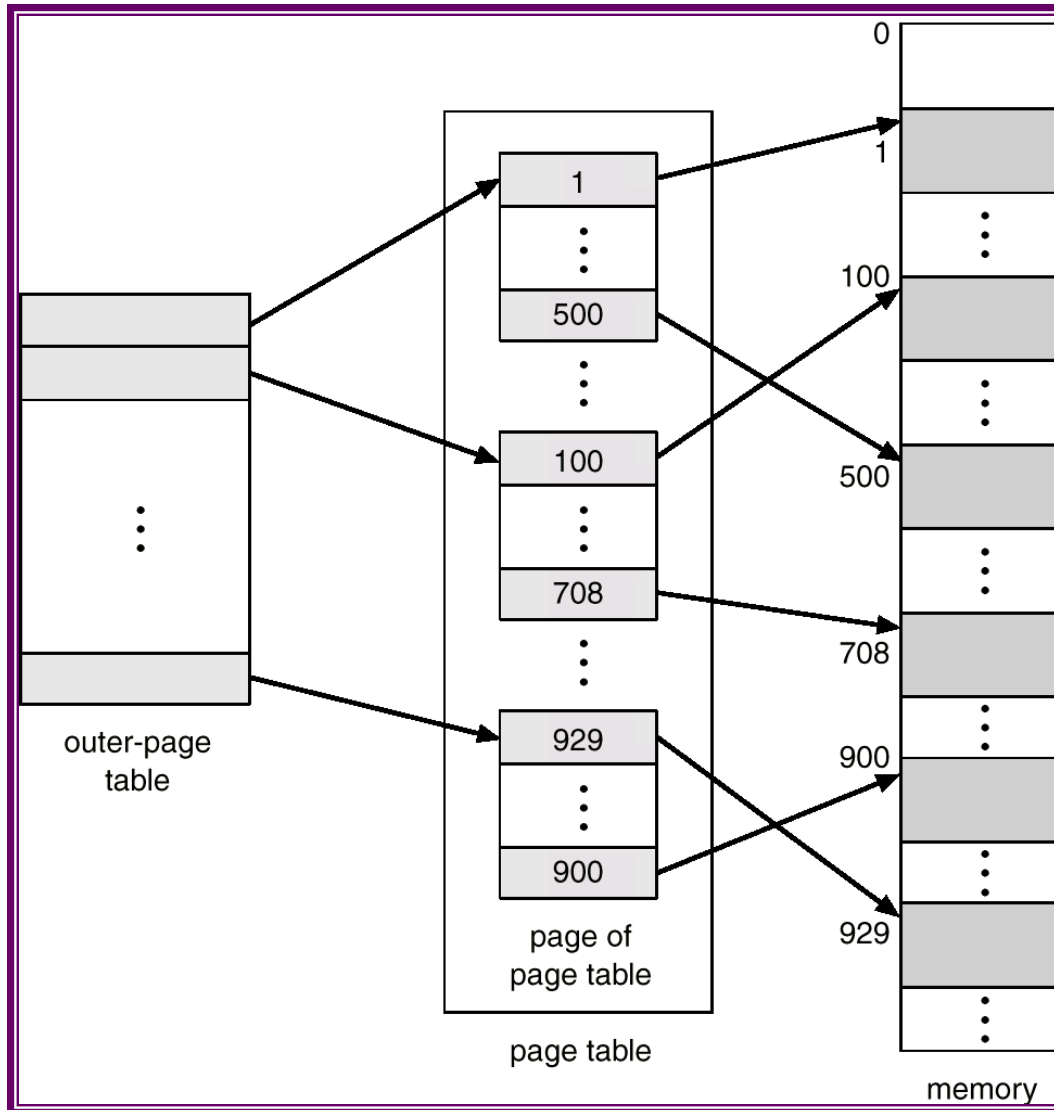
- A logical address (on 32-bit machine with 4K page size) is divided into:
  - ☞ a page number consisting of 20 bits.
  - ☞ a page offset consisting of 12 bits.
- Since the page table is paged, the page number is further divided into:
  - ☞ a 10-bit page number.
  - ☞ a 10-bit page offset.
- Thus, a logical address is as follows:



where  $p_1$  is an index into the outer page table, and  $p_2$  is the displacement within the page of the outer page table.

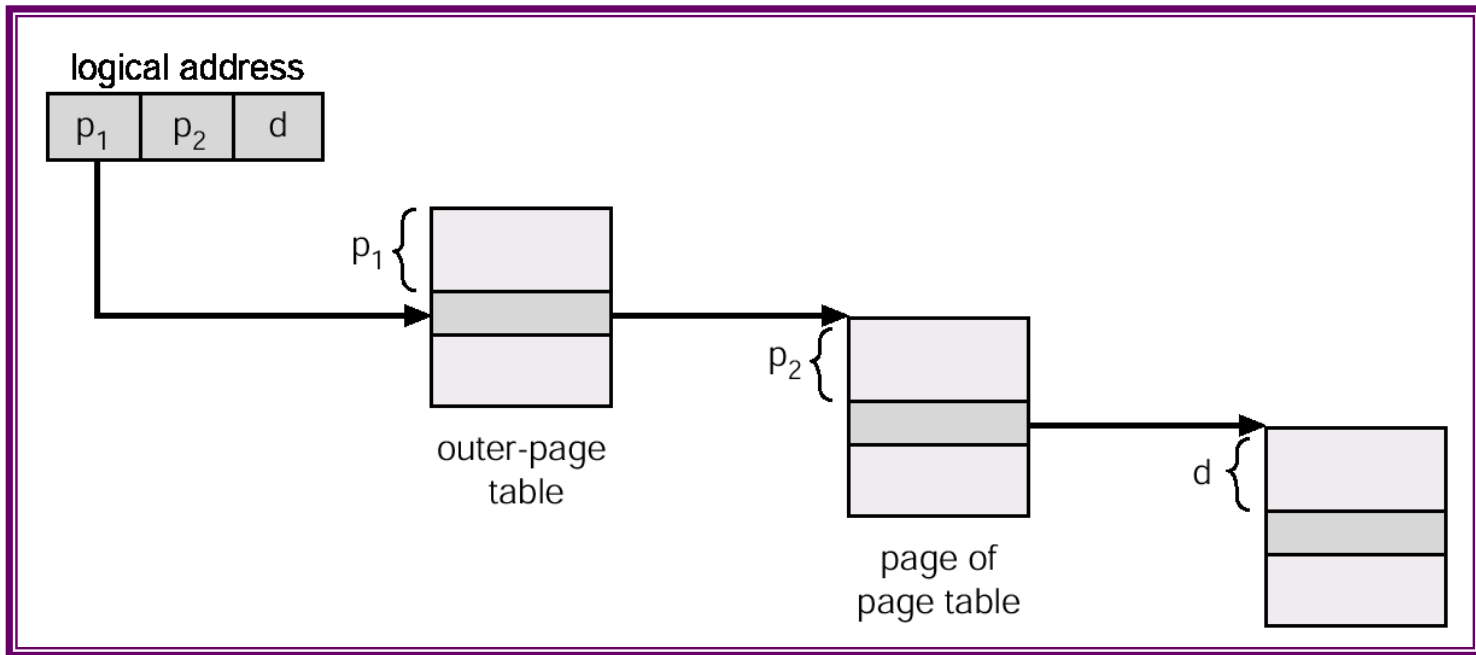


# Two-Level Page-Table Scheme



# Address-Translation Scheme

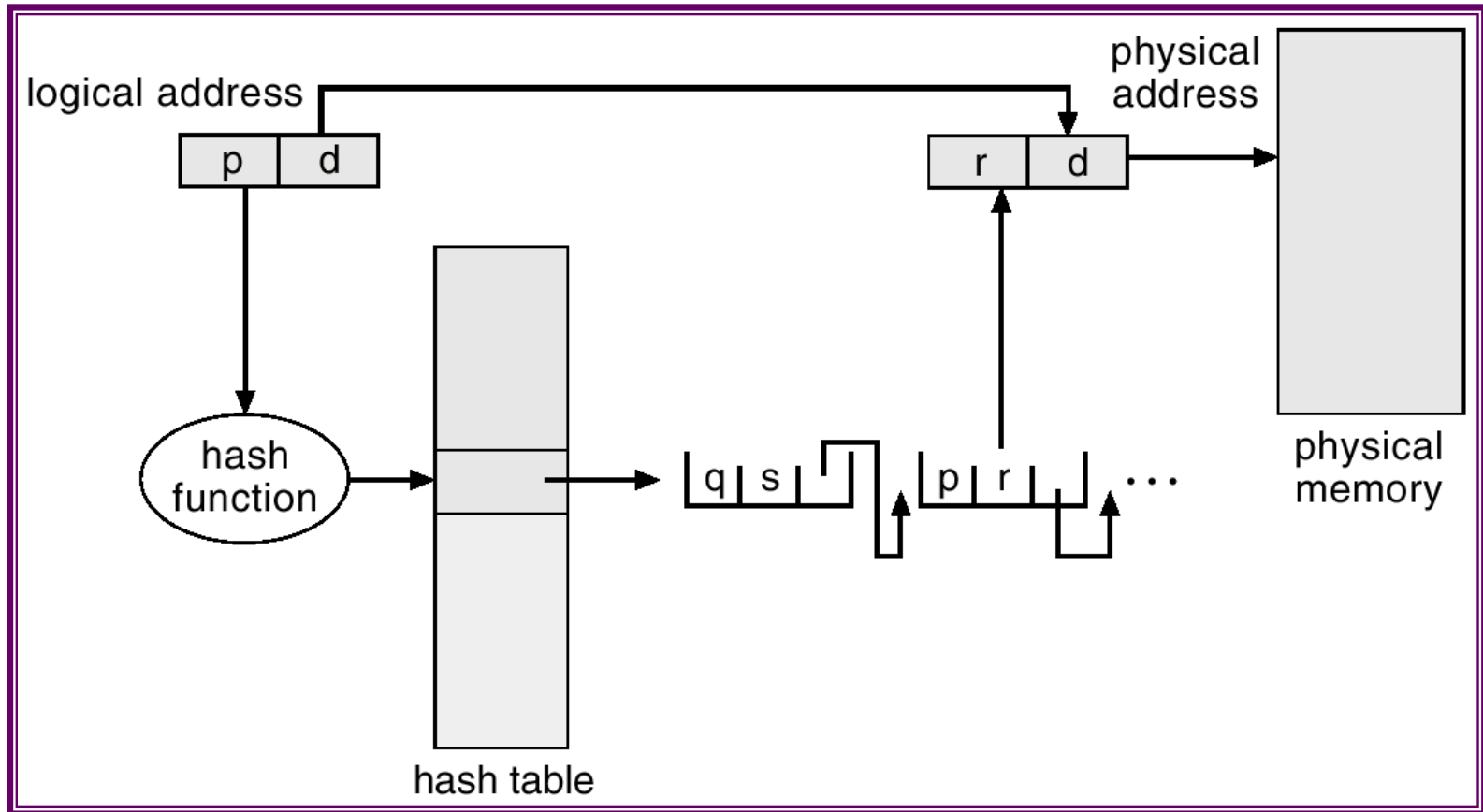
- Address-translation scheme for a two-level 32-bit paging architecture



# Hashed Page Tables

- Common in address spaces  $> 32$  bits.
- The virtual page number is hashed into a page table. This page table contains a chain of elements hashing to the same location.
- Virtual page numbers are compared in this chain searching for a match. If a match is found, the corresponding physical frame is extracted.

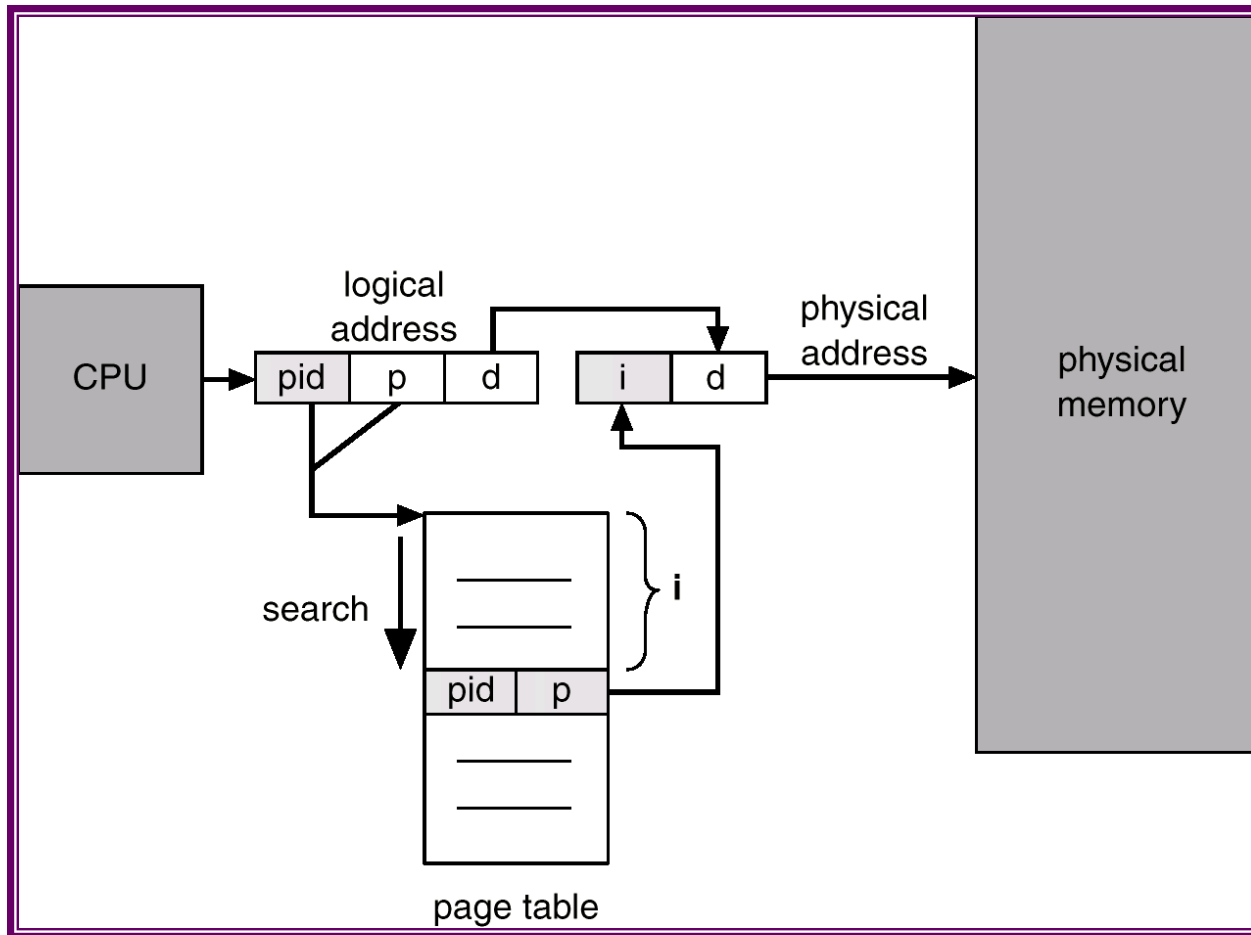
# Hashed Page Table



# Inverted Page Table

- One entry for each real page of memory.
- Entry consists of the virtual address of the page stored in that real memory location, with information about the process that owns that page.
- Decreases memory needed to store each page table, but increases time needed to search the table when a page reference occurs.
- Use hash table to limit the search to one — or at most a few — page-table entries.

# Inverted Page Table Architecture



# Shared Pages

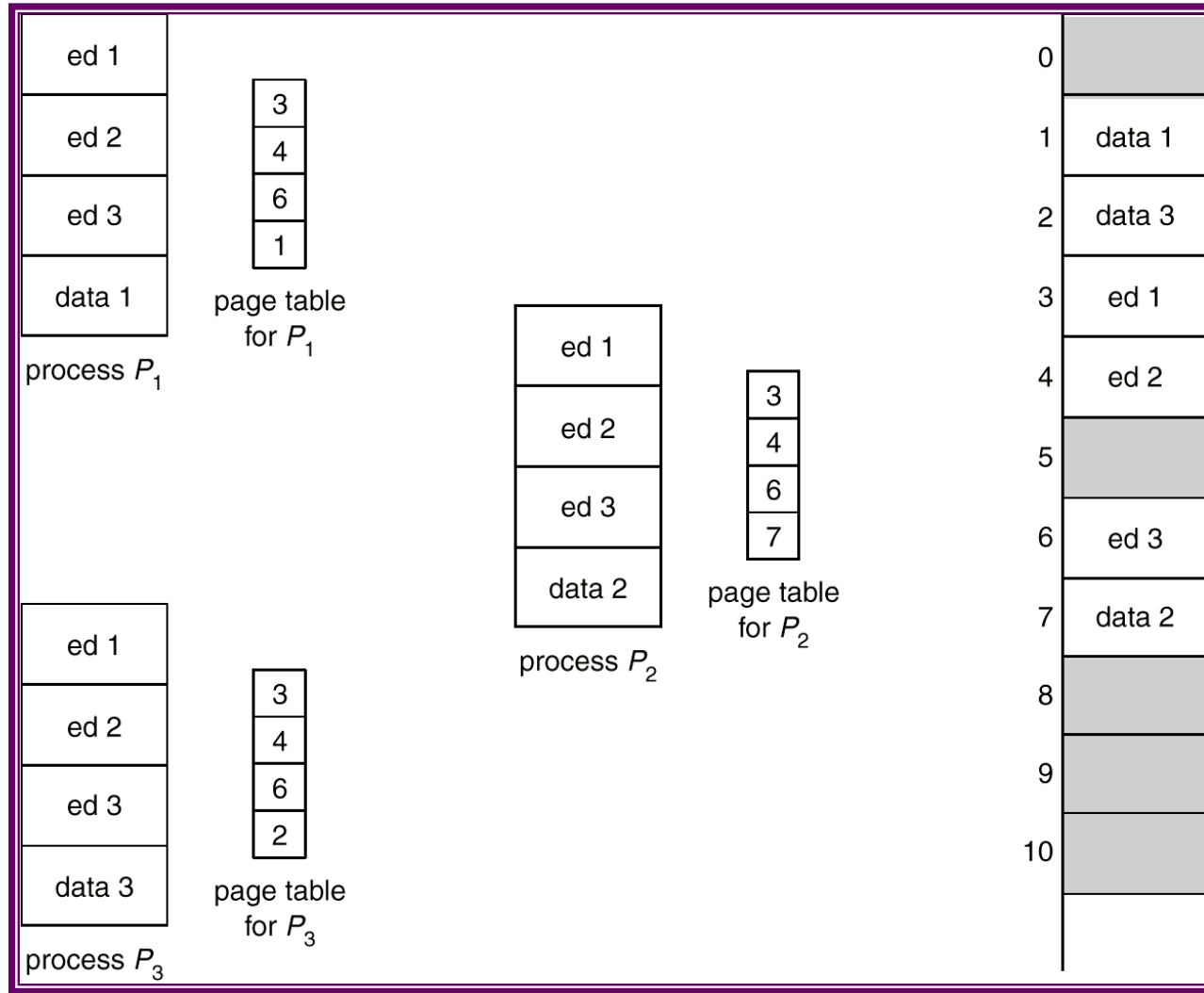
## ■ Shared code

- ☞ One copy of read-only (reentrant) code shared among processes (i.e., text editors, compilers, window systems).
- ☞ Shared code must appear in same location in the logical address space of all processes.

## ■ Private code and data

- ☞ Each process keeps a separate copy of the code and data.
- ☞ The pages for the private code and data can appear anywhere in the logical address space.

# Shared Pages Example

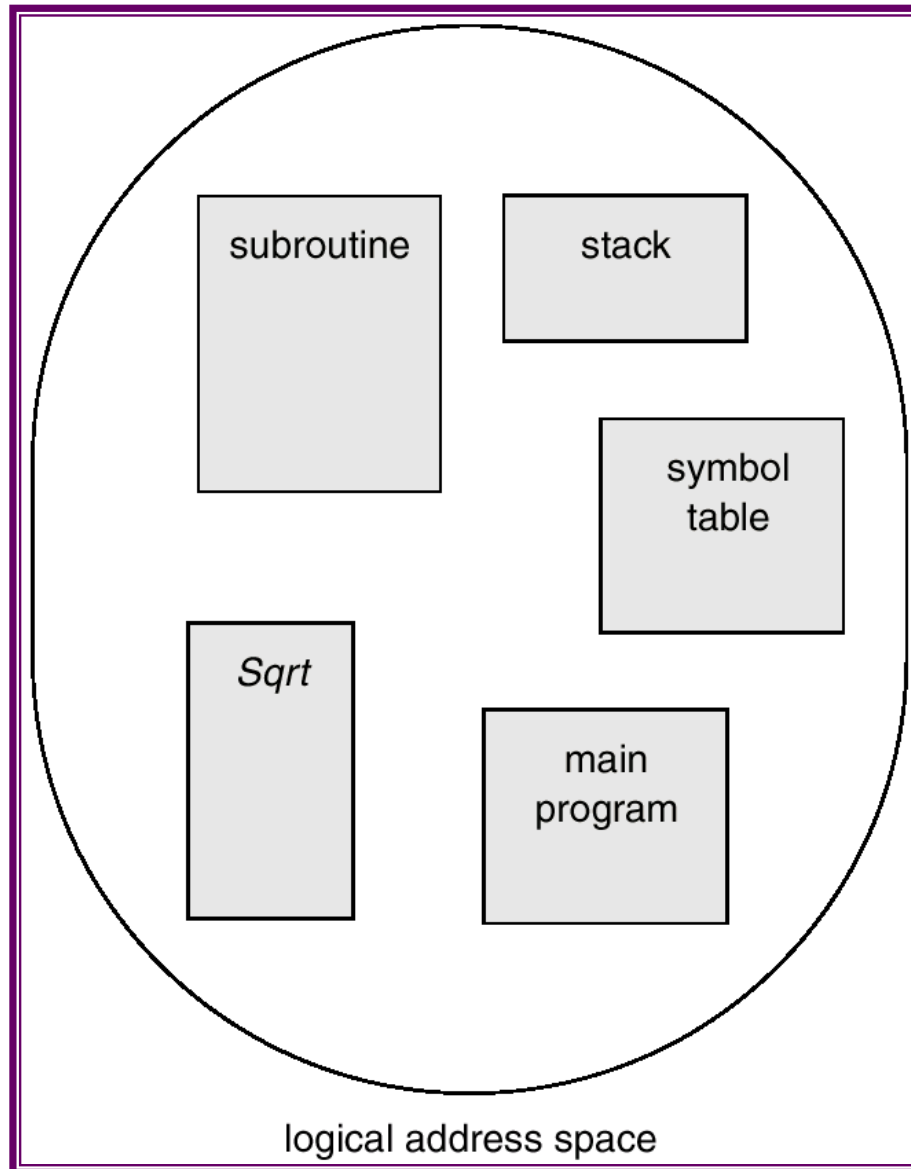




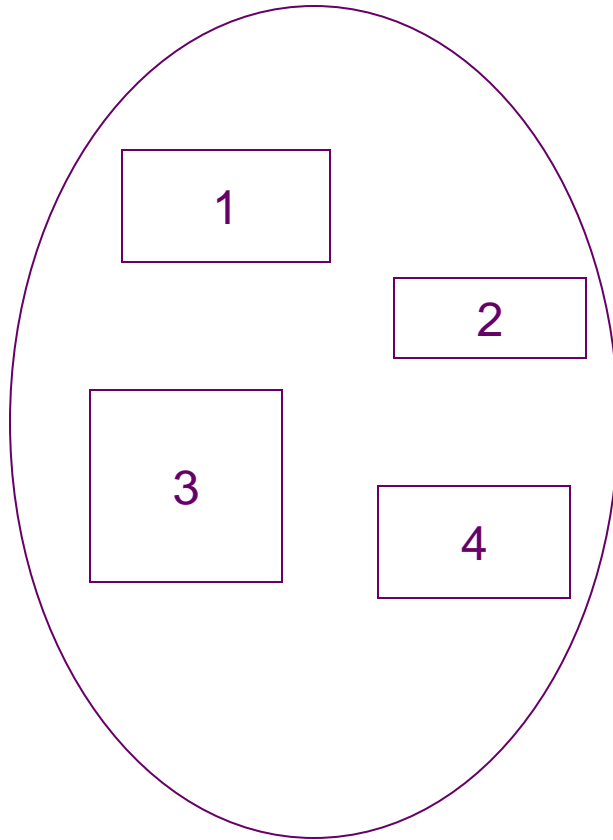
# Segmentation

- Memory-management scheme that supports user view of memory.
- A program is a collection of segments. A segment is a logical unit such as:
  - main program,
  - procedure,
  - function,
  - method,
  - object,
  - local variables, global variables,
  - common block,
  - stack,
  - symbol table, arrays

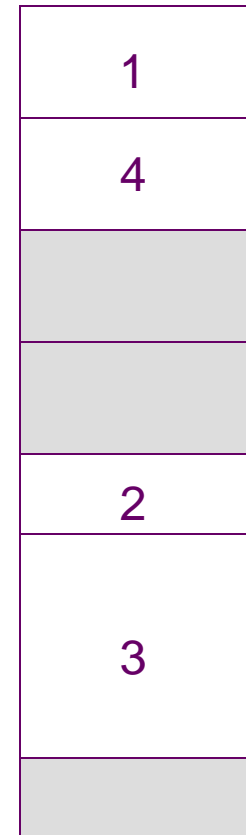
# User's View of a Program



# Logical View of Segmentation



user space



physical memory space

# Segmentation Architecture

- Logical address consists of a two tuple:  
    <segment-number, offset>,
- *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.
- *Segment-table length register (STLR)* indicates number of segments used by a program;  
    segment number  $s$  is legal if  $s < \text{STLR}$ .

# Segmentation Architecture (Cont.)

## ■ Relocation.

- ☞ dynamic
- ☞ by segment table

## ■ Sharing.

- ☞ shared segments
- ☞ same segment number

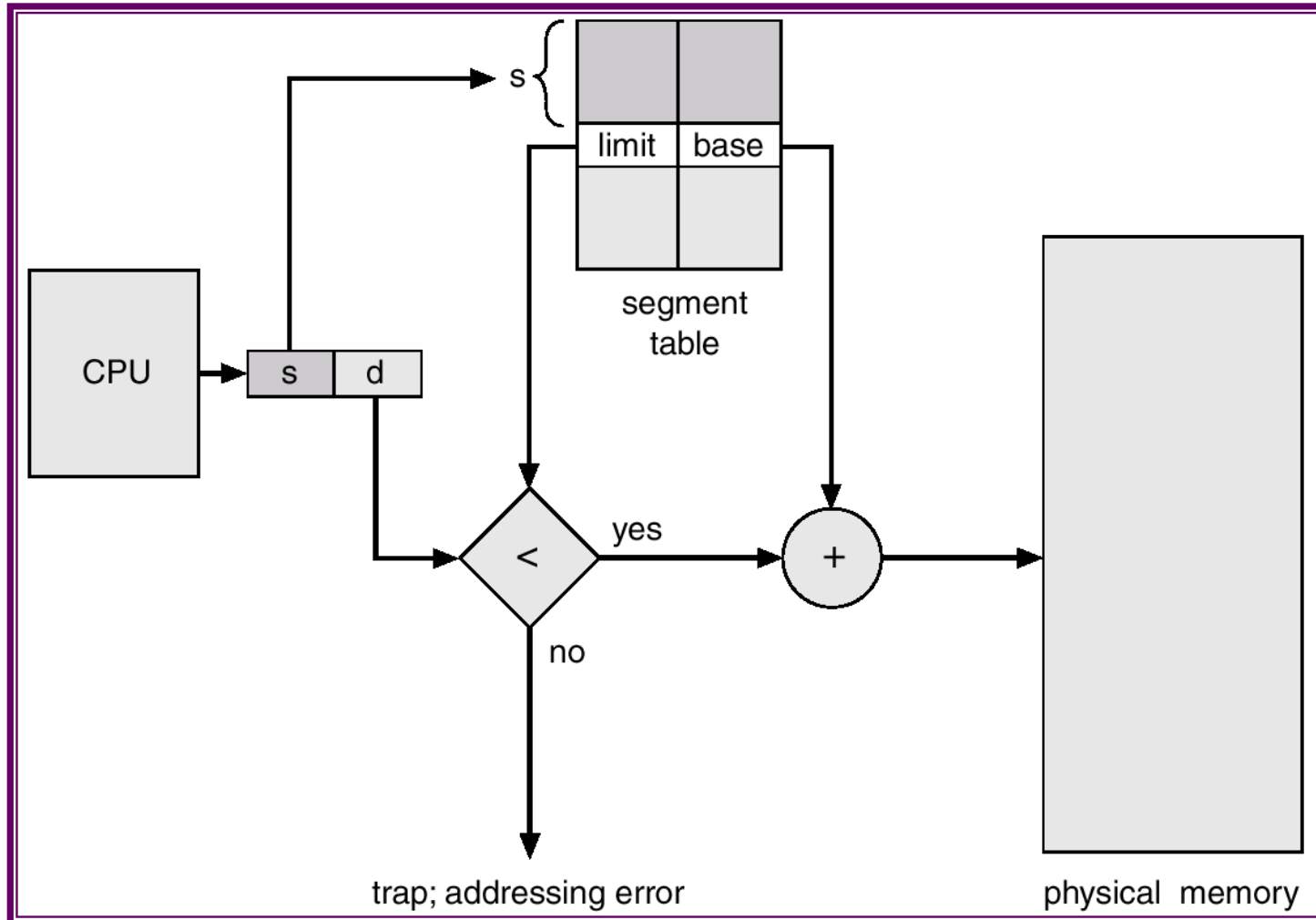
## ■ Allocation.

- ☞ first fit/best fit
- ☞ external fragmentation

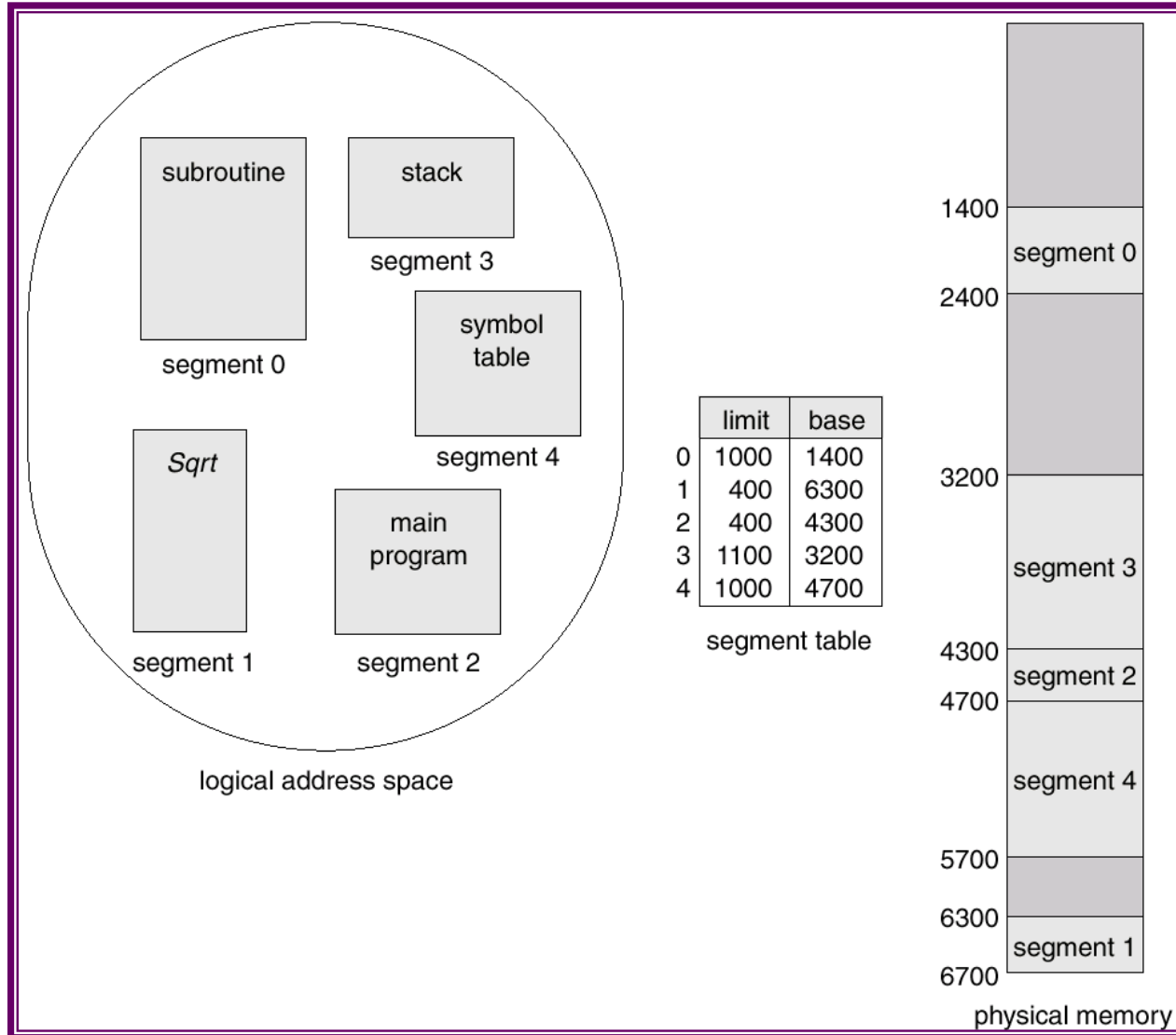
# Segmentation Architecture (Cont.)

- Protection. With each entry in segment table associate:
  - ☞ validation bit = 0  $\Rightarrow$  illegal segment
  - ☞ read/write/execute privileges
- Protection bits associated with segments; code sharing occurs at segment level.
- Since segments vary in length, memory allocation is a dynamic storage-allocation problem.
- A segmentation example is shown in the following diagram

# Segmentation Hardware

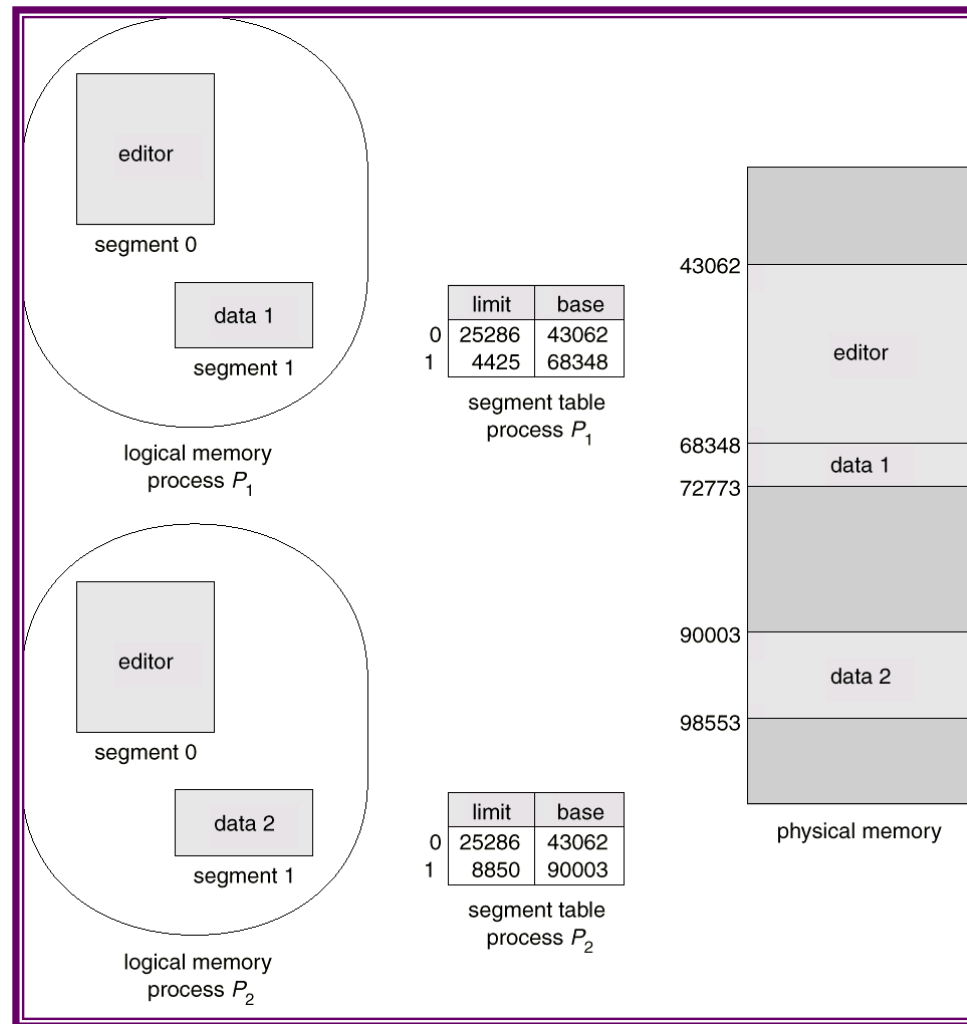


# Example of Segmentation





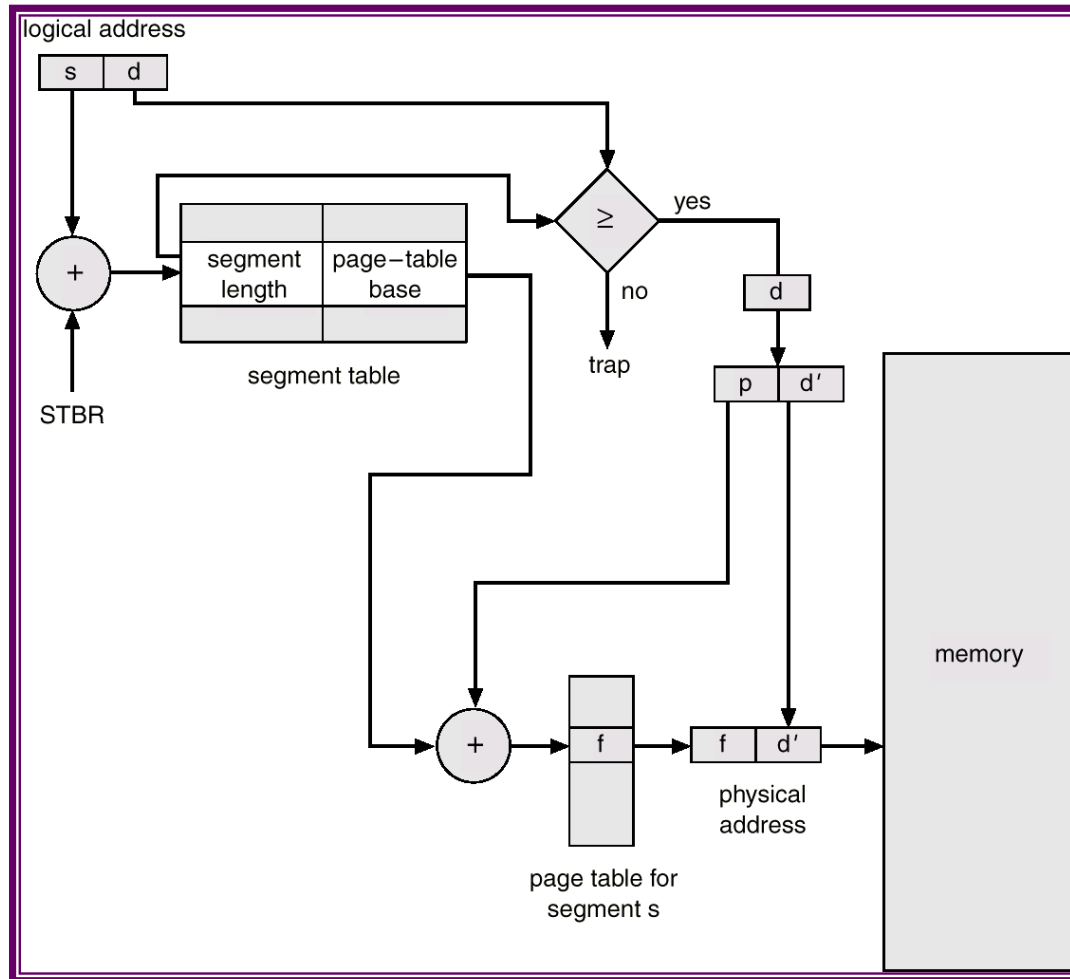
# Sharing of Segments



# Segmentation with Paging – MULTICS

- The MULTICS system solved problems of external fragmentation and lengthy search times by paging the segments.
- Solution differs from pure segmentation in that the segment-table entry contains not the base address of the segment, but rather the base address of a *page table* for this segment.

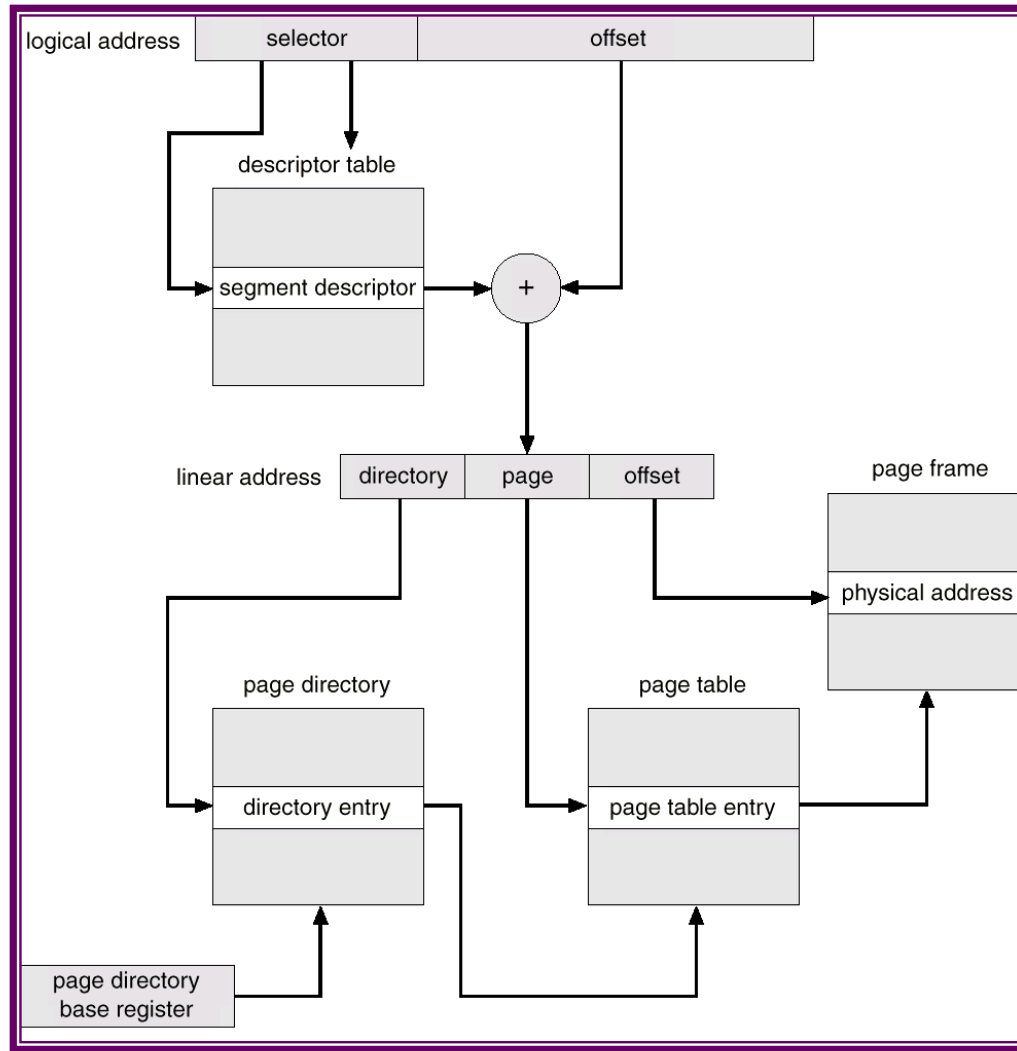
# MULTICS Address Translation Scheme



# Segmentation with Paging – Intel 386

- As shown in the following diagram, the Intel 386 uses segmentation with paging for memory management with a two-level paging scheme.

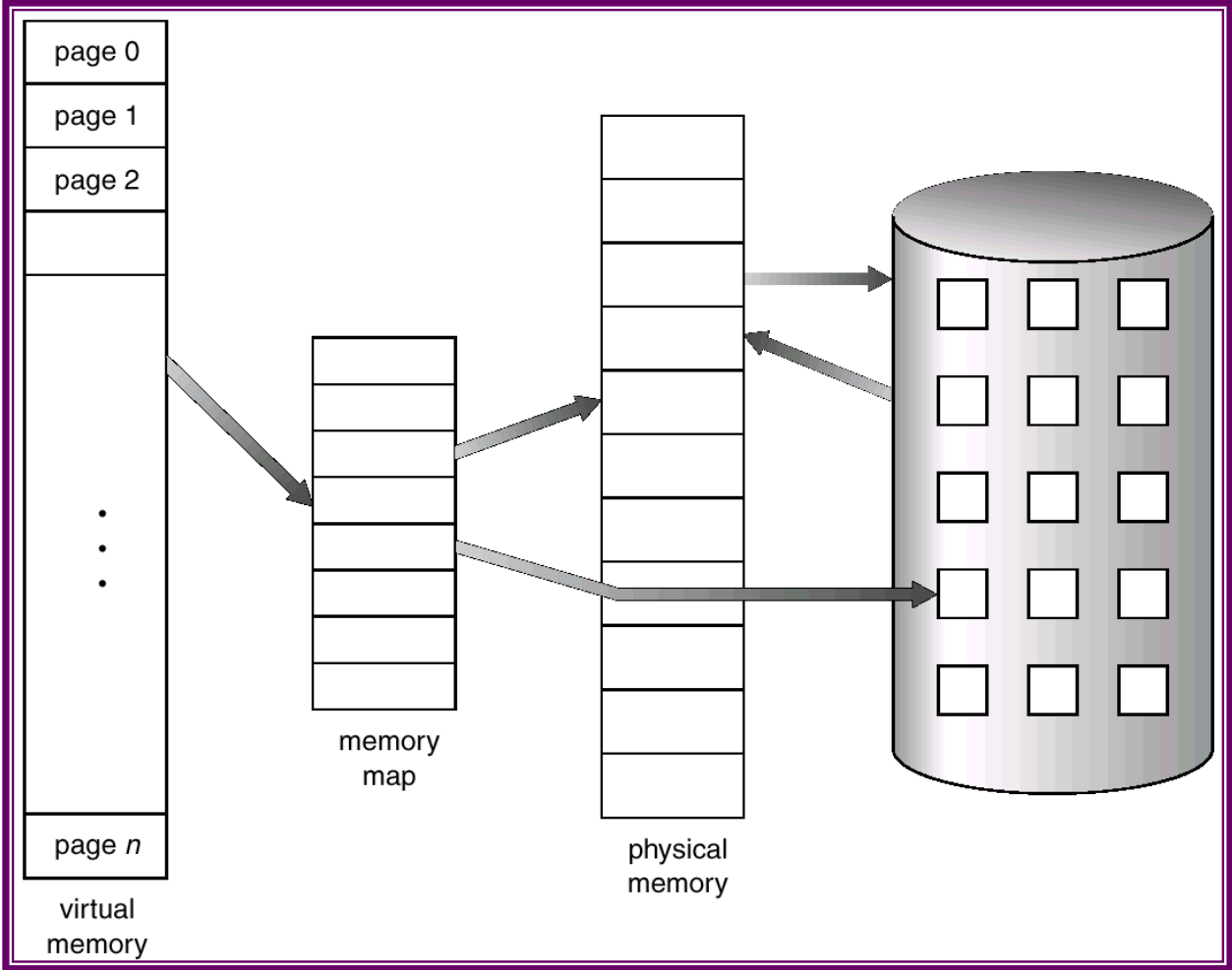
# Intel 30386 Address Translation



# Background

- **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 address space.
  - ☞ Allows address spaces to be shared by several processes.
  - ☞ Allows for more efficient process creation.
- Virtual memory can be implemented via:
  - ☞ Demand paging
  - ☞ Demand segmentation

# Virtual Memory That is Larger Than Physical Memory

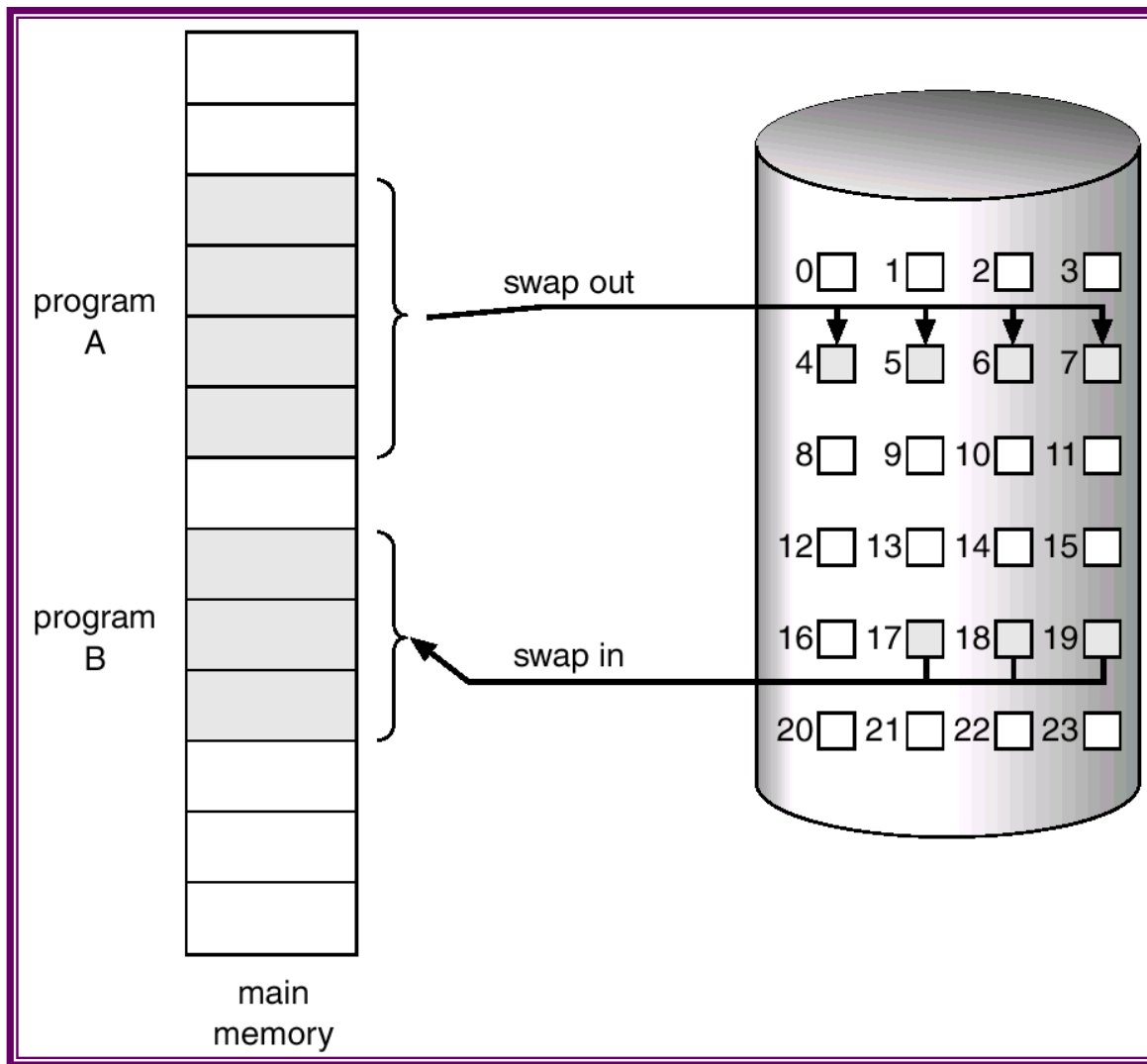


# Demand Paging

- Bring a page into memory only when it is needed.
  - ☞ Less I/O needed
  - ☞ Less memory needed
  - ☞ Faster response
  - ☞ More users
- Page is needed  $\Rightarrow$  reference to it
  - ☞ invalid reference  $\Rightarrow$  abort
  - ☞ not-in-memory  $\Rightarrow$  bring to memory



# Transfer of a Paged Memory to Contiguous Disk Space



# Valid-Invalid Bit

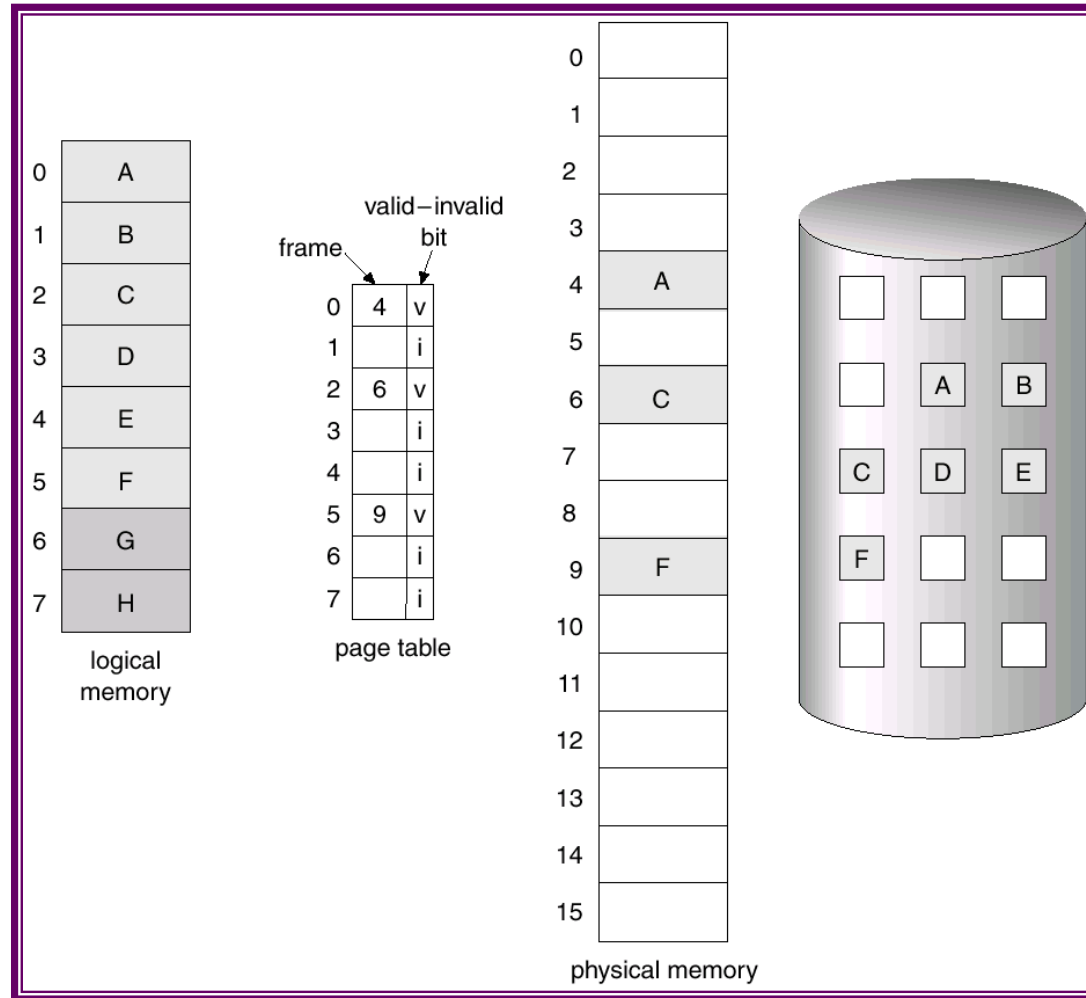
- With each page table entry a valid–invalid bit is associated  
(1  $\Rightarrow$  in-memory, 0  $\Rightarrow$  not-in-memory)
- Initially valid–invalid bit is set to 0 on all entries.
- Example of a page table snapshot.

Frame #	valid-invalid bit
	1
	1
	1
	1
	0
⋮	
	0
	0

page table

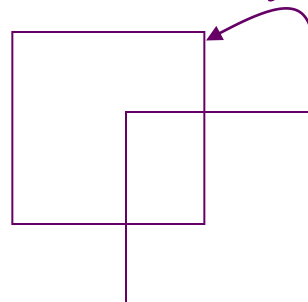
- During address translation, if valid–invalid bit in page table entry is 0  $\Rightarrow$  page fault.

# Page Table When Some Pages Are Not in Main Memory



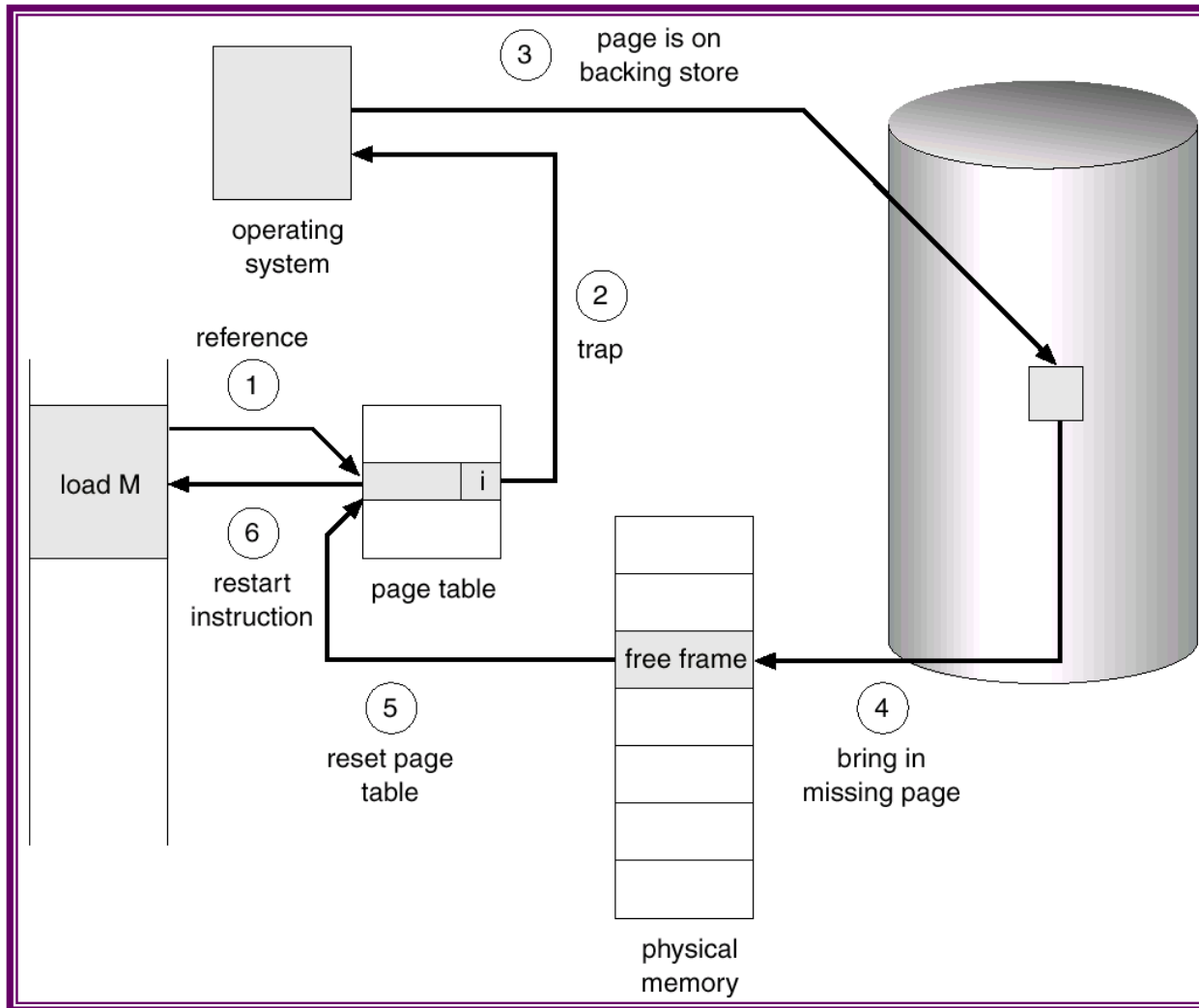
# Page Fault

- If there is ever a reference to a page, first reference will trap to OS  $\Rightarrow$  page fault
- OS looks at another table to decide:
  - ☞ Invalid reference  $\Rightarrow$  abort.
  - ☞ Just not in memory.
- Get empty frame.
- Swap page into frame.
- Reset tables, validation bit = 1.
- Restart instruction: Least Recently Used
  - ☞ block move



- ☞ auto increment/decrement location

# Steps in Handling a Page Fault



# What happens if there is no free frame?

- Page replacement – find some page in memory, but not really in use, swap it out.
  - ☞ algorithm
  - ☞ performance – want an algorithm which will result in minimum number of page faults.
- Same page may be brought into memory several times.

# Performance of Demand Paging

- Page Fault Rate  $0 \leq p \leq 1.0$

- ☞ if  $p = 0$  no page faults

- ☞ if  $p = 1$ , every reference is a fault

- Effective Access Time (EAT)

$$\begin{aligned} \text{EAT} = & (1 - p) \times \text{memory access} \\ & + p (\text{page fault overhead} \\ & + [\text{swap page out}] \\ & + \text{swap page in} \\ & + \text{restart overhead}) \end{aligned}$$

# Demand Paging Example

- Memory access time = 1 microsecond
- 50% of the time the page that is being replaced has been modified and therefore needs to be swapped out.

- Swap Page Time = 10 msec = 10,000 msec

$$EAT = (1 - p) \times 1 + p (15000)$$

$$1 + 15000P \quad (\text{in msec})$$



# Process Creation

- Virtual memory allows other benefits during process creation:
  - Copy-on-Write
  - Memory-Mapped Files

# Copy-on-Write

- Copy-on-Write (COW) allows both parent and child processes to initially *share* the same pages in memory.

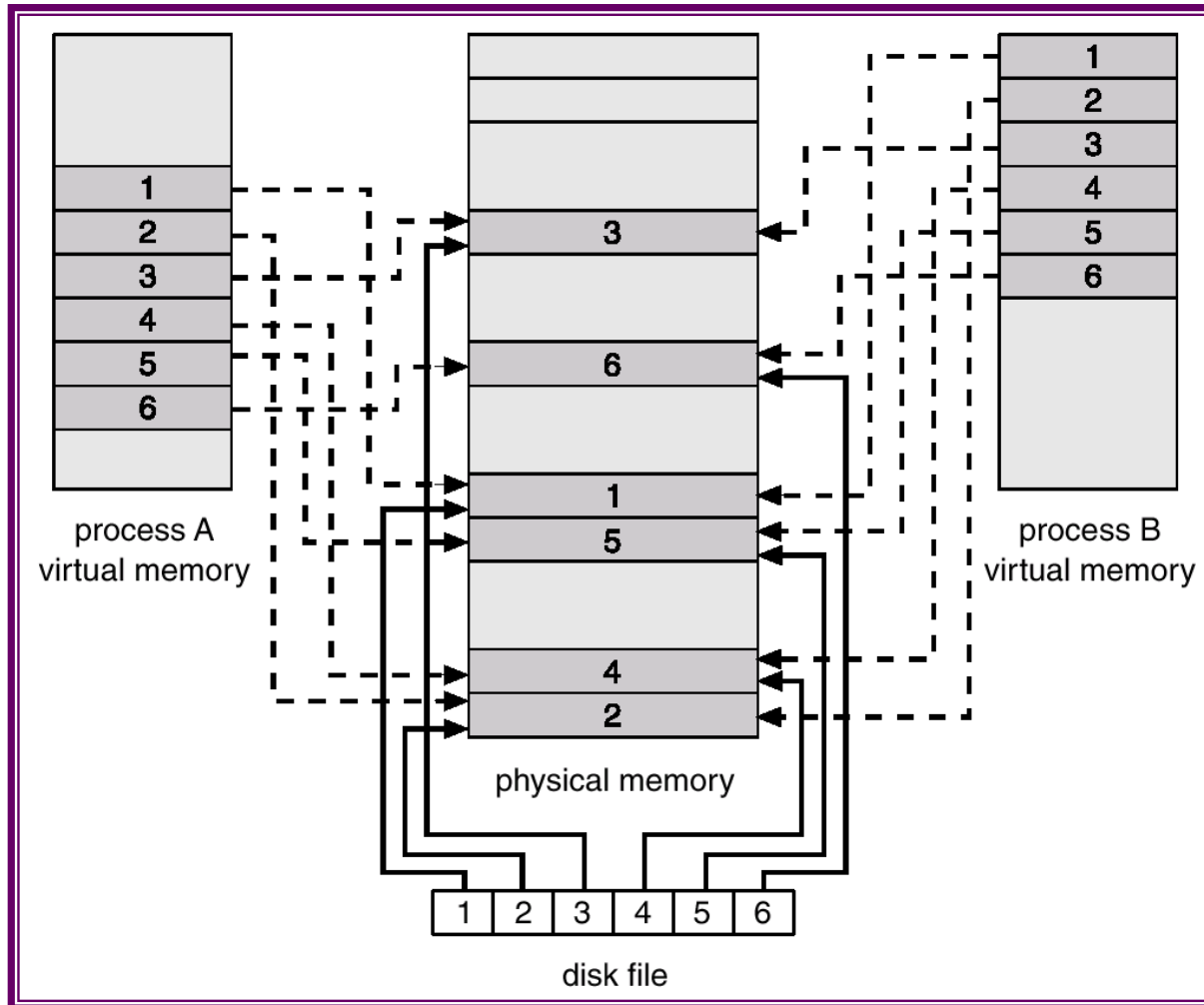
If either process modifies a shared page, only then is the page copied.

- COW allows more efficient process creation as only modified pages are copied.
- Free pages are allocated from a *pool* of zeroed-out pages.

# Memory-Mapped Files

- Memory-mapped file I/O allows file I/O to be treated as routine memory access by *mapping* a disk block to a page in memory.
- A file is initially read using demand paging. A page-sized portion of the file is read from the file system into a physical page. Subsequent reads/writes to/from the file are treated as ordinary memory accesses.
- Simplifies file access by treating file I/O through memory rather than **read()** **write()** system calls.
- Also allows several processes to map the same file allowing the pages in memory to be shared.

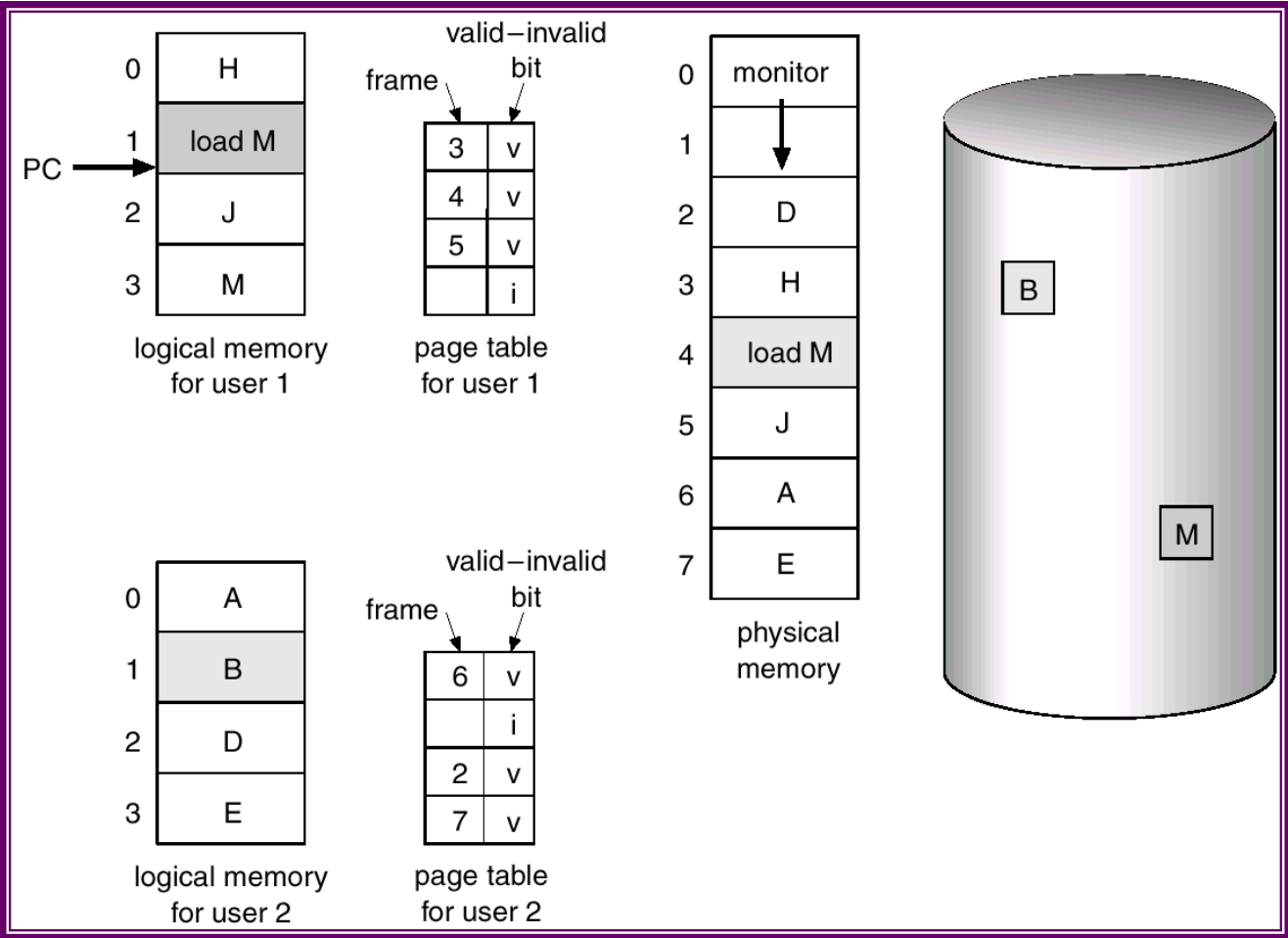
# Memory Mapped Files



# Page Replacement

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

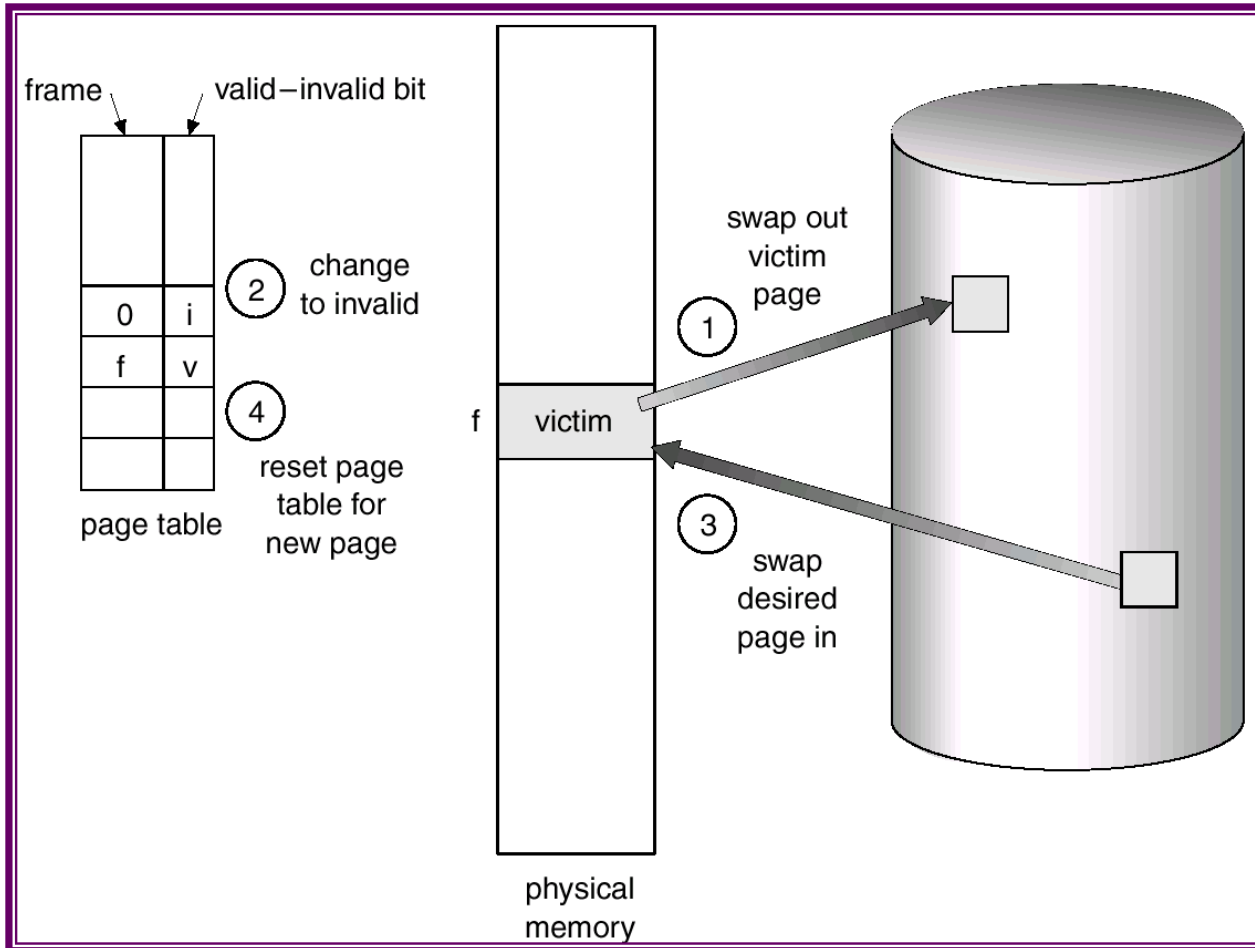
# Need For Page Replacement



# Basic Page Replacement

1. Find the location of the desired page on disk.
2. Find a free frame:
  - If there is a free frame, use it.
  - If there is no free frame, use a page replacement algorithm to select a *victim* frame.
3. Read the desired page into the (newly) free frame.  
Update the page and frame tables.
4. Restart the process.

# Page Replacement

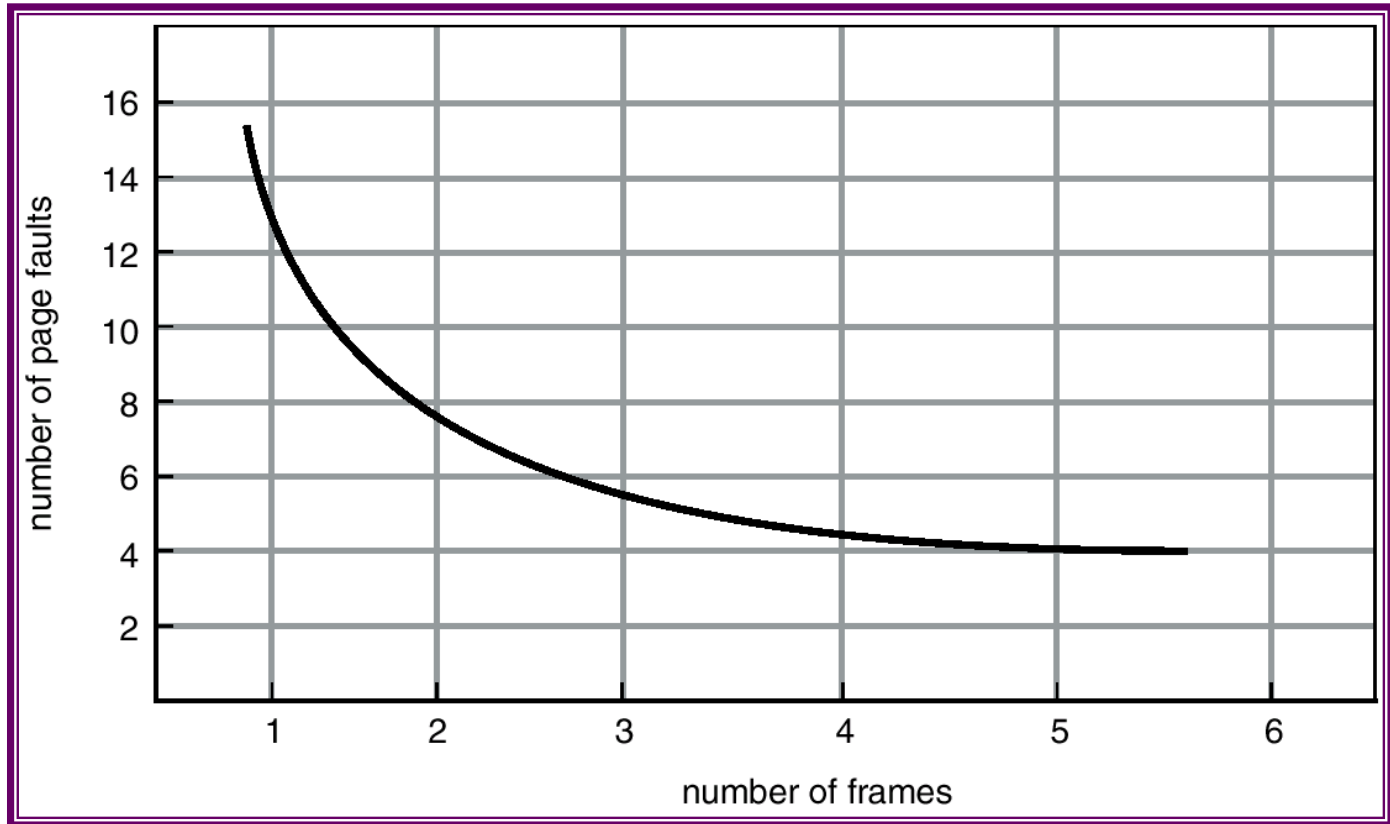




# Page Replacement Algorithms

- Want lowest page-fault rate.
- Evaluate algorithm by running it on a particular string of memory references (reference string) and computing the number of page faults on that string.
- In all our examples, the reference string is  
1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5.

# Graph of Page Faults Versus The Number of Frames



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

- Reference string: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5
- 3 frames (3 pages can be in memory at a time per process)

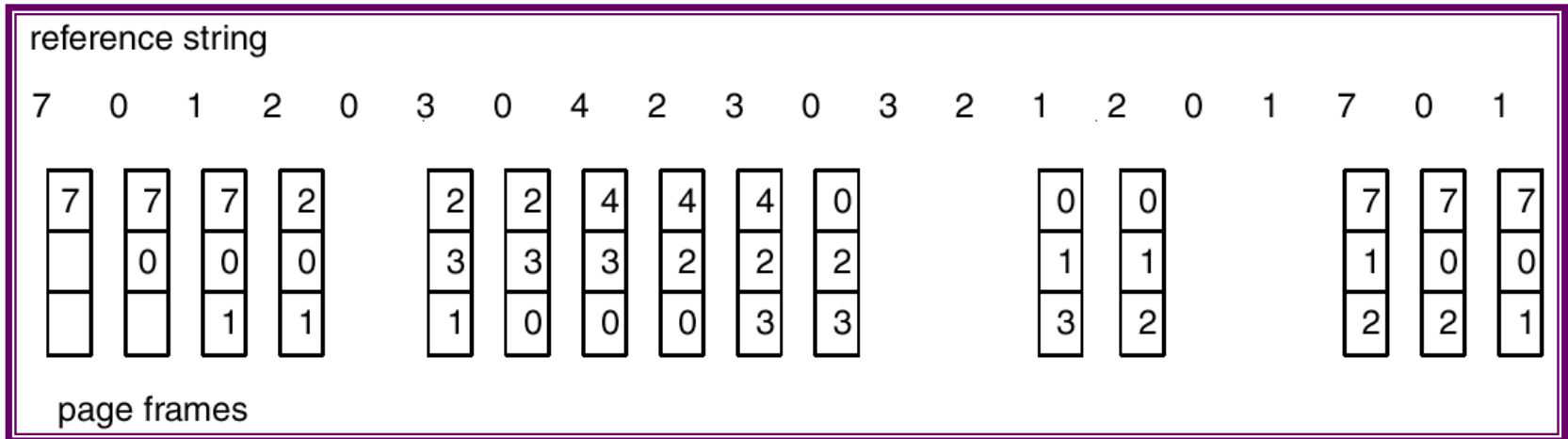
1	1	4	5	
2	2	1	3	9 page faults
3	3	2	4	

- 4 frames

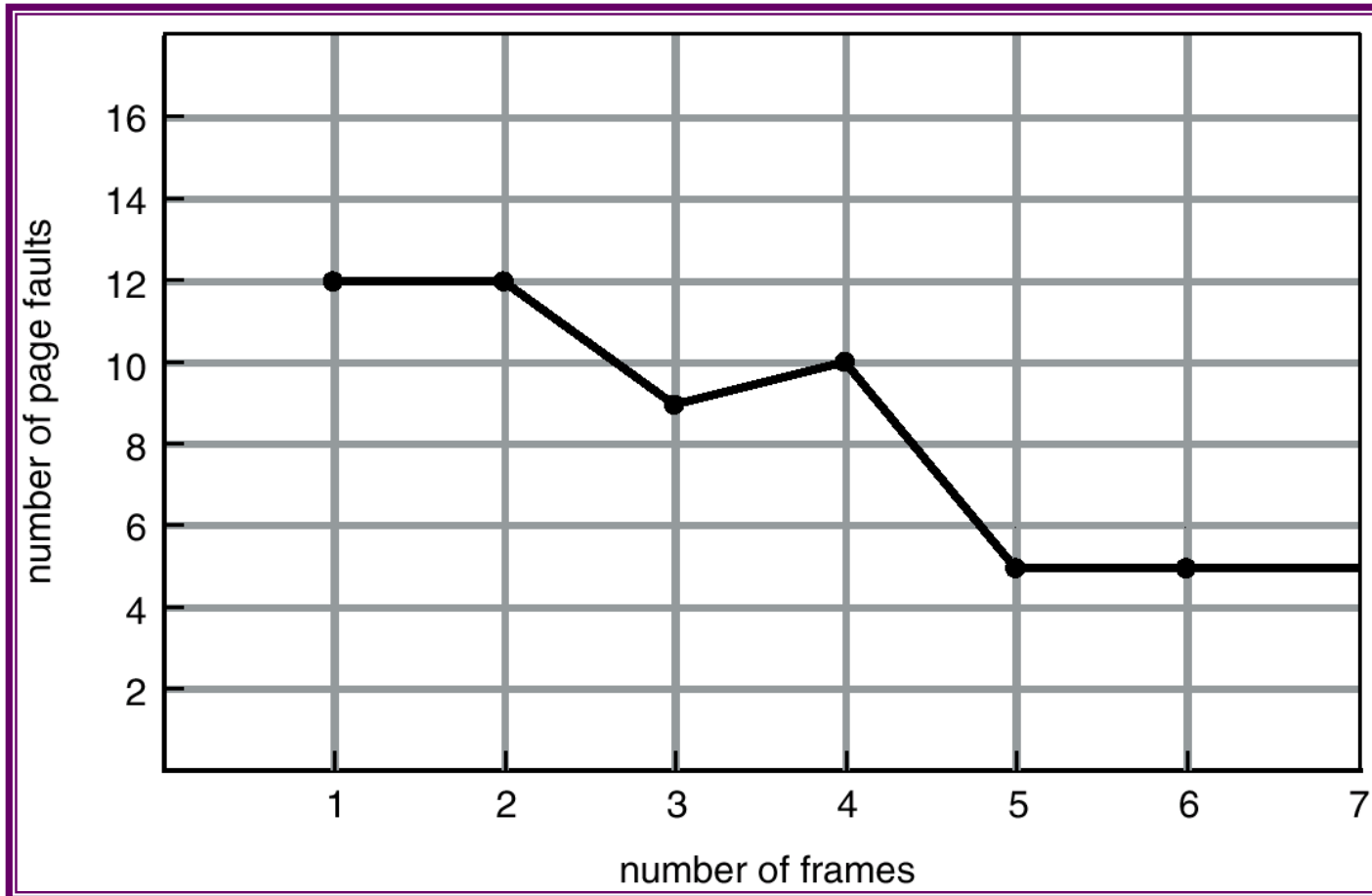
1	1	5	4	
2	2	1	5	10 page faults
3	3	2		
4	4	3		

- FIFO Replacement – Belady’s Anomaly  
☞ more frames  $\Rightarrow$  less page faults

# FIFO Page Replacement

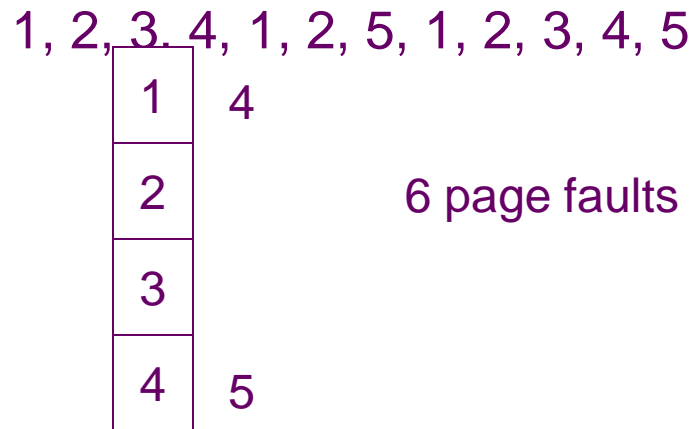


# FIFO Illustrating Belady's Anomaly



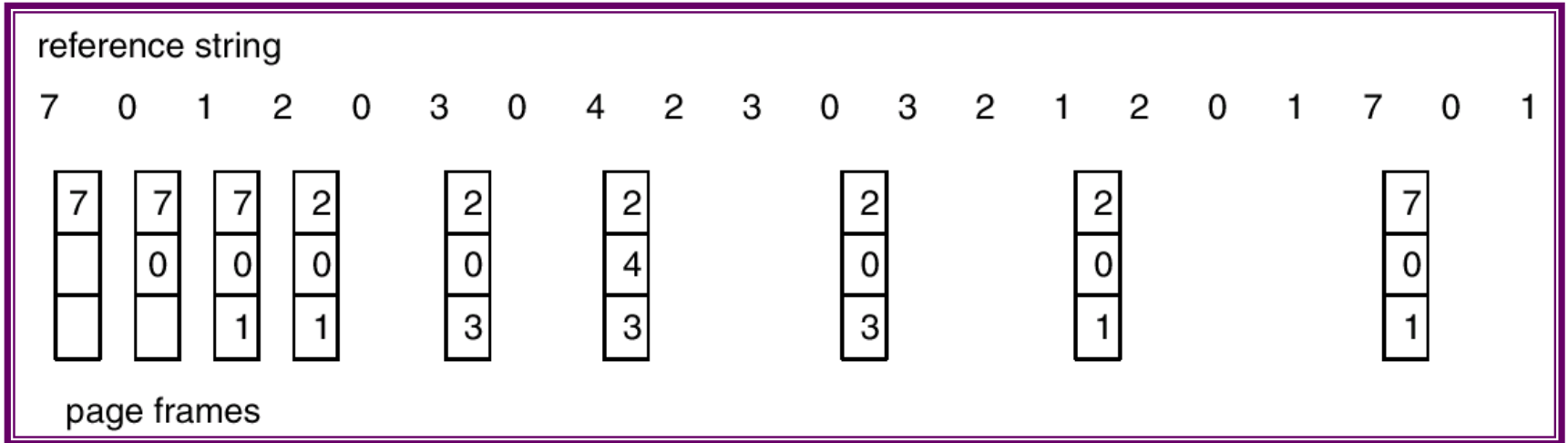
# Optimal Algorithm

- Replace page that will not be used for longest period of time.
- 4 frames example



- How do you know this?
- Used for measuring how well your algorithm performs.

# Optimal Page Replacement



# Least Recently Used (LRU) Algorithm

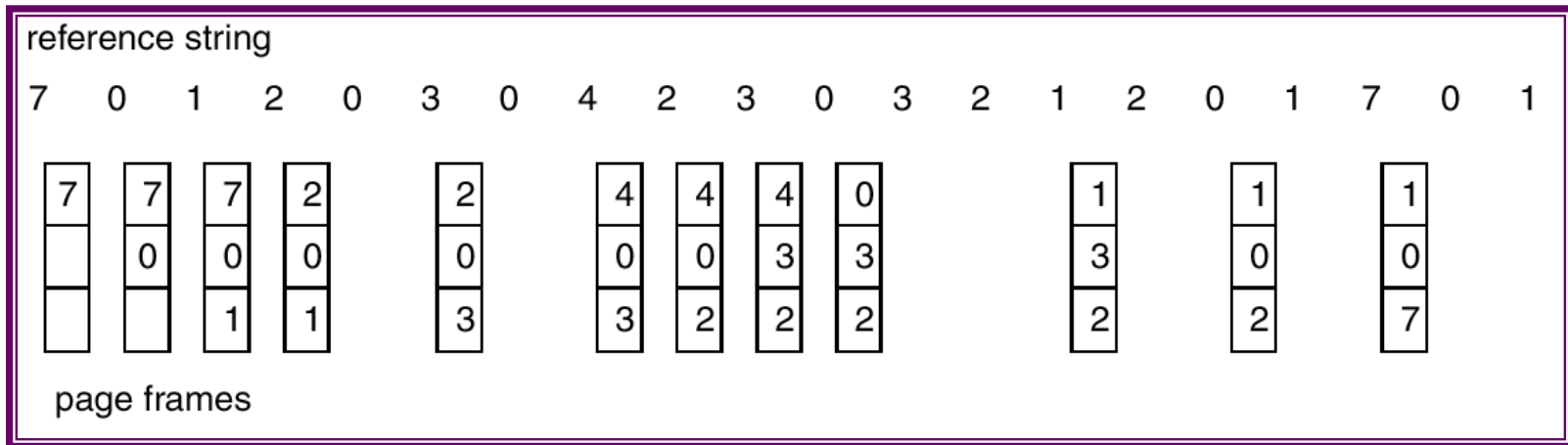
- Reference string: 1, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5
- |   |   |   |
|---|---|---|
| 1 | 3 | 4 |
| 2 |   |   |
| 3 | 5 | 4 |
| 4 | 3 |   |

- Counter implementation

- ☞ Every page entry has a counter; every time page is referenced through this entry, copy the clock into the counter.
- ☞ When a page needs to be changed, look at the counters to determine which are to change.



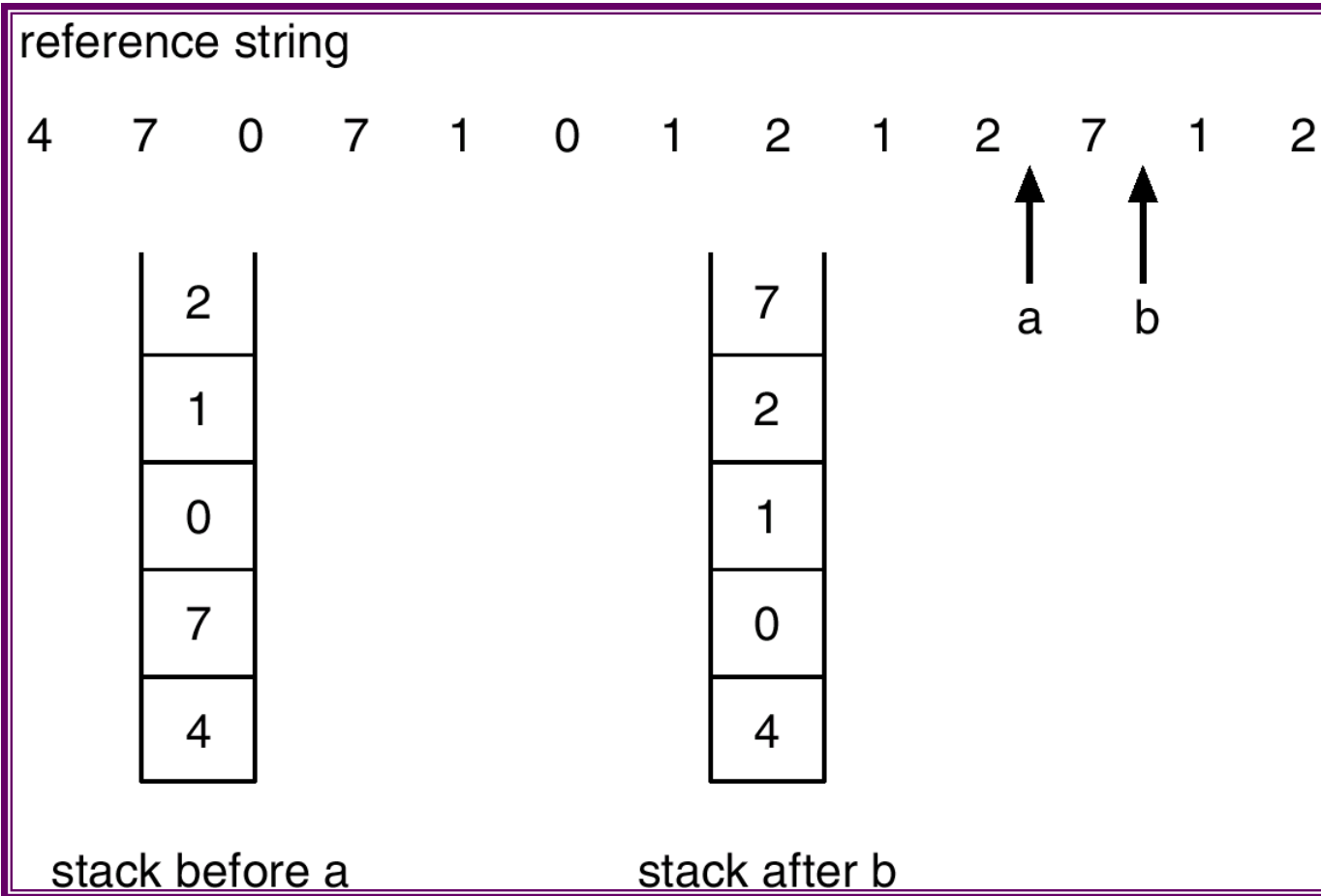
# LRU Page Replacement



# LRU Algorithm (Cont.)

- Stack implementation – keep a stack of page numbers in a double link form:
  - ☞ Page referenced:
    - 📄 move it to the top
    - 📄 requires 6 pointers to be changed
  - ☞ No search for replacement

# Use Of A Stack To Record The Most Recent Page References



# LRU Approximation Algorithms

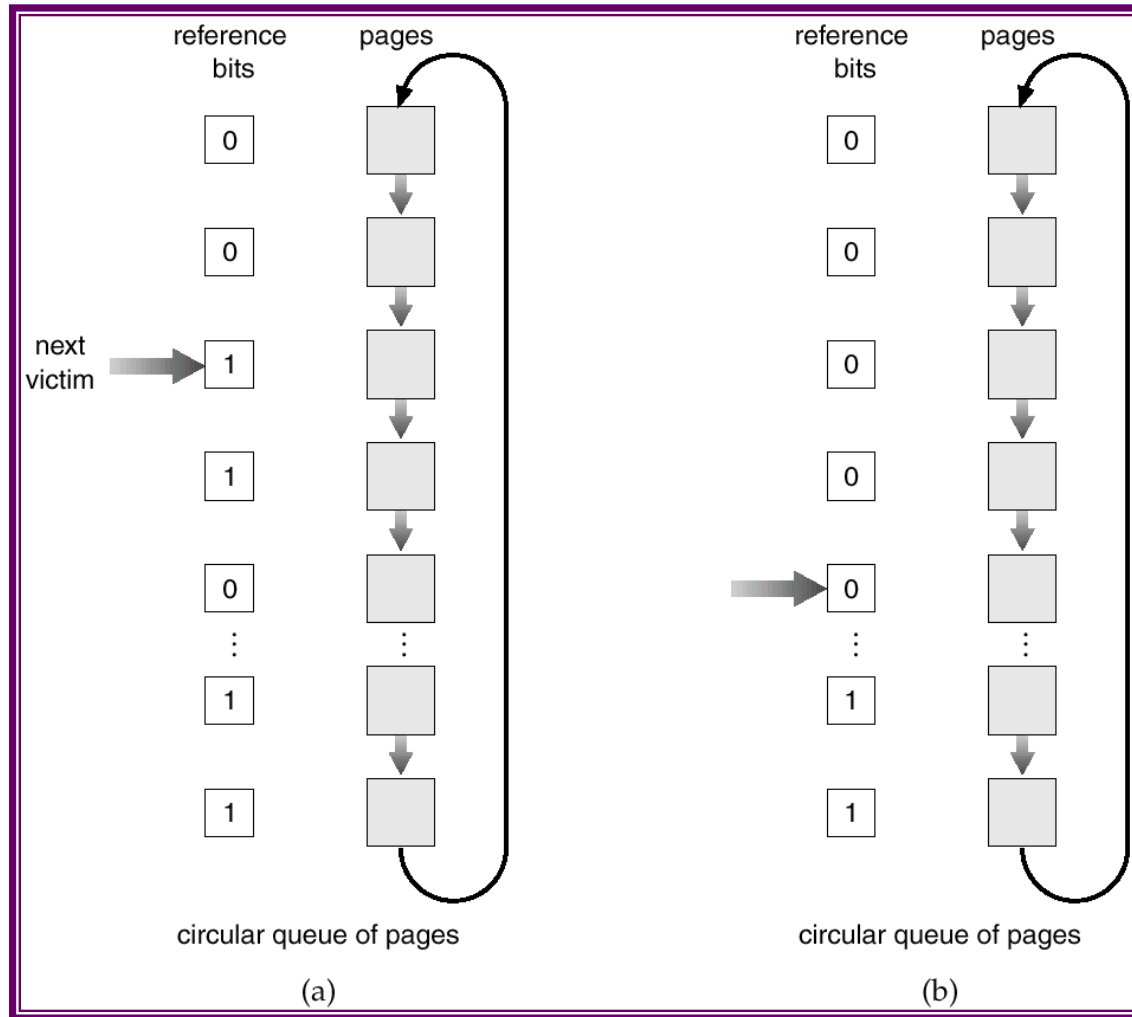
## ■ Reference bit

- ☞ With each page associate a bit, initially = 0
- ☞ When page is referenced bit set to 1.
- ☞ Replace the one which is 0 (if one exists). We do not know the order, however.

## ■ Second chance

- ☞ Need reference bit.
- ☞ Clock replacement.
- ☞ If page to be replaced (in clock order) has reference bit = 1. then:
  - 📄 set reference bit 0.
  - 📄 leave page in memory.
  - 📄 replace next page (in clock order), subject to same rules.

# Second-Chance (clock) Page-Replacement Algorithm



# Counting Algorithms

- Keep a counter of the number of references that have been made to each page.
- LFU Algorithm: replaces page with smallest count.
- MFU Algorithm: based on the argument that the page with the smallest count was probably just brought in and has yet to be used.

# Allocation of Frames

- Each process needs **minimum** number of pages.
- 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**.
- Two major allocation schemes.
  - ☞ fixed allocation
  - ☞ priority allocation

# Fixed Allocation

- Equal allocation – e.g., if 100 frames and 5 processes, give each 20 pages.
- Proportional allocation – Allocate according to the size of process.

–  $s_i$  = size of process  $p_i$

–  $S = \sum s_i$

–  $m$  = total number of frames

–  $a_i$  = allocation for  $p_i = \frac{s_i}{S} \times m$

$$m = 64$$

$$s_1 = 10$$

$$s_2 = 127$$

$$a_1 = \frac{10}{137} \times 64 \approx 5$$

$$a_2 = \frac{127}{137} \times 64 \approx 59$$



# Priority Allocation

- Use a proportional allocation scheme using priorities rather than size.
- 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.

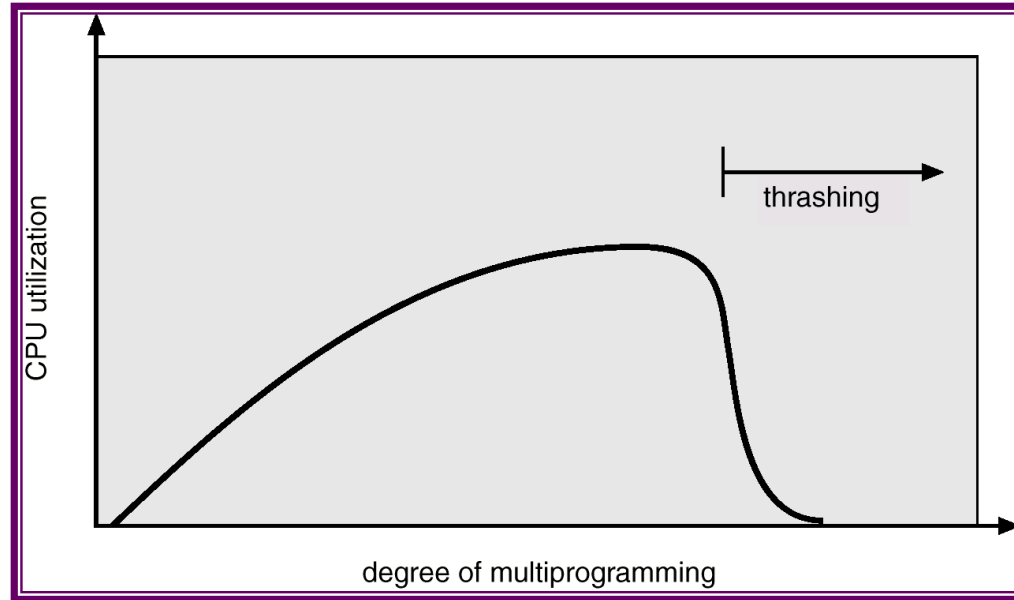
# Global vs. Local Allocation

- **Global** replacement – process selects a replacement frame from the set of all frames; one process can take a frame from another.
- **Local** replacement – each process selects from only its own set of allocated frames.

# Thrashing

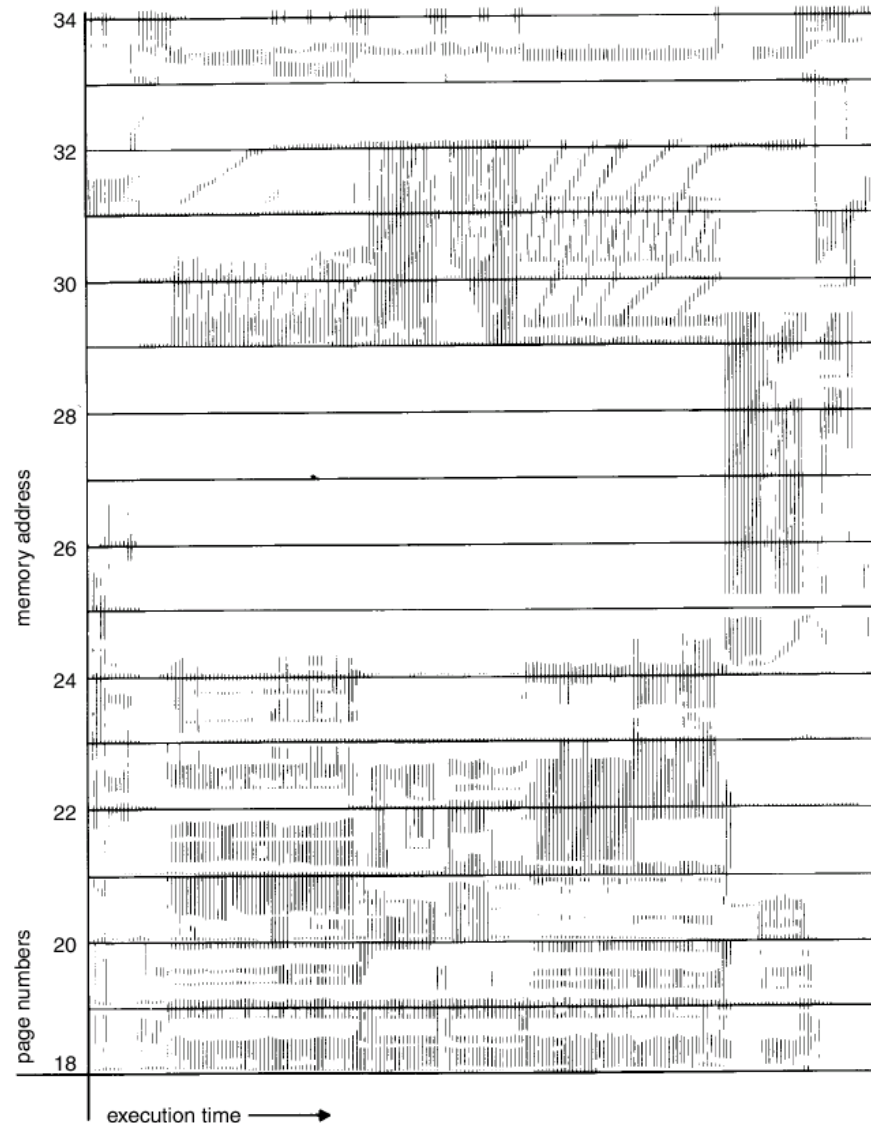
- If a process does not have “enough” pages, the page-fault rate is very high. This leads to:
  - ☞ low CPU utilization.
  - ☞ operating system thinks that it needs to increase the degree of multiprogramming.
  - ☞ another process added to the system.
- **Thrashing**  $\equiv$  a process is busy swapping pages in and out.

# Thrashing



- Why does paging work?  
Locality model
  - ☞ Process migrates from one locality to another.
  - ☞ Localities may overlap.
- Why does thrashing occur?  
 $\Sigma$  size of locality > total memory size

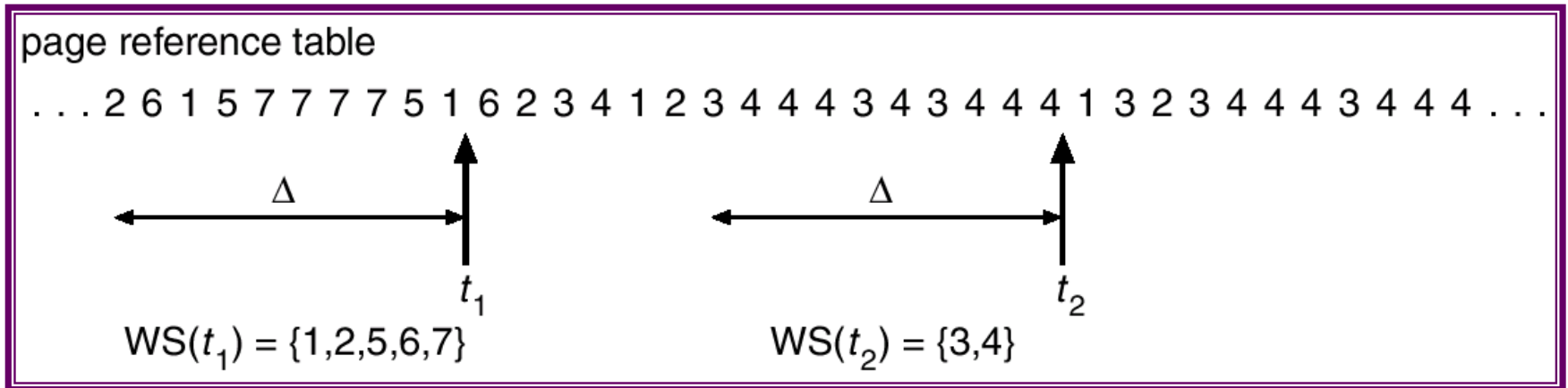
# Locality In A Memory-Reference Pattern



# Working-Set Model

- $\Delta \equiv$  working-set window  $\equiv$  a fixed number of page references  
Example: 10,000 instruction
- $WSS_i$  (working set of Process  $P_i$ ) = total number of pages referenced in the most recent  $\Delta$  (varies in time)
  - ☞ if  $\Delta$  too small will not encompass entire locality.
  - ☞ if  $\Delta$  too large will encompass several localities.
  - ☞ if  $\Delta = \infty \Rightarrow$  will encompass entire program.
- $D = \sum WSS_i \equiv$  total demand frames
- if  $D > m \Rightarrow$  Thrashing
- Policy if  $D > m$ , then suspend one of the processes.

# Working-set model

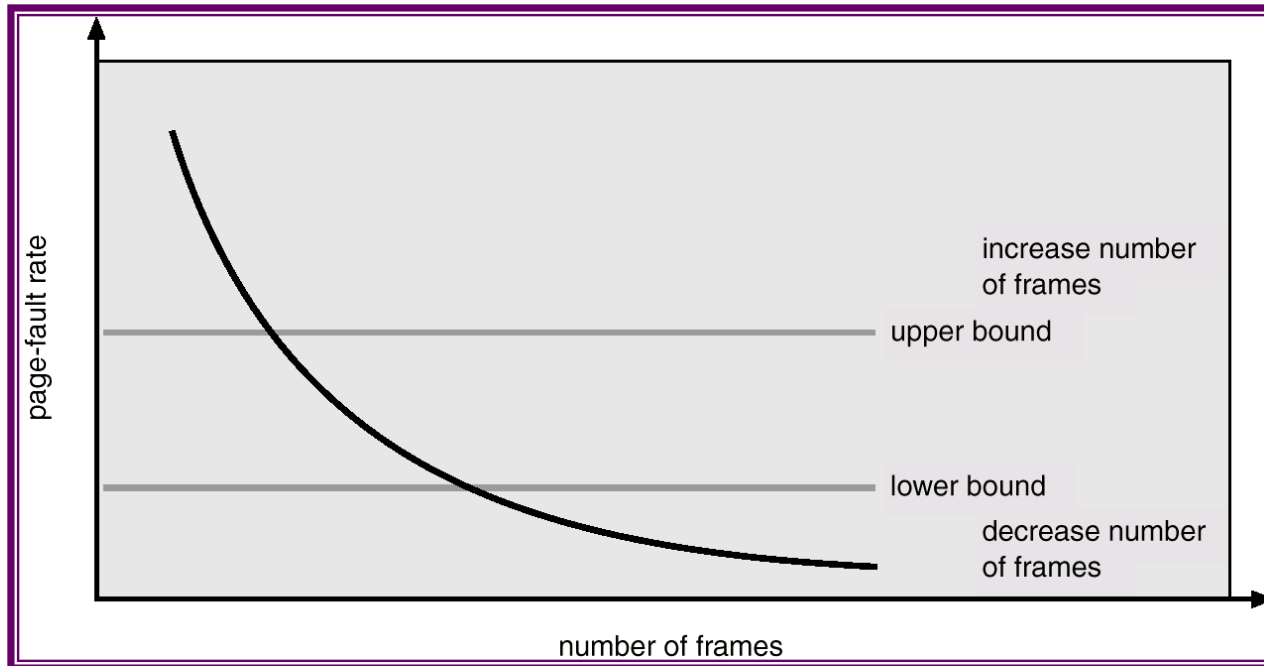


# Keeping Track of the Working Set

- Approximate with interval timer + a reference bit
- Example:  $\Delta = 10,000$ 
  - ☞ Timer interrupts after every 5000 time units.
  - ☞ Keep in memory 2 bits for each page.
  - ☞ Whenever a timer interrupts copy and sets the values of all reference bits to 0.
  - ☞ If one of the bits in memory = 1  $\Rightarrow$  page in working set.
- Why is this not completely accurate?
- Improvement = 10 bits and interrupt every 1000 time units.



# Page-Fault Frequency Scheme



- Establish “acceptable” page-fault rate.
  - ☞ If actual rate too low, process loses frame.
  - ☞ If actual rate too high, process gains frame.

# Other Considerations

- Prepaging
- Page size selection
  - ☞ fragmentation
  - ☞ table size
  - ☞ I/O overhead
  - ☞ locality

# Other Considerations (Cont.)

- **TLB Reach** - The amount of memory accessible from the TLB.
- $\text{TLB Reach} = (\text{TLB Size}) \times (\text{Page Size})$
- Ideally, the working set of each process is stored in the TLB. Otherwise there is a high degree of page faults.

# Increasing the Size of the TLB

- **Increase the Page Size.** This may lead to an increase in fragmentation as not all applications require a large page size.
- **Provide Multiple Page Sizes.** This allows applications that require larger page sizes the opportunity to use them without an increase in fragmentation.

# Other Considerations (Cont.)

## ■ Program structure

☞ `int A[][] = new int[1024][1024];`

☞ Each row is stored in one page

☞ Program 1

```
for (j = 0; j < A.length; j++)  
    for (i = 0; i < A.length; i++)  
        A[i,j] = 0;
```

1024 x 1024 page faults

☞ Program 2

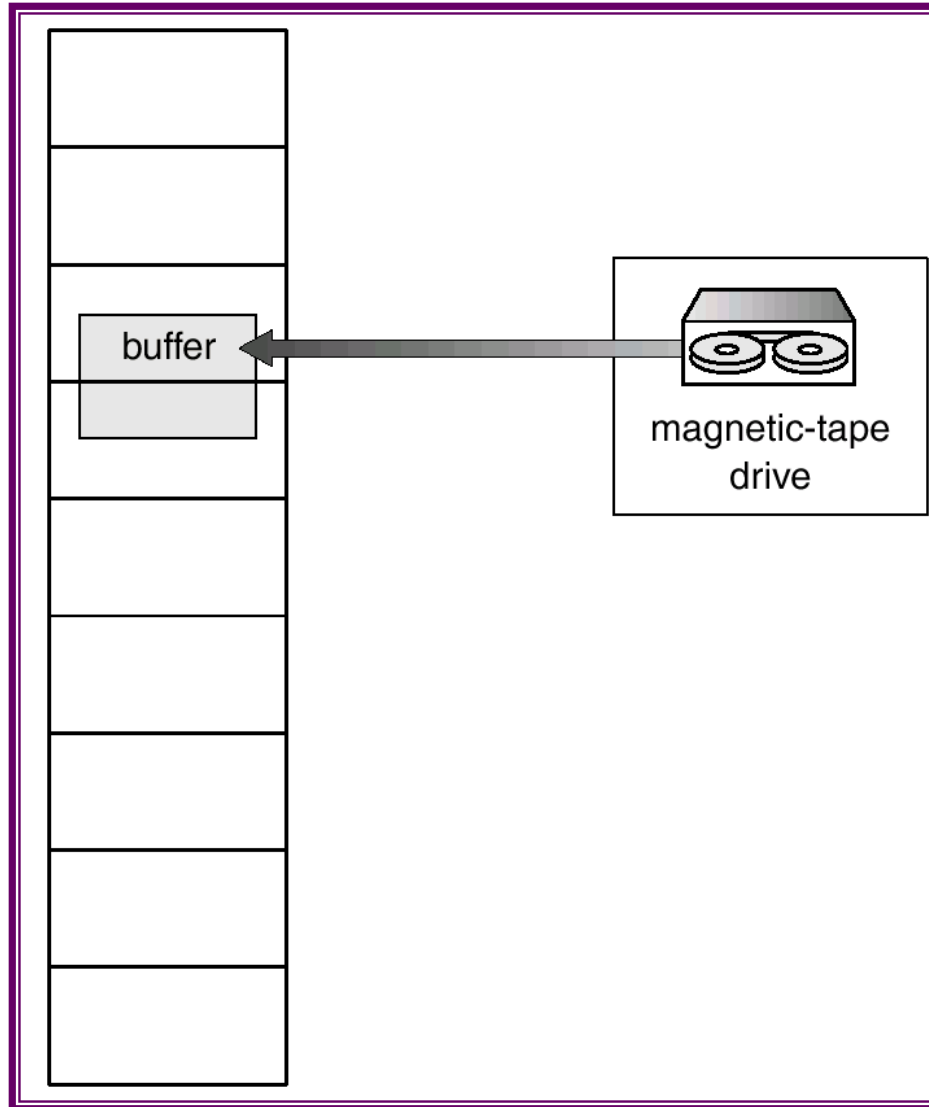
```
for (i = 0; i < A.length; i++)  
    for (j = 0; j < A.length; j++)  
        A[i,j] = 0;
```

1024 page faults

# Other Considerations (Cont.)

- **I/O Interlock** – Pages must sometimes be locked into memory.
- Consider I/O. Pages that are used for copying a file from a device must be locked from being selected for eviction by a page replacement algorithm.

# Reason Why Frames Used For I/O Must Be In Memory



# Operating System Examples

- Windows NT
- Solaris 2



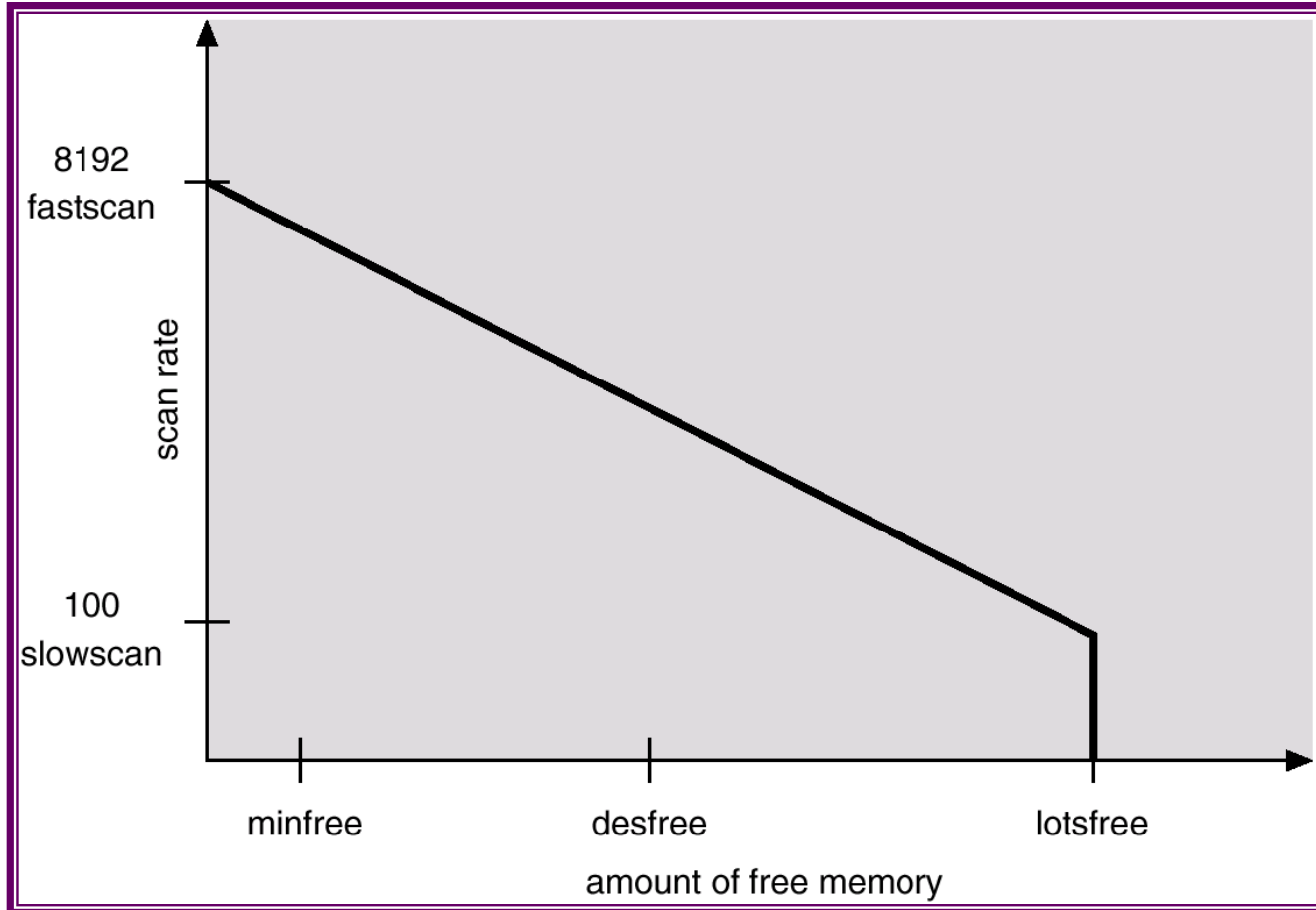
# Windows NT

- Uses demand paging with **clustering**. Clustering brings in pages surrounding the faulting page.
- Processes are assigned **working set minimum** and **working set maximum**.
- Working set minimum is the minimum number of pages the process is guaranteed to have in memory.
- A process may be assigned as many pages up to its working set maximum.
- When the amount of free memory in the system falls below a threshold, **automatic working set trimming** is performed to restore the amount of free memory.
- Working set trimming removes pages from processes that have pages in excess of their working set minimum.

# Solaris 2

- Maintains a list of free pages to assign faulting processes.
- **Lotsfree** – threshold parameter to begin paging.
- Paging is performed by *pageout* process.
- Pageout scans pages using modified clock algorithm.
- **Scanrate** is the rate at which pages are scanned. This ranged from **slowscan** to **fastscan**.
- Pageout is called more frequently depending upon the amount of free memory available.

# Solar Page Scanner

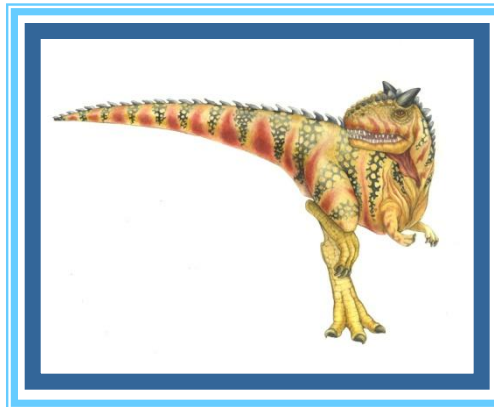




# UNIT IV

## File-System Interface

---



# Chapter 11: File-System Interface

- File Concept
- Access Methods
- Disk and Directory Structure
- File-System Mounting
- File Sharing
- Protection

# Objectives

- To explain the function of file systems
- To describe the interfaces to file systems
- To discuss file-system design tradeoffs, including access methods, file sharing, file locking, and directory structures
- To explore file-system protection

# File Concept

- Contiguous logical address space

- Types:

  - ☞ Data

    - 📄 numeric

    - 📄 character

    - 📄 binary

  - ☞ Program

- Contents defined by file's creator

  - ☞ Many types

    - 📄 Consider **text file, source file, executable file**



# File Attributes

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

# File info Window on Mac OS X



# File Operations

- File is an **abstract data type**
- **Create**
- **Write** – at **write pointer** location
- **Read** – at **read pointer** location
- **Reposition within file - seek**
- **Delete**
- **Truncate**
- **$Open(F_i)$**  – search the directory structure on disk for entry  $F_i$ , and move the content of entry to memory
- **$Close(F_i)$**  – move the content of entry  $F_i$  in memory to directory structure on disk

# Open Files

- Several pieces of data are needed to manage open files:
  - ☞ **Open-file table:** tracks open files
  - ☞ File pointer: pointer to last read/write location, per process that has the file open
  - ☞ **File-open count:** counter of number of times a file is open – to allow removal of data from open-file table when last processes closes it
  - ☞ Disk location of the file: cache of data access information
  - ☞ Access rights: per-process access mode information

# Open File Locking

- Provided by some operating systems and file systems
  - ☞ Similar to reader-writer locks
  - ☞ **Shared lock** similar to reader lock – several processes can acquire concurrently
  - ☞ **Exclusive lock** similar to writer lock
- Mediates access to a file
- Mandatory or advisory:
  - ☞ **Mandatory** – access is denied depending on locks held and requested
  - ☞ **Advisory** – processes can find status of locks and decide what to do

# File Locking Example – Java API

```
import java.io.*;
import java.nio.channels.*;
public class LockingExample {
    public static final boolean EXCLUSIVE = false;
    public static final boolean SHARED = true;
    public static void main(String arsg[]) throws IOException {
        FileLock sharedLock = null;
        FileLock exclusiveLock = null;
        try {
            RandomAccessFile raf = new RandomAccessFile("file.txt", "rw");
            // get the channel for the file
            FileChannel ch = raf.getChannel();
            // this locks the first half of the file - exclusive
            exclusiveLock = ch.lock(0, raf.length()/2, EXCLUSIVE);
            /** Now modify the data . . . */
            // release the lock
            exclusiveLock.release();
        }
    }
}
```

# File Locking Example – Java API (Cont.)

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

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

# File Types – Name, Extension

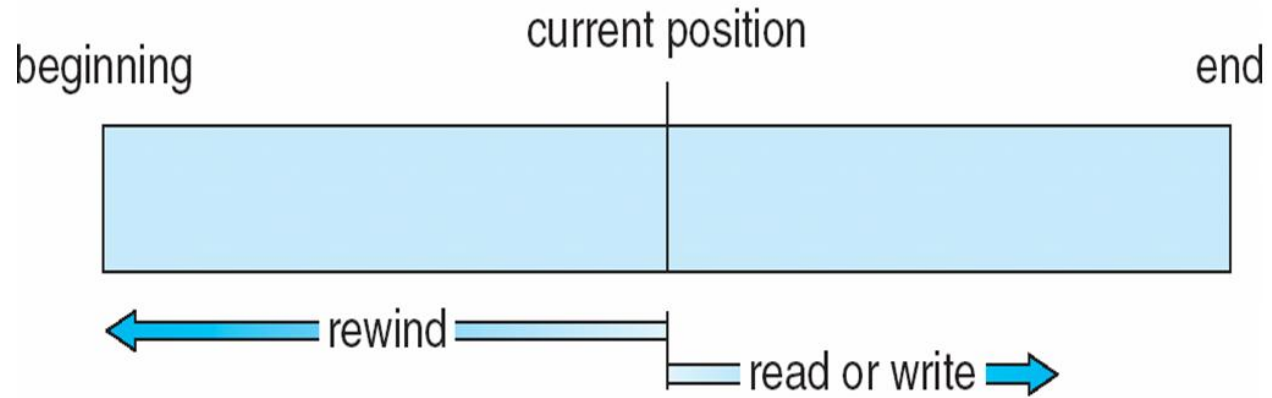
file type	usual extension	function
executable	exe, com, bin or none	ready-to-run machine- language program
object	obj, o	compiled, machine language, not linked
source code	c, cc, java, pas, asm, a	source code in various languages
batch	bat, sh	commands to the command interpreter
text	txt, doc	textual 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 com- pressed, for archiving or storage
multimedia	mpeg, mov, rm, mp3, avi	binary file containing audio or A/V information



# File Structure

- None - sequence of words, bytes
- Simple record structure
  - ☞ Lines
  - ☞ Fixed length
  - ☞ Variable length
- Complex Structures
  - ☞ Formatted document
  - ☞ Relocatable load file
- Can simulate last two with first method by inserting appropriate control characters
- Who decides:
  - ☞ Operating system
  - ☞ Program

# Sequential-access File



# Access Methods

- **Sequential Access**

```
read next
write next
reset
no read after last write
    (rewrite)
```

- **Direct Access** – file is fixed length **logical records**

```
read n
write n
position to n
    read next
    write next
rewrite n
```

*n* = relative block number

- Relative block numbers allow OS to decide where file should be placed

- ☞ See allocation problem in Ch 12

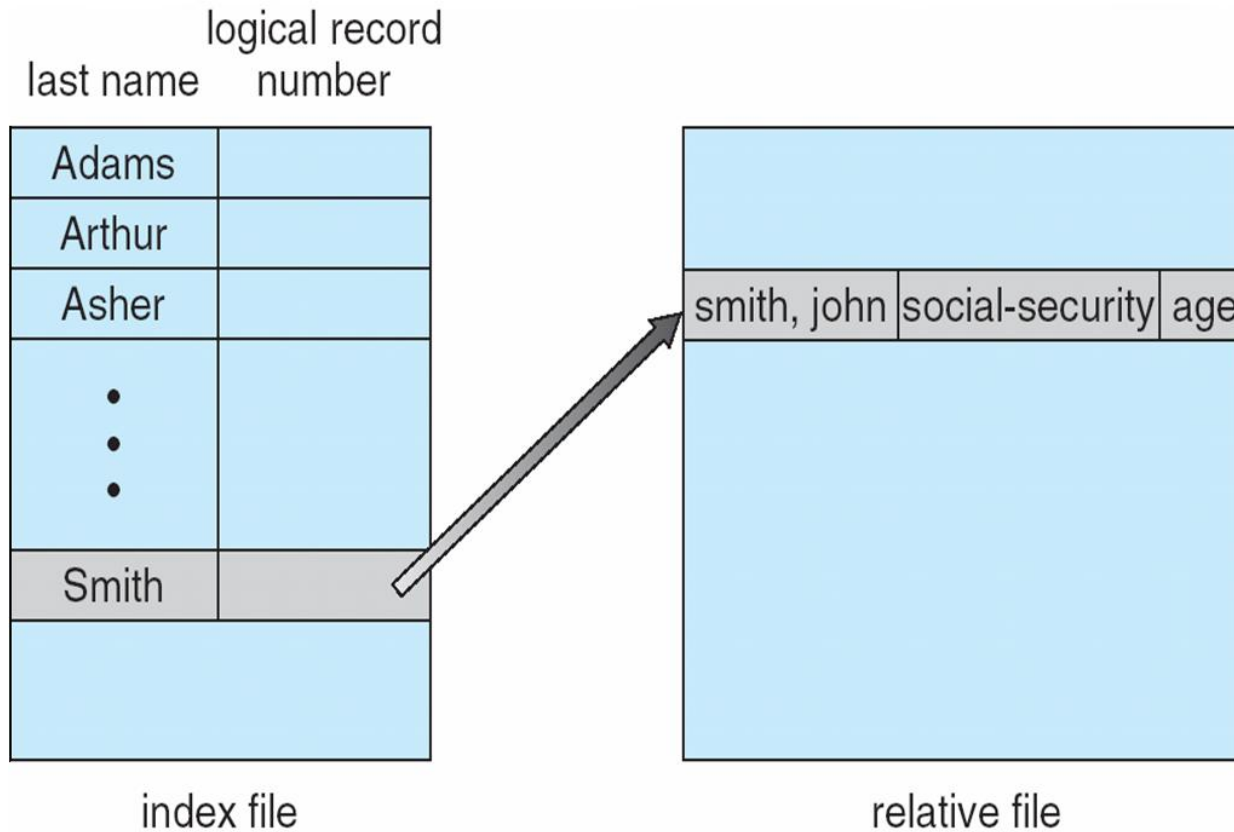
## Simulation of Sequential Access on Direct-access File

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>

# Other Access Methods

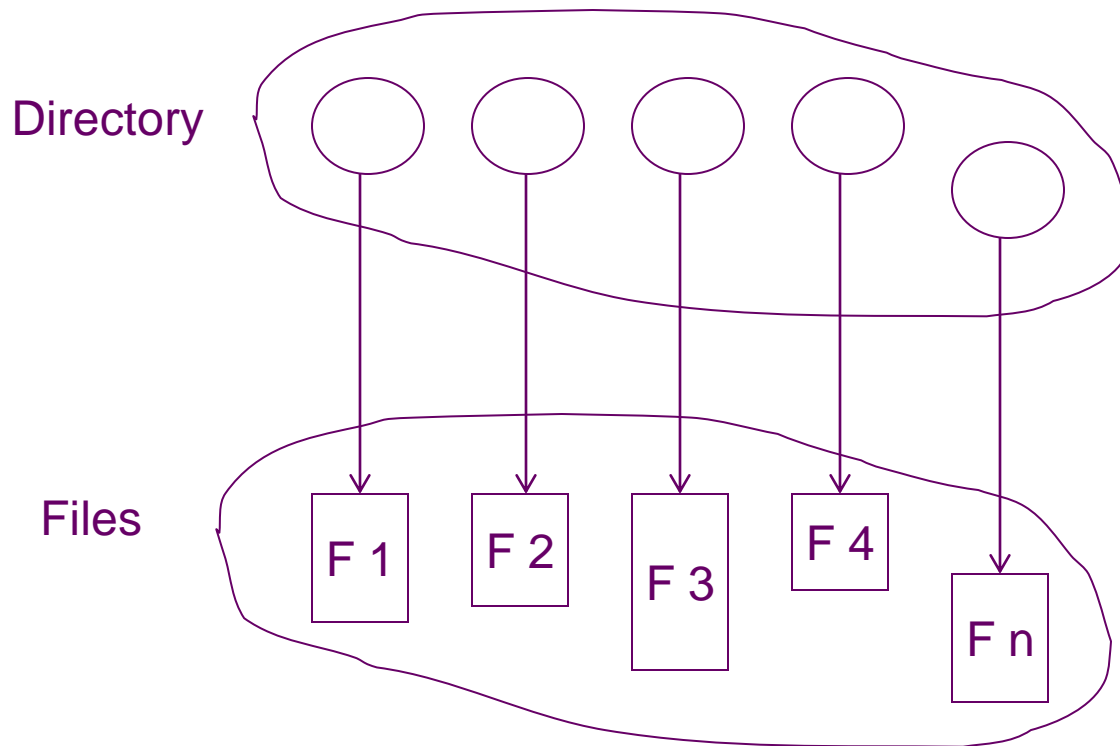
- Can be built on top of base methods
- General involve creation of an **index** for the file
- Keep index in memory for fast determination of location of data to be operated on (consider UPC code plus record of data about that item)
- If too large, index (in memory) of the index (on disk)
- IBM indexed sequential-access method (ISAM)
  - ☞ Small master index, points to disk blocks of secondary index
  - ☞ File kept sorted on a defined key
  - ☞ All done by the OS
- VMS operating system provides index and relative files as another example (see next slide)

# Example of Index and Relative Files



# Directory Structure

- A collection of nodes containing information about all files



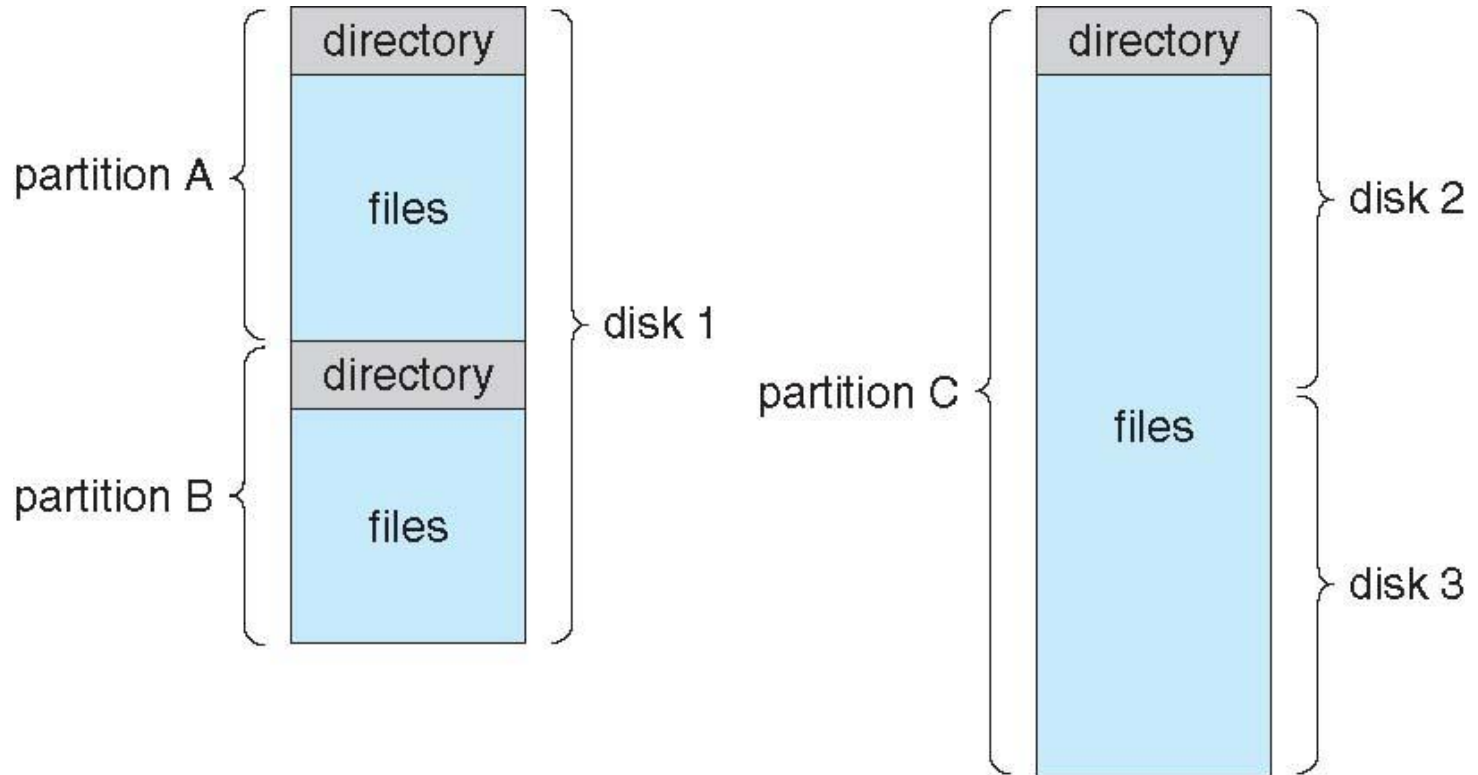
Both the directory structure and the files reside on disk

# Disk Structure

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



# A Typical File-system Organization



# Types of File Systems

- We mostly talk of general-purpose file systems
- But systems frequently have many file systems, some general- and some special- purpose
- Consider Solaris has
  - ☞ tmpfs – memory-based volatile FS for fast, temporary I/O
  - ☞ objfs – interface into kernel memory to get kernel symbols for debugging
  - ☞ ctfs – contract file system for managing daemons
  - ☞ lofs – loopback file system allows one FS to be accessed in place of another
  - ☞ procfs – kernel interface to process structures
  - ☞ ufs, zfs – general purpose file systems

# Operations Performed on Directory

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

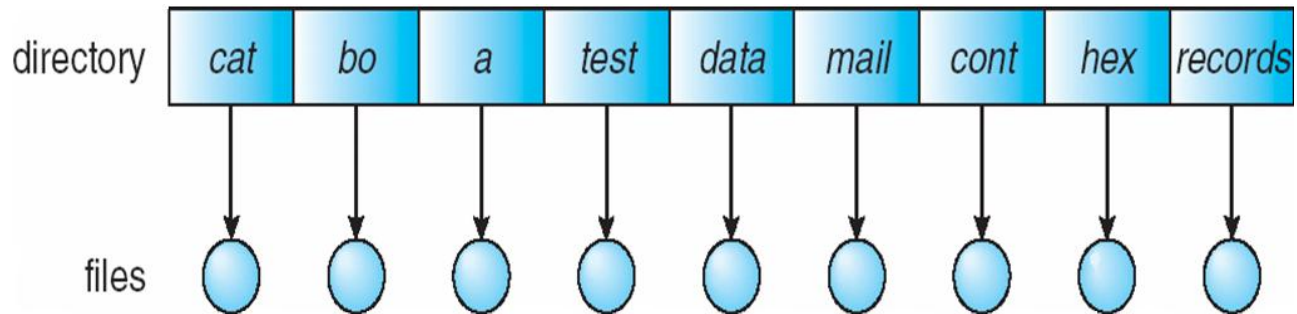
# Directory Organization

The directory is organized logically to obtain

- Efficiency – locating a file quickly
- Naming – convenient to users
  - ☞ Two users can have same name for different files
  - ☞ The same file can have several different names
- Grouping – logical grouping of files by properties, (e.g., all Java programs, all games, ...)

# Single-Level Directory

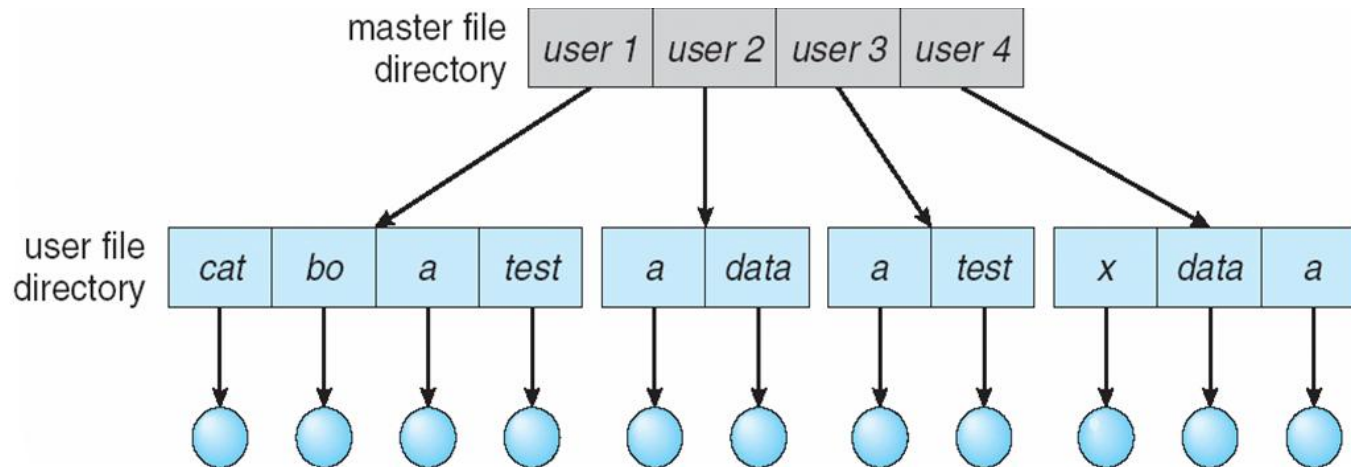
- A single directory for all users



- Naming problem
- Grouping problem

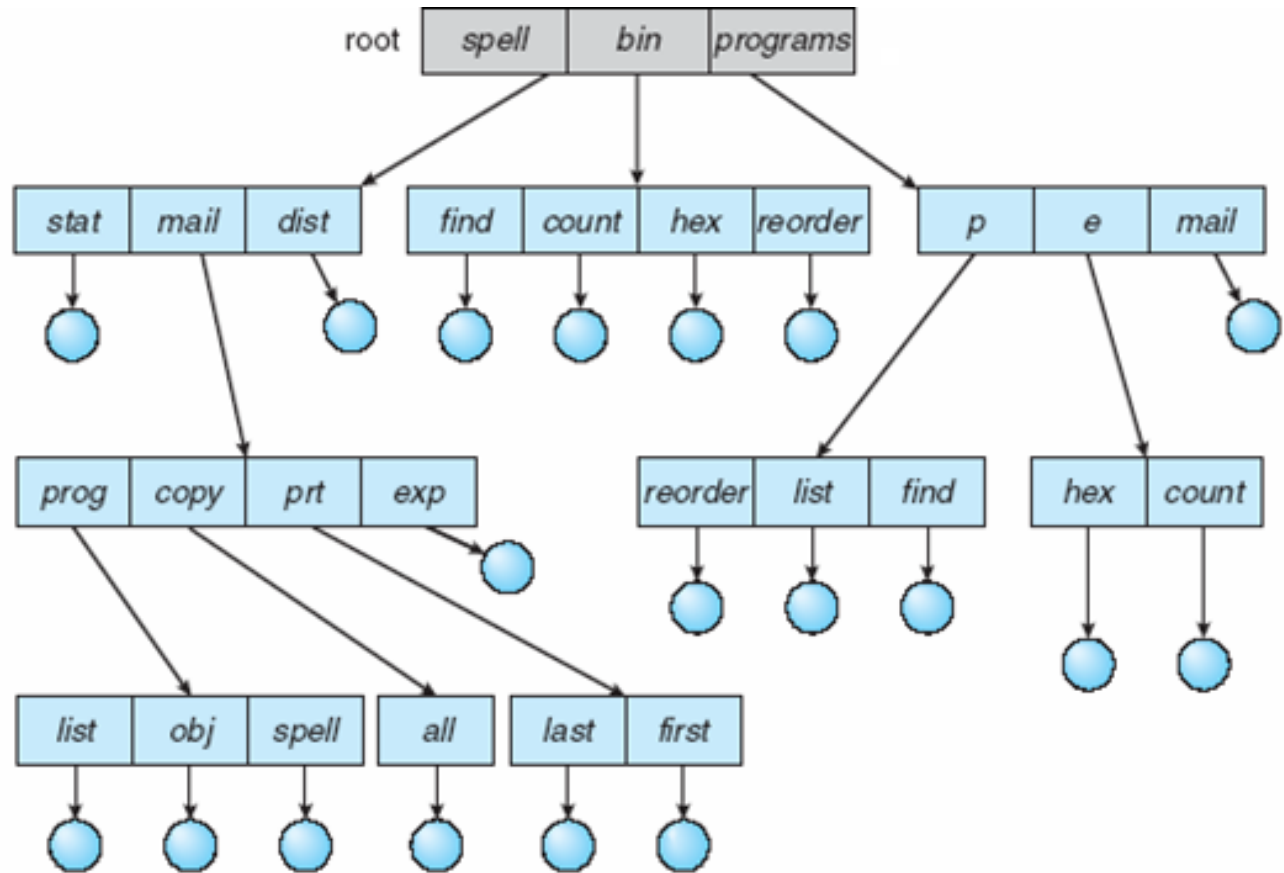
# Two-Level Directory

- Separate directory for each user



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

# Tree-Structured Directories



# Tree-Structured Directories (Cont.)

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



# Tree-Structured Directories (Cont)

- **Absolute** or **relative** path name
- Creating a new file is done in current directory
- Delete a file

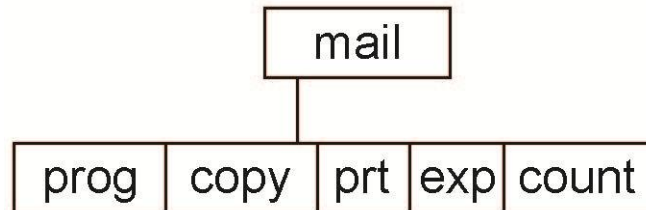
```
rm <file-name>
```

- Creating a new subdirectory is done in current directory

```
mkdir <dir-name>
```

Example: if in current directory `/mail`

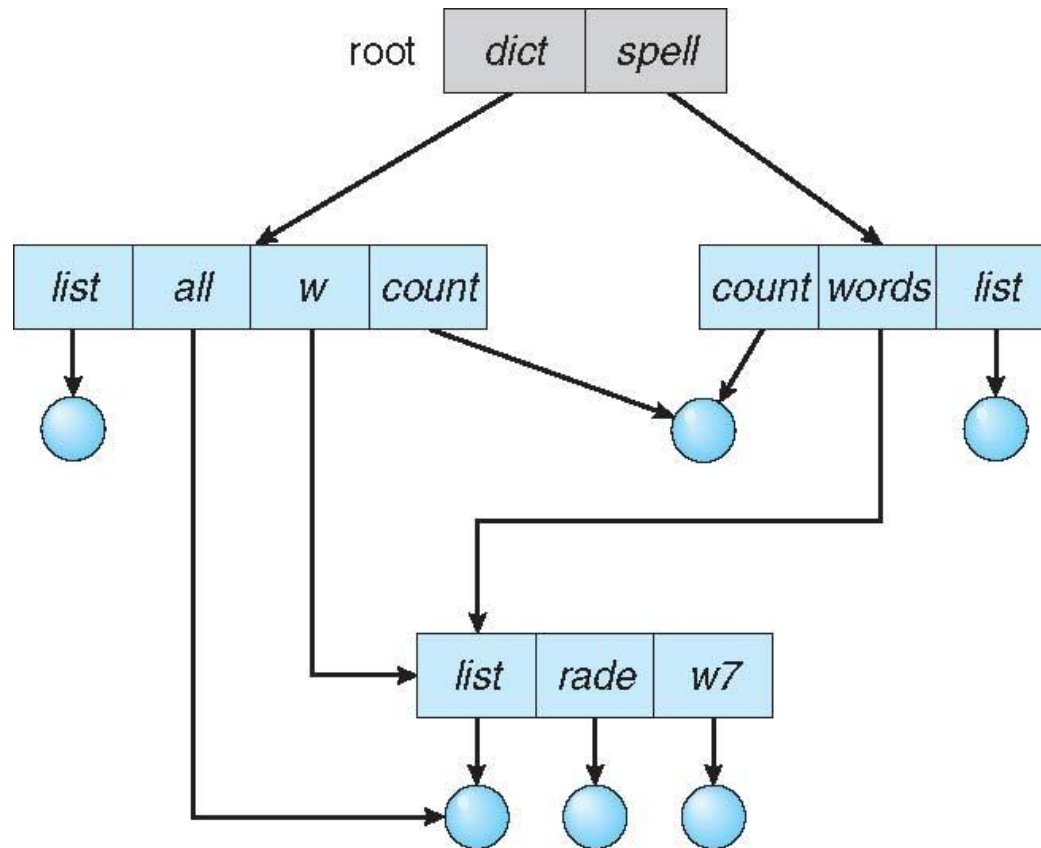
```
mkdir count
```



Deleting “mail”  $\Rightarrow$  deleting the entire subtree rooted by “mail”

# Acyclic-Graph Directories

- Have shared subdirectories and files



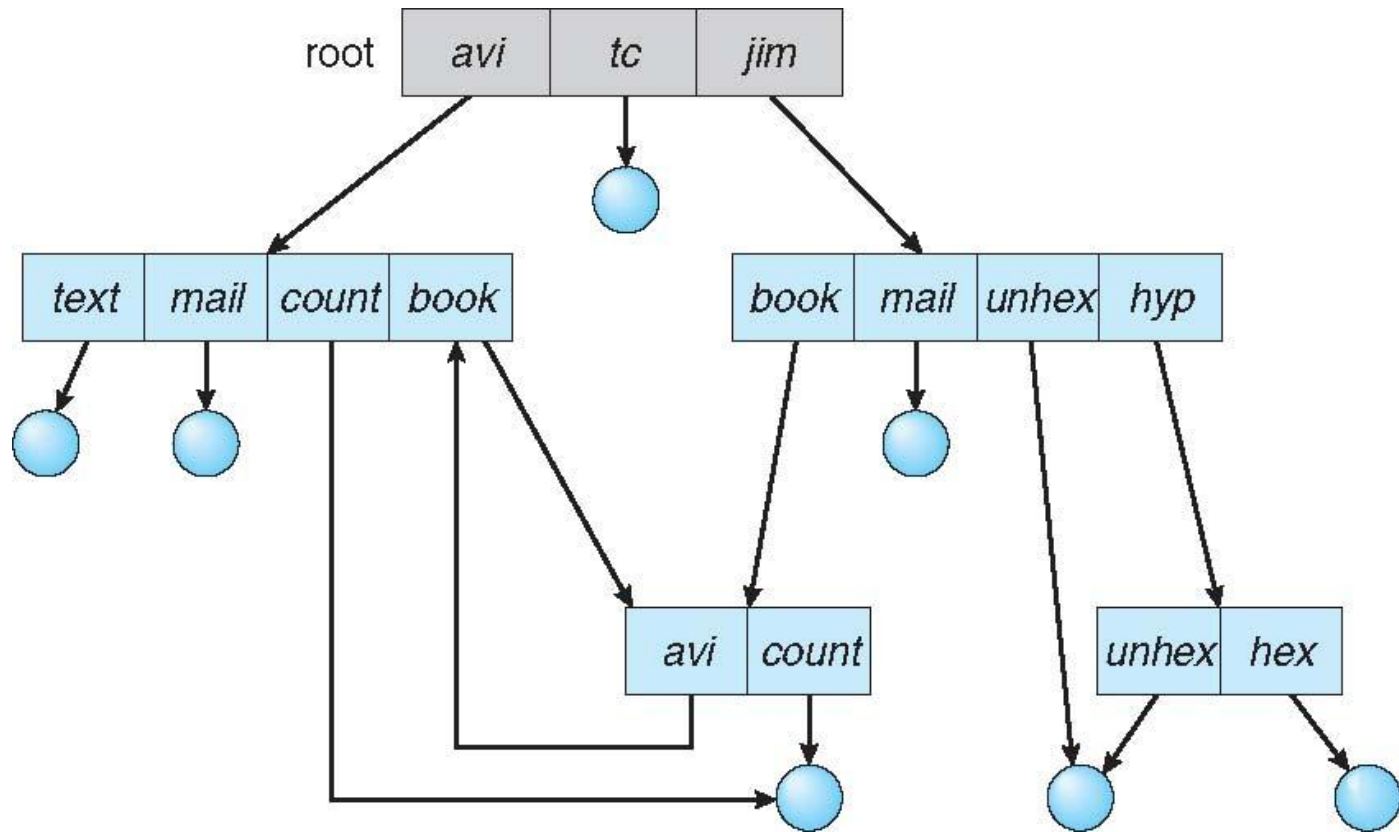
# Acyclic-Graph Directories (Cont.)

- Two different names (aliasing)
- If *dict* deletes *list*  $\Rightarrow$  dangling pointer

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 entry type
  - ☞ **Link** – another name (pointer) to an existing file
  - ☞ **Resolve the link** – follow pointer to locate the file

# General Graph Directory

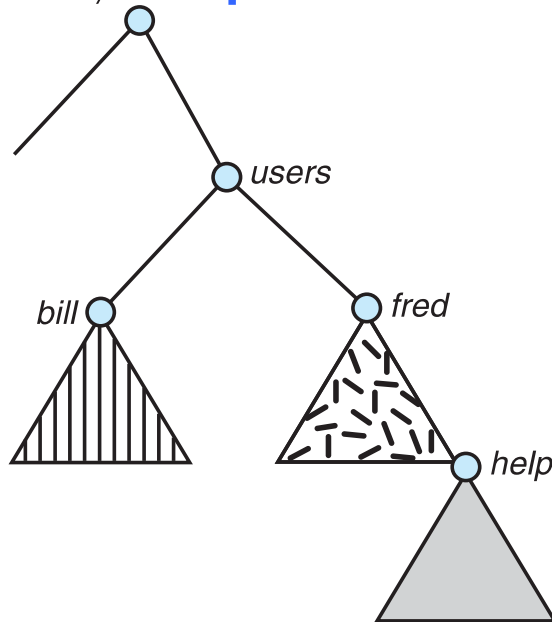


# General Graph Directory (Cont.)

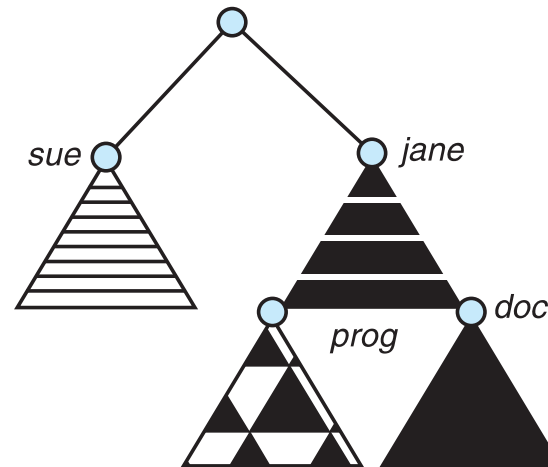
- How do we guarantee no cycles?
  - ☞ Allow only links to file not subdirectories
  - ☞ **Garbage collection**
  - ☞ Every time a new link is added use a cycle detection algorithm to determine whether it is OK

# File System Mounting

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

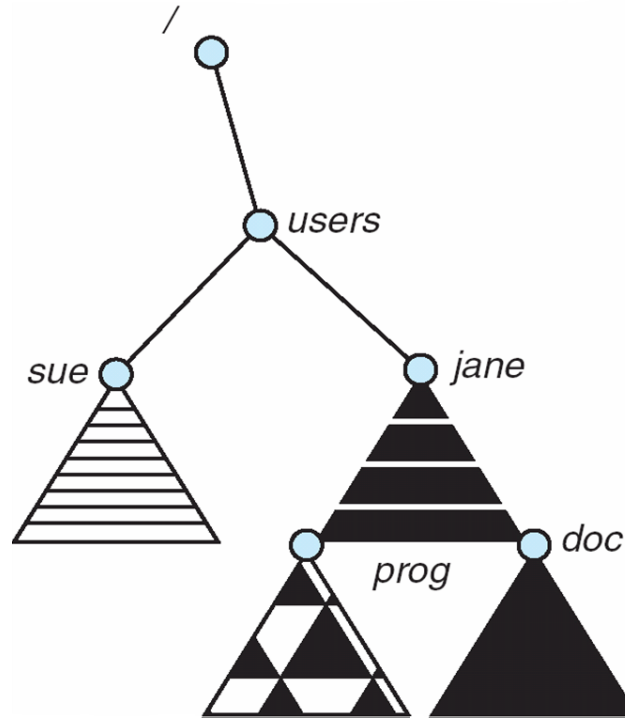


(a)



(b)

# Mount Point



# File Sharing

- Sharing of files on multi-user systems is desirable
- Sharing may be done through a **protection** scheme
- On distributed systems, files may be shared across a network
- Network File System (NFS) is a common distributed file-sharing method
- If multi-user system
  - ☞ **User IDs** identify users, allowing permissions and protections to be per-user
  - Group IDs** allow users to be in groups, permitting group access rights
  - ☞ Owner of a file / directory
  - ☞ Group of a file / directory



# File Sharing – Remote File Systems

- 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**
- **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
- Distributed Information Systems (**distributed naming services**) such as LDAP, DNS, NIS, Active Directory implement unified access to information needed for remote computing

# File Sharing – Failure Modes

- All file systems have failure modes
  - ☞ For example corruption of directory structures or other non-user data, called **metadata**
- Remote file systems add new failure modes, due to network failure, server failure
- Recovery from failure can involve **state information** about status of each remote request
- **Stateless** protocols such as NFS v3 include all information in each request, allowing easy recovery but less security

# File Sharing – Consistency Semantics

- Specify how multiple users are to access a shared file simultaneously
  - ☞ Similar to Ch 5 process synchronization algorithms
    - 📄 Tend to be less complex due to disk I/O and network latency (for remote file systems)
  - ☞ Andrew File System (AFS) implemented complex remote file sharing semantics
  - ☞ Unix file system (UFS) implements:
    - 📄 Writes to an open file visible immediately to other users of the same open file
    - 📄 Sharing file pointer to allow multiple users to read and write concurrently
  - ☞ AFS has session semantics
    - 📄 Writes only visible to sessions starting after the file is closed

# Protection

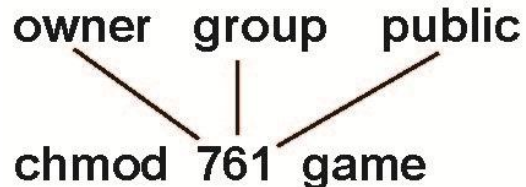
- File owner/creator should be able to control:
  - ☞ what can be done
  - ☞ by whom
- Types of access
  - ☞ **Read**
  - ☞ **Write**
  - ☞ **Execute**
  - ☞ **Append**
  - ☞ **Delete**
  - ☞ **List**

# Access Lists and Groups

- Mode of access: read, write, execute
- Three classes of users on Unix / Linux

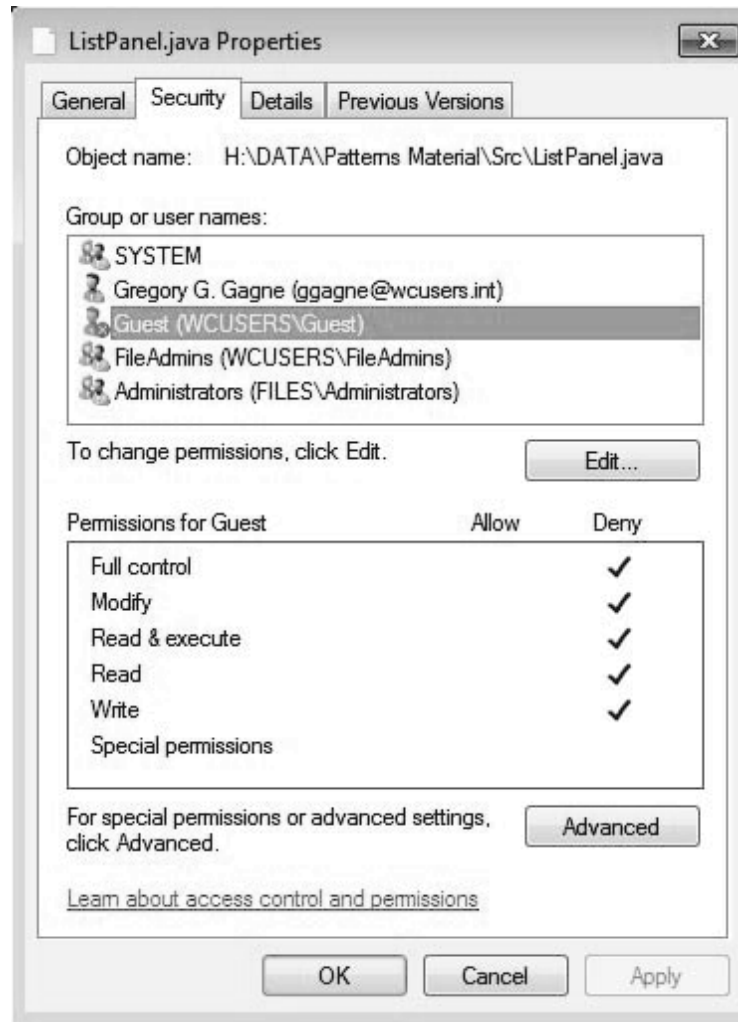
a) <b>owner access</b>	7	⇒	RWX 1 1 1
b) <b>group access</b>	6	⇒	RWX 1 1 0
c) <b>public access</b>	1	⇒	RWX 0 0 1

- Ask manager to create a group (unique name), say G, and add some users to the group.
- For a particular file (say *game*) or subdirectory, define an appropriate access.



Attach a group to a file  
chgrp G game

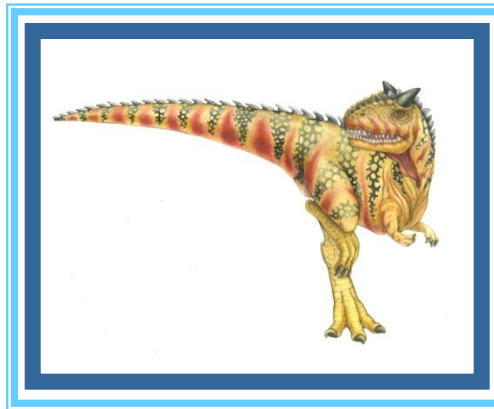
# Windows 7 Access-Control List Management



# A Sample UNIX Directory Listing

```
-rw-rw-r--  1 pbg  staff    31200  Sep 3 08:30  intro.ps
drwx-----  5 pbg  staff      512  Jul 8 09:33  private/
drwxrwxr-x  2 pbg  staff      512  Jul 8 09:35  doc/
drwxrwx---  2 pbg  student    512  Aug 3 14:13  student-proj/
-rw-r--r--  1 pbg  staff    9423  Feb 24 2003  program.c
-rwxr-xr-x  1 pbg  staff   20471  Feb 24 2003  program
drwx--x--x  4 pbg  faculty    512  Jul 31 10:31  lib/
drwx-----  3 pbg  staff    1024  Aug 29 06:52  mail/
drwxrwxrwx  3 pbg  staff      512  Jul 8 09:35  test/
```

# UNIT IV: File System Implementation





# File System Implementation

- File-System Structure
- File-System Implementation
- Directory Implementation
- Allocation Methods
- Free-Space Management
- Efficiency and Performance

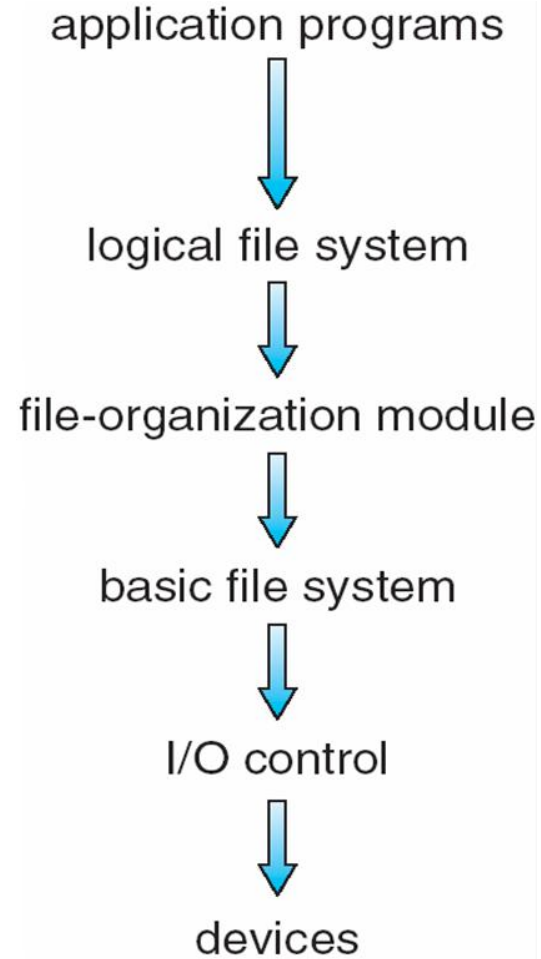
# Objectives

- To describe the details of implementing local file systems and directory structures
- To describe the implementation of remote file systems
- To discuss block allocation and free-block algorithms and trade-offs

# File-System Structure

- File structure
  - ☞ Logical storage unit
  - ☞ Collection of related information
- **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
- Disk provides in-place rewrite and random access
  - ☞ I/O transfers performed in **blocks** of **sectors** (usually 512 bytes)
- **File control block** – storage structure consisting of information about a file
- **Device driver** controls the physical device
- File system organized into layers

# Layered File System



# File System Layers

- **Device drivers** manage I/O devices at the I/O control layer
  - Given commands like “read drive1, cylinder 72, track 2, sector 10, into memory location 1060” outputs low-level hardware specific commands to hardware controller
- **Basic file system** given command like “retrieve block 123” translates to device driver
- Also manages memory buffers and caches (allocation, freeing, replacement)
  - Buffers hold data in transit
  - Caches hold frequently used data
- **File organization module** understands files, logical address, and physical blocks
- Translates logical block # to physical block #
- Manages free space, disk allocation

# File System Layers (Cont.)

- **Logical file system** manages metadata information
  - ☞ Translates file name into file number, file handle, location by maintaining file control blocks (**inodes** in UNIX)
  - ☞ Directory management
  - ☞ Protection
- Layering useful for reducing complexity and redundancy, but adds overhead and can decrease performance
  - ☞ Translates file name into file number, file handle, location by maintaining file control blocks (**inodes** in UNIX)
  - ☞ Logical layers can be implemented by any coding method according to OS designer

# File System Layers (Cont.)

- Many file systems, sometimes many within an operating system
  - Each with its own format (CD-ROM is ISO 9660; Unix has **UFS**, FFS; Windows has FAT, FAT32, NTFS as well as floppy, CD, DVD Blu-ray, Linux has more than 40 types, with **extended file system** ext2 and ext3 leading; plus distributed file systems, etc.)
  - New ones still arriving – ZFS, GoogleFS, Oracle ASM, FUSE

# File-System Implementation

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



# File-System Implementation (Cont.)

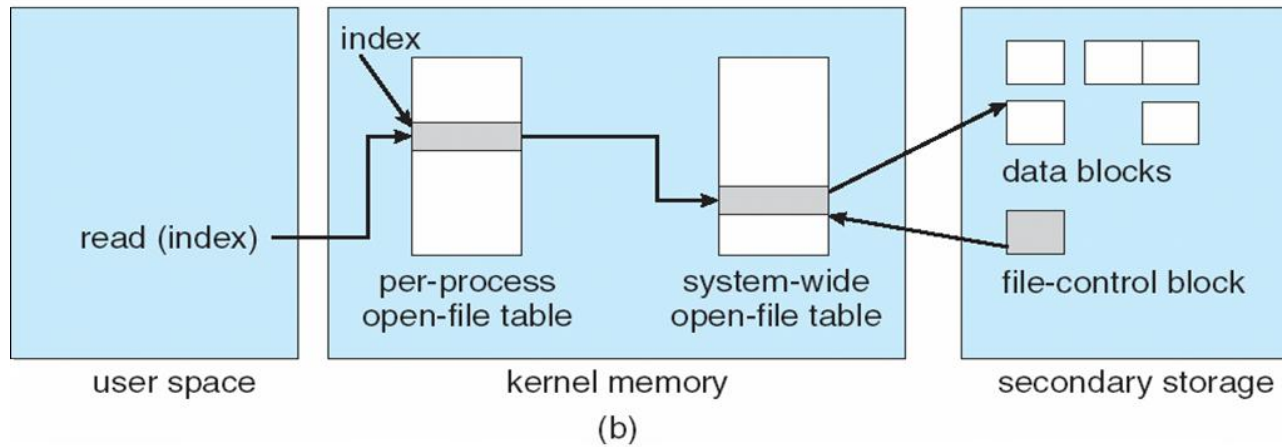
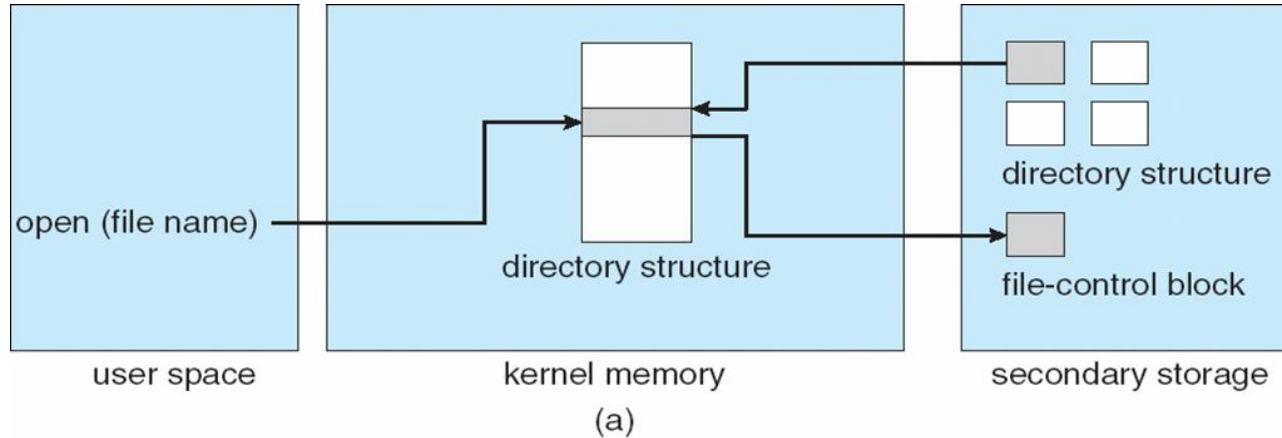
- Per-file **File Control Block (FCB)** contains many details about the file
  - ☞ inode number, permissions, size, dates
  - ☞ NFTS stores into in master file table using relational DB structures

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

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

# In-Memory File System Structures



# Partitions and Mounting

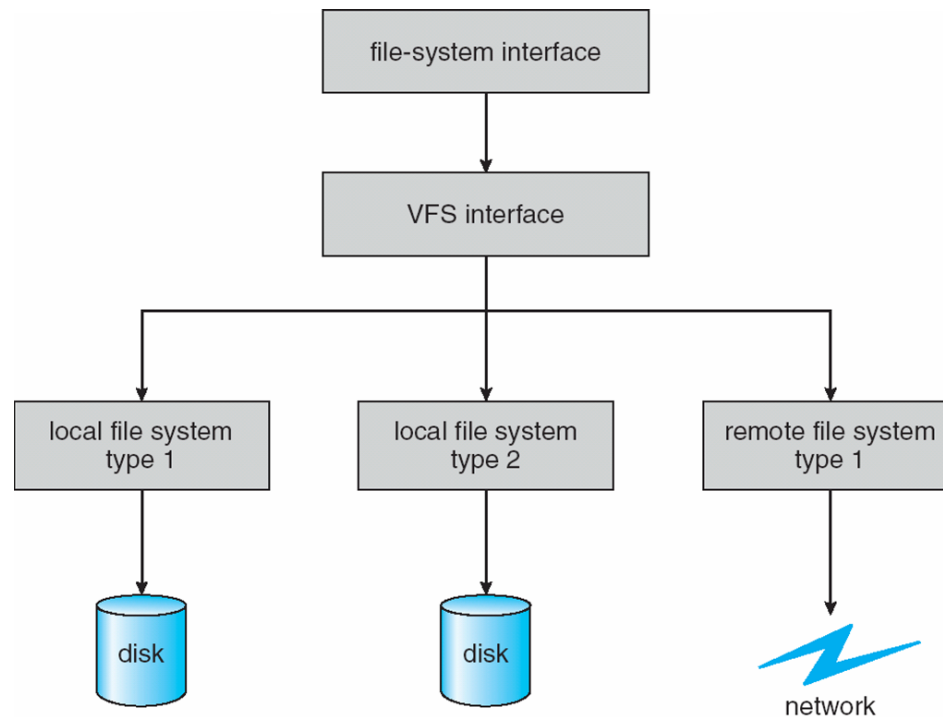
- Partition can be a volume containing a file system (“cooked”) or **raw** – just a sequence of blocks with no file system
- Boot block can point to boot volume or boot loader set of blocks that contain enough code to know how to load the kernel from the file system
  - ☞ Or a boot management program for multi-os booting
- **Root partition** contains the OS, other partitions can hold other Oses, other file systems, or be raw
  - ☞ Mounted at boot time
  - ☞ Other partitions can mount automatically or manually
- At mount time, file system consistency checked
  - ☞ Is all metadata correct?
    - 📄 If not, fix it, try again
    - 📄 If yes, add to mount table, allow access

# Virtual File Systems

- **Virtual File Systems (VFS)** on Unix provide an object-oriented way of implementing file systems
- 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 file system
    - 📄 Implements **vnodes** which hold inodes or network file details
  - ☞ Then dispatches operation to appropriate file system implementation routines

# Virtual File Systems (Cont.)

- The API is to the VFS interface, rather than any specific type of file system



# Virtual File System Implementation

- For example, Linux has four object types:
  - ☞ inode, file, superblock, dentry
- VFS defines set of operations on the objects that must be implemented
  - ☞ Every object has a pointer to a function table
    - ☞ Function table has addresses of routines to implement that function on that object
    - ☞ For example:
      - ☞ `• int open(. . .)`—Open a file
      - ☞ `• int close(. . .)`—Close an already-open file
      - ☞ `• ssize_t read(. . .)`—Read from a file
      - ☞ `• ssize_t write(. . .)`—Write to a file
      - ☞ `• int mmap(. . .)`—Memory-map a file

# Directory Implementation

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

- ☞ Simple to program

- ☞ Time-consuming to execute

- 📄 Linear search time

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

- **Hash Table** – linear list with hash data structure

- ☞ 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

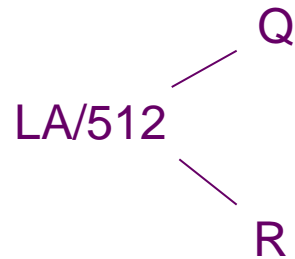


# Allocation Methods - Contiguous

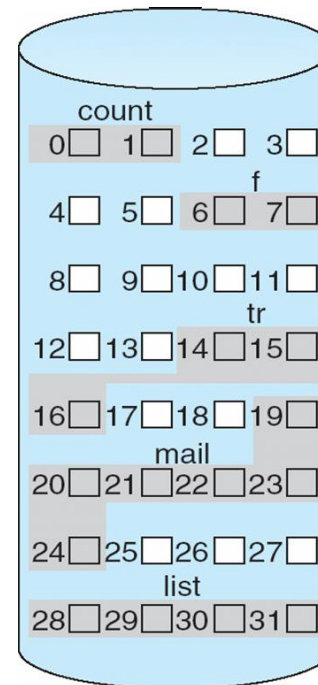
- An allocation method refers to how disk blocks are allocated for files:
- **Contiguous allocation** – each file occupies set of contiguous blocks
  - ☞ Best performance in most cases
  - ☞ Simple – only starting location (block #) and length (number of blocks) are required
  - ☞ Problems include finding space for file, knowing file size, external fragmentation, need for **compaction off-line (downtime)** or **on-line**

# Contiguous Allocation

- Mapping from logical to physical



Block to be accessed =  $Q +$   
starting address  
Displacement into block =  $R$



directory

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

# Extent-Based Systems

- Many newer file systems (i.e., Veritas File System) use a modified contiguous allocation scheme
- Extent-based file systems allocate disk blocks in extents
- An **extent** is a contiguous block of disks
  - ☞ Extents are allocated for file allocation
  - ☞ A file consists of one or more extents

# Allocation Methods - Linked

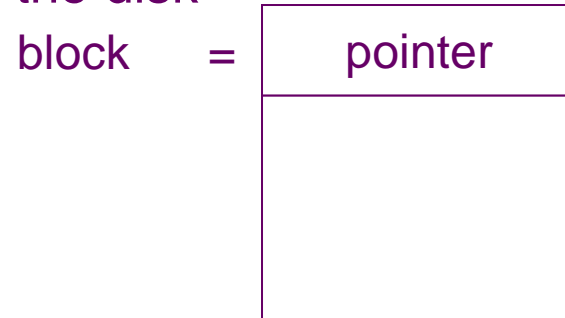
- **Linked allocation** – each file a linked list of blocks
  - ☞ File ends at nil pointer
  - ☞ No external fragmentation
  - ☞ Each block contains pointer to next block
  - ☞ No compaction, external fragmentation
  - ☞ Free space management system called when new block needed
  - ☞ Improve efficiency 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

# Allocation Methods – Linked (Cont.)

- 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

# Linked Allocation

- Each file is a linked list of disk blocks: blocks may be scattered anywhere on the disk



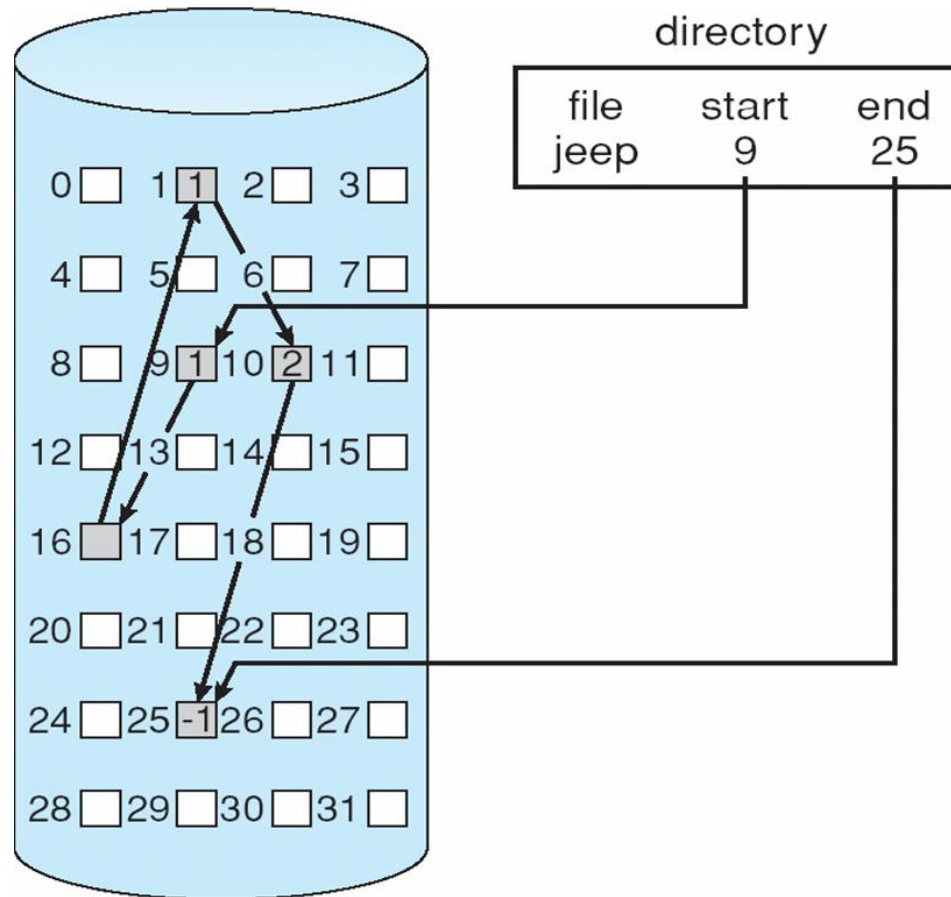
- Mapping



Block to be accessed is the Qth block in the linked chain of blocks representing the file.

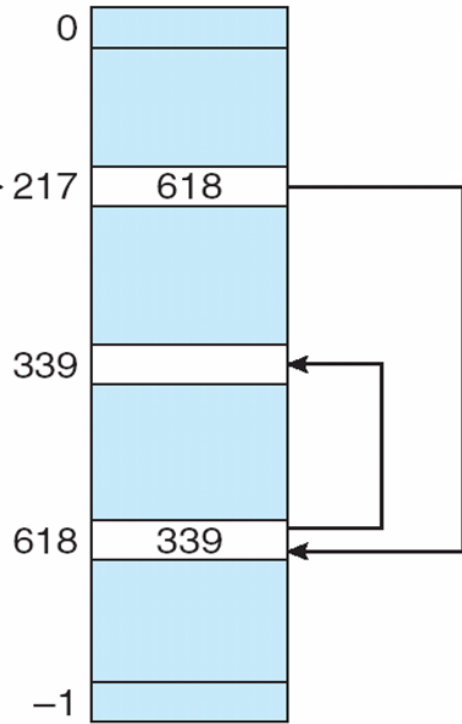
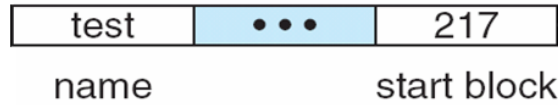
$$\text{Displacement into block} = R + 1$$

# Linked Allocation



# File-Allocation Table

directory entry



no. of disk blocks

FAT

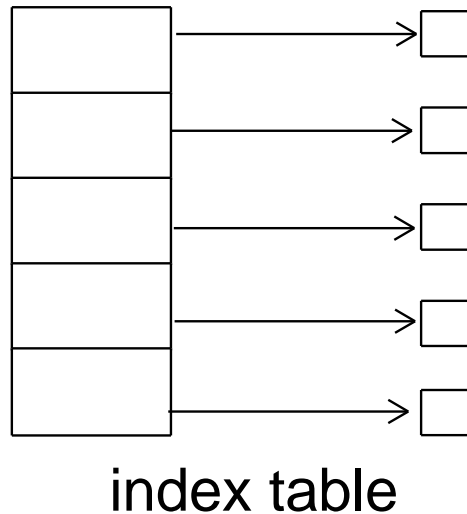


# Allocation Methods - Indexed

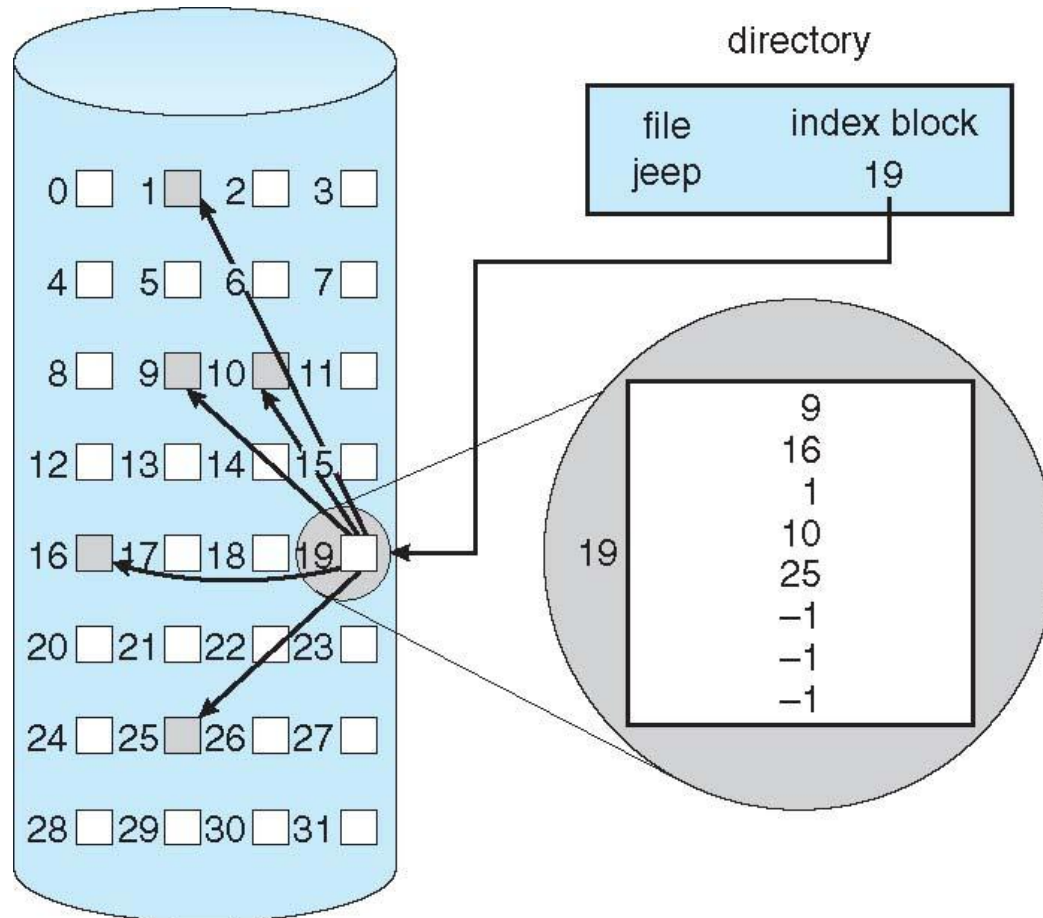
## ■ Indexed allocation

- ☞ Each file has its own **index block**(s) of pointers to its data blocks

## ■ Logical view



# Example of Indexed Allocation



# Indexed Allocation (Cont.)

- Need index table
- Random access
- Dynamic access without external fragmentation, but have overhead of index block
- Mapping from logical to physical in a file of maximum size of 256K bytes and block size of 512 bytes. We need only 1 block for index table



Q = displacement into index table

R = displacement into block

# Indexed Allocation – Mapping (Cont.)

- Mapping from logical to physical in a file of unbounded length (block size of 512 words)
- Linked scheme – Link blocks of index table (no limit on size)

$$LA / (512 \times 511) \begin{cases} Q_1 \\ R_1 \end{cases}$$

$Q_1$  = block of index table

$R_1$  is used as follows:

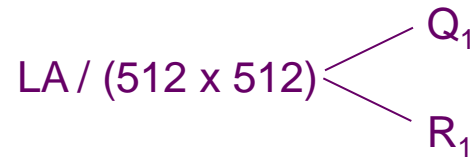
$$R_1 / 512 \begin{cases} Q_2 \\ R_2 \end{cases}$$

$Q_2$  = displacement into block of index table

$R_2$  displacement into block of file:

# Indexed Allocation – Mapping (Cont.)

- Two-level index (4K blocks could store 1,024 four-byte pointers in outer index -> 1,048,567 data blocks and file size of up to 4GB)

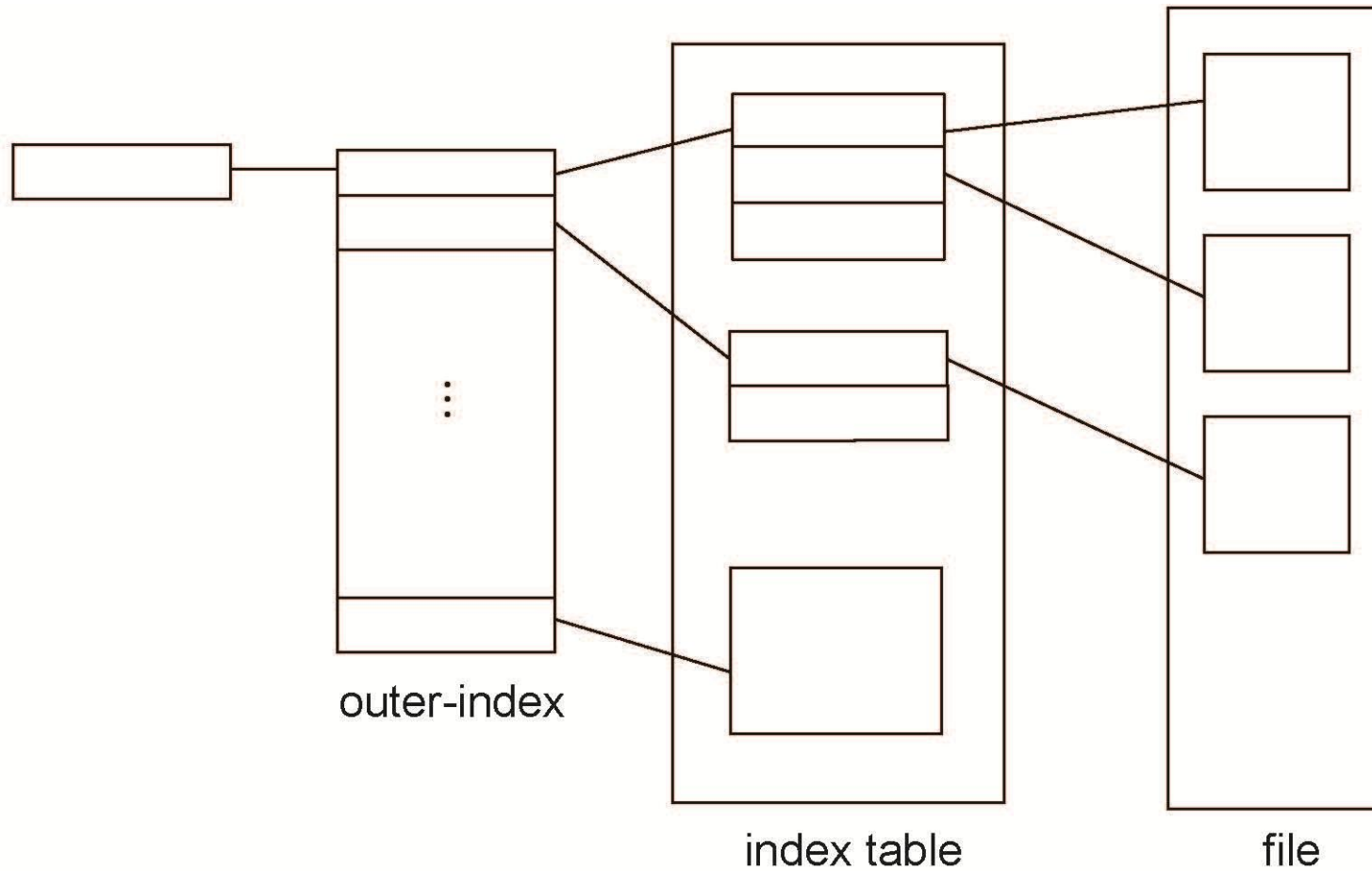


$Q_1$  = displacement into outer-index  
 $R_1$  is used as follows:



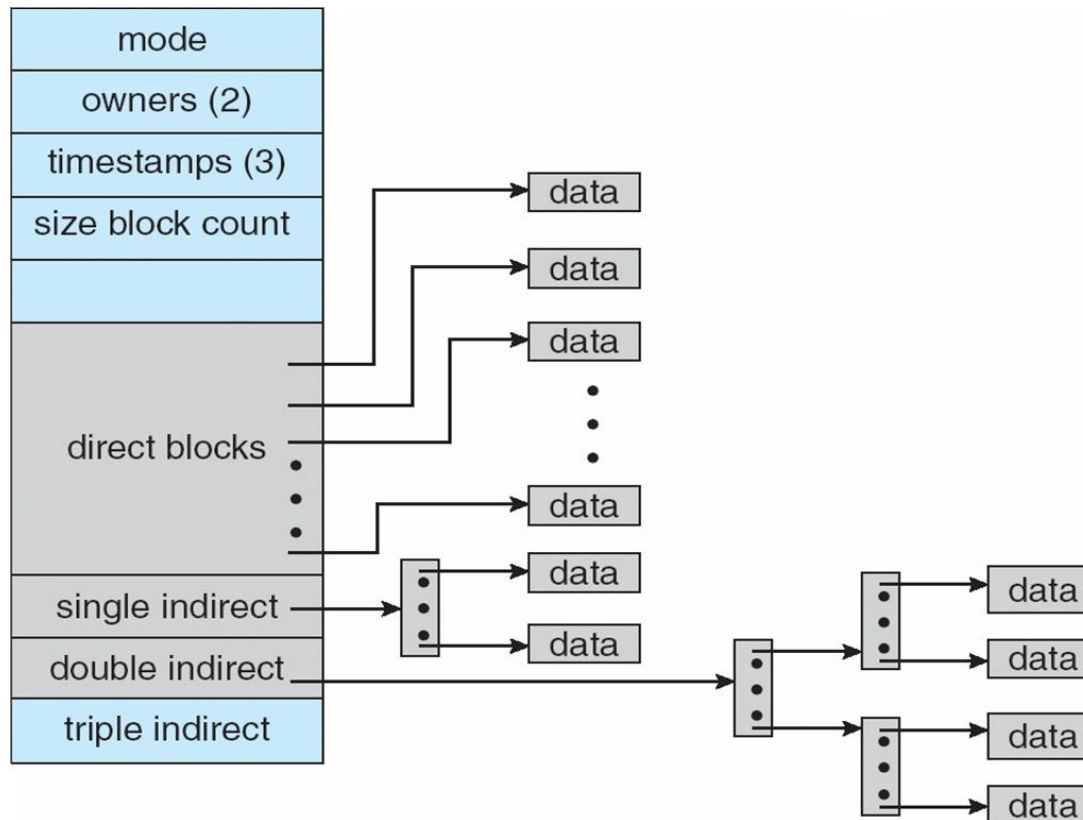
$Q_2$  = displacement into block of index table  
 $R_2$  displacement into block of file:

# Indexed Allocation – Mapping (Cont.)



# Combined Scheme: UNIX UFS

4K bytes per block, 32-bit addresses



More index blocks than can be addressed with 32-bit file pointer

# Performance

- Best method depends on file access type
  - ☞ Contiguous great for sequential and random
- Linked good for sequential, not random
- Declare access type at creation -> select either contiguous or linked
- Indexed more complex
  - ☞ Single block access could require 2 index block reads then data block read
  - ☞ Clustering can help improve throughput, reduce CPU overhead



# Performance (Cont.)

- Adding instructions to the execution path to save one disk I/O is reasonable
  - ☞ Intel Core i7 Extreme Edition 990x (2011) at 3.46Ghz = 159,000 MIPS
    - 📄 [http://en.wikipedia.org/wiki/Instructions\\_per\\_second](http://en.wikipedia.org/wiki/Instructions_per_second)
  - ☞ Typical disk drive at 250 I/Os per second
    - 📄  $159,000 \text{ MIPS} / 250 = 630 \text{ million instructions during one disk I/O}$
  - ☞ Fast SSD drives provide 60,000 IOPS
    - 📄  $159,000 \text{ MIPS} / 60,000 = 2.65 \text{ millions instructions during one disk I/O}$

# Free-Space Management

- File system maintains **free-space list** to track available blocks/clusters
  - ☞ (Using term “block” for simplicity)

- **Bit vector** or **bit map** ( $n$  blocks)



$$\text{bit}[i] = \begin{cases} 1 \Rightarrow \text{block}[i] \text{ free} \\ 0 \Rightarrow \text{block}[i] \text{ occupied} \end{cases}$$

Block number calculation

$$\begin{aligned} & (\text{number of bits per word}) * \\ & (\text{number of 0-value words}) + \\ & \text{offset of first 1 bit} \end{aligned}$$

CPUs have instructions to return offset within word of first “1” bit

# Free-Space Management (Cont.)

- Bit map requires extra space

☞ Example:

block size = 4KB =  $2^{12}$  bytes

disk size =  $2^{40}$  bytes (1 terabyte)

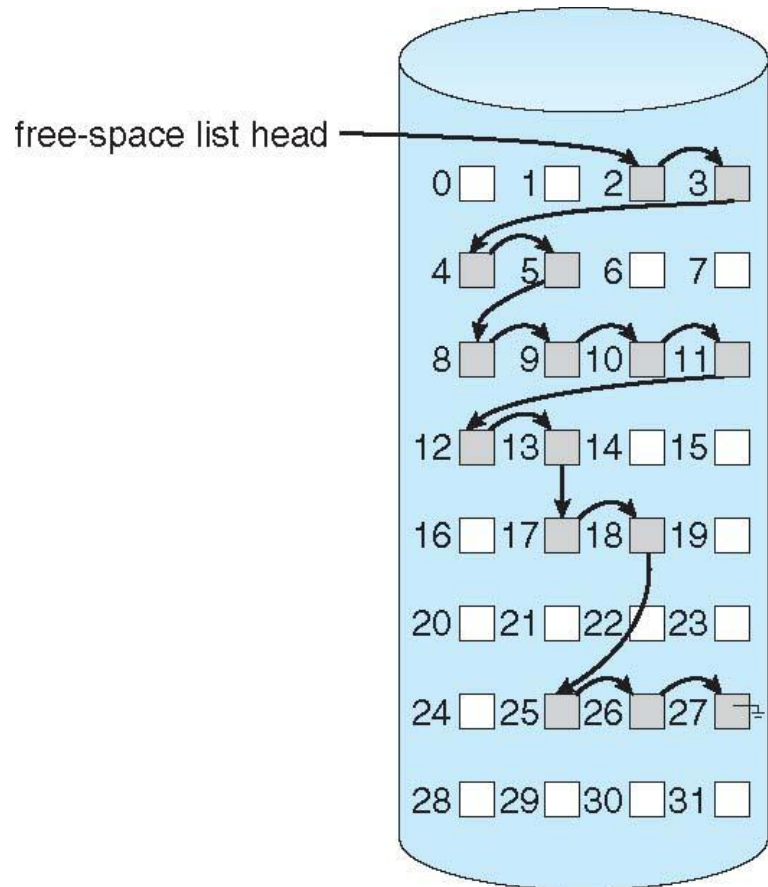
$n = 2^{40}/2^{12} = 2^{28}$  bits (or 32MB)

if clusters of 4 blocks -> 8MB of memory

- Easy to get contiguous files

# Linked Free Space List on Disk

- Linked list (free list)
  - Cannot get contiguous space easily
  - No waste of space
  - No need to traverse the entire list (if # free blocks recorded)



# Free-Space Management (Cont.)

## ■ Grouping

- ☞ 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

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

# Free-Space Management (Cont.)

## ■ Space Maps

- 👉 Used in **ZFS**
- 👉 Consider meta-data I/O on very large file systems
  - 📄 Full data structures like bit maps couldn't fit in memory -> thousands of I/Os
- 👉 Divides device space into **metaslab** units and manages metaslabs
  - 📄 Given volume can contain hundreds of metaslabs
- 👉 Each metaslab has associated space map
  - 📄 Uses counting algorithm
- 👉 But records to log file rather than file system
  - 📄 Log of all block activity, in time order, in counting format
- 👉 Metaslab activity -> load space map into memory in balanced-tree structure, indexed by offset
  - 📄 Replay log into that structure
  - 📄 Combine contiguous free blocks into single entry

# Efficiency and Performance

- Efficiency dependent on:
  - ☞ Disk allocation and directory algorithms
  - ☞ Types of data kept in file's directory entry
  - ☞ Pre-allocation or as-needed allocation of metadata structures
  - ☞ Fixed-size or varying-size data structures

# Efficiency and Performance (Cont.)

## ■ Performance

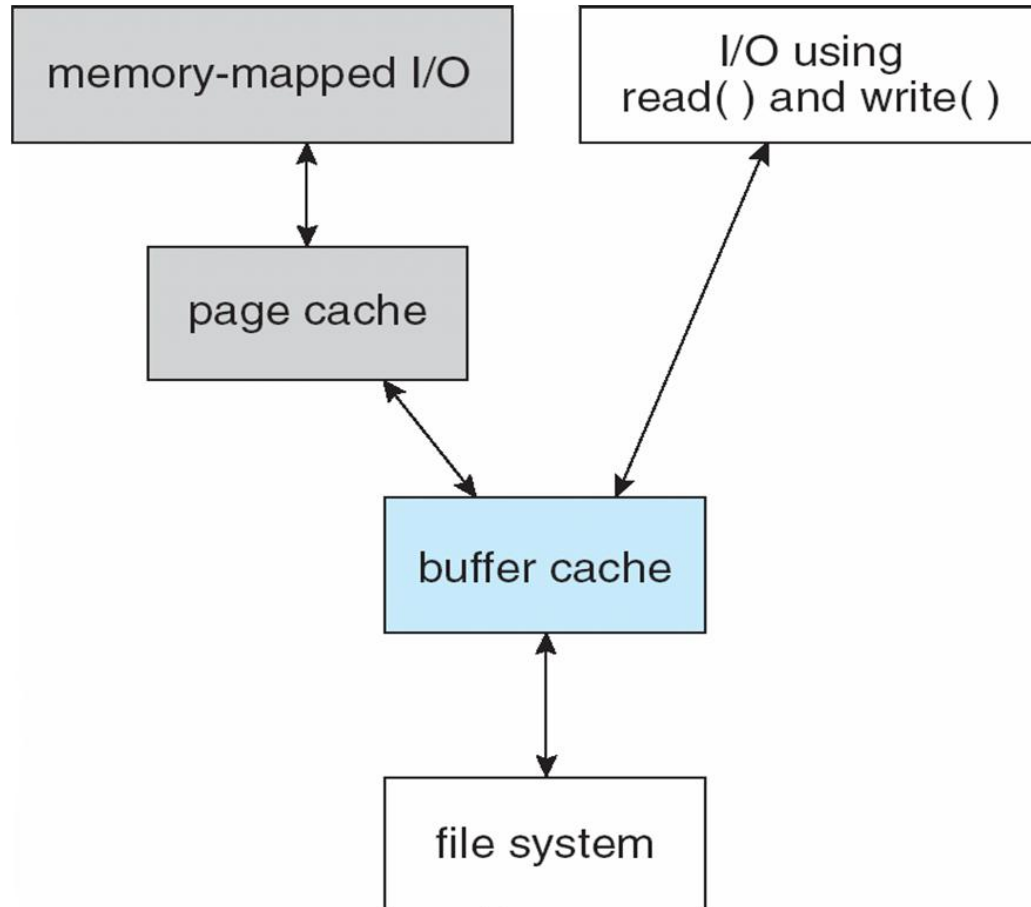
- ☞ Keeping data and metadata close together
- ☞ **Buffer cache** – separate section of main memory for frequently used blocks
- ☞ **Synchronous** writes sometimes requested by apps or needed by OS
  - 📄 No buffering / caching – writes must hit disk before acknowledgement
  - 📄 **Asynchronous** writes more common, buffer-able, faster
- ☞ **Free-behind** and **read-ahead** – techniques to optimize sequential access
- ☞ Reads frequently slower than writes



# Page Cache

- A **page cache** caches pages rather than disk blocks using virtual memory techniques and addresses
- Memory-mapped I/O uses a page cache
- Routine I/O through the file system uses the buffer (disk) cache
- This leads to the following figure

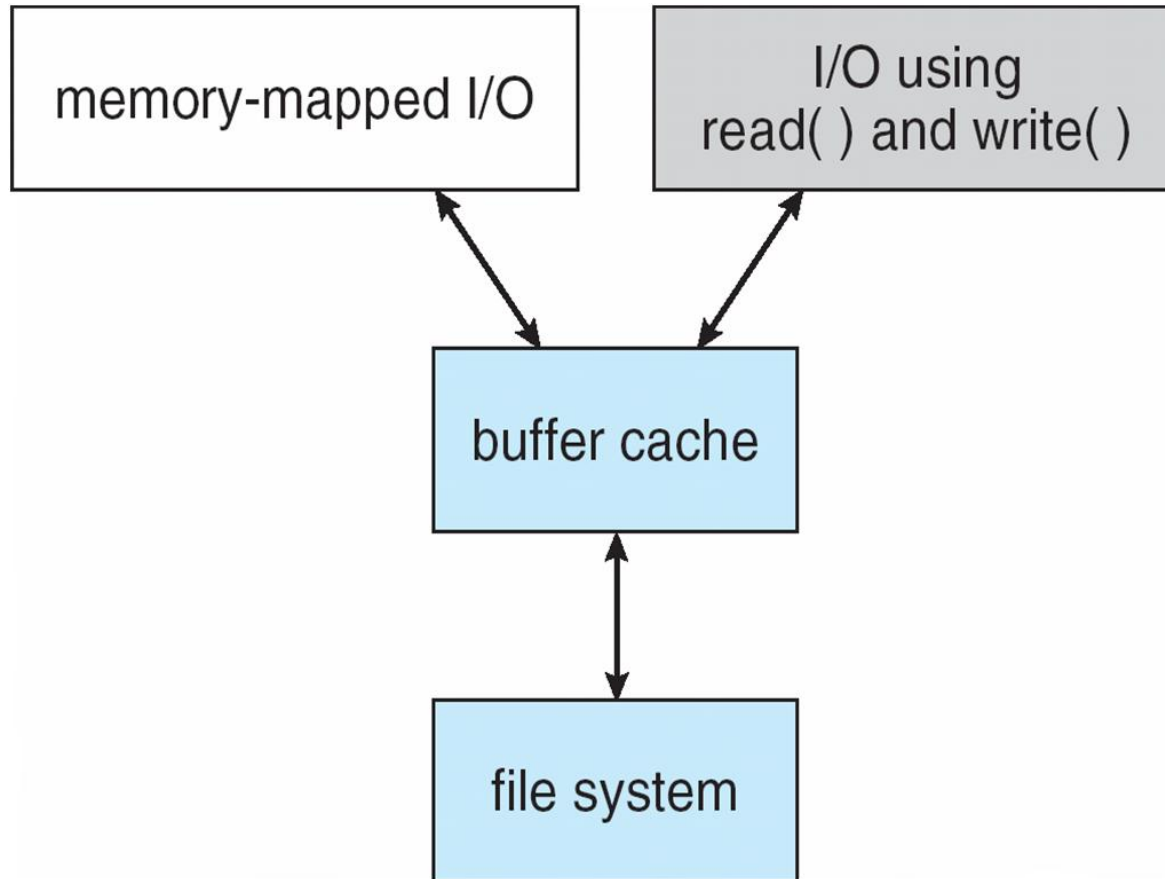
# I/O Without a Unified Buffer Cache



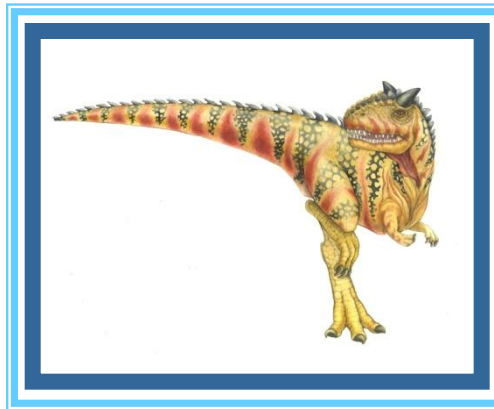
# Unified Buffer Cache

- A **unified buffer cache** uses the same page cache to cache both memory-mapped pages and ordinary file system I/O to avoid **double caching**
- But which caches get priority, and what replacement algorithms to use?

# I/O Using a Unified Buffer Cache



# UNIT IV: Mass-Storage Systems



# Mass-Storage Systems

- Overview of Mass Storage Structure
- Disk Structure
- Disk Attachment
- Disk Scheduling
- Disk Management
- Swap-Space Management

# Objectives

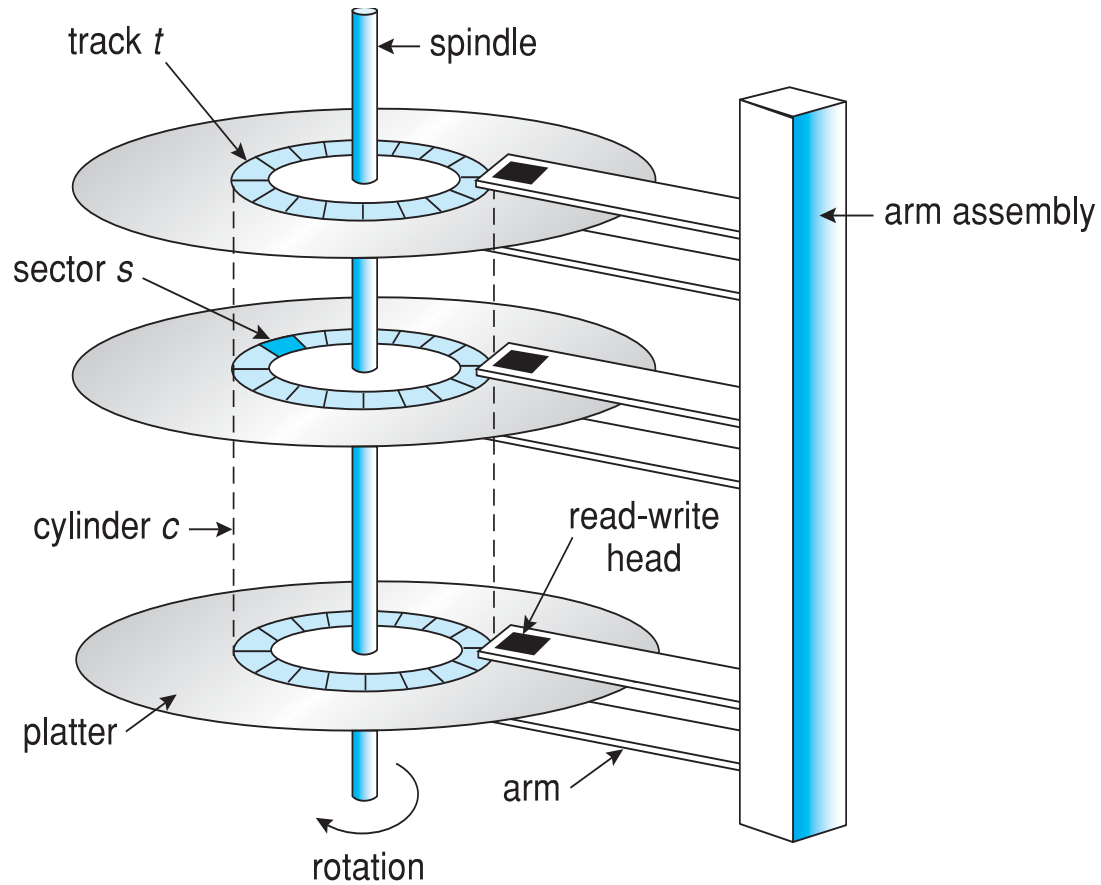
- To describe the physical structure of secondary storage devices and its effects on the uses of the devices
- To explain the performance characteristics of mass-storage devices
- To evaluate disk scheduling algorithms
- To discuss operating-system services provided for mass storage, including RAID

# Overview of Mass Storage Structure

- **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 -- That's bad
- Disks can be removable
- 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



# Moving-head Disk Mechanism



# Hard Disks

- Platters range from .85” to 14” (historically)
  - ☞ Commonly 3.5”, 2.5”, and 1.8”
- Range from 30GB to 3TB per drive
- Performance
  - ☞ Transfer Rate – theoretical – 6 Gb/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
    - 📄  $1 / (\text{RPM} / 60) = 60 / \text{RPM}$
  - ☞ Average latency =  $\frac{1}{2}$  latency

Spindle [rpm]	Average latency [ms]
4200	7.14
5400	5.56
7200	4.17
10000	3
15000	2

(From Wikipedia)

# Hard Disk Performance

- **Access Latency = Average access time** = average seek time + average latency
  - ☞ For fastest disk  $3\text{ms} + 2\text{ms} = 5\text{ms}$
  - ☞ For slow disk  $9\text{ms} + 5.56\text{ms} = 14.56\text{ms}$
- Average I/O time = average access time + (amount to transfer / transfer rate) + controller overhead
- For example to transfer a 4KB block on a 7200 RPM disk with a 5ms average seek time, 1Gb/sec transfer rate with a .1ms controller overhead =
  - ☞  $5\text{ms} + 4.17\text{ms} + 0.1\text{ms} + \text{transfer time} =$
  - ☞ Transfer time =  $4\text{KB} / 1\text{Gb/s} * 8\text{Gb} / \text{GB} * 1\text{GB} / 1024^2\text{KB} = 32 / (1024^2) = 0.031 \text{ ms}$
  - ☞ Average I/O time for 4KB block =  $9.27\text{ms} + .031\text{ms} = 9.301\text{ms}$

# The First Commercial Disk Drive



1956

IBM RAMDAC computer  
included the IBM Model 350  
disk storage system

5M (7 bit) characters

50 x 24" platters

Access time = < 1 second

# Solid-State Disks

- Nonvolatile memory used like a hard drive
  - ☞ Many technology variations
- Can be more reliable than HDDs
- More expensive per MB
- Maybe have shorter life span
- Less capacity
- But much faster
- Busses can be too slow -> connect directly to PCI for example
- No moving parts, so no seek time or rotational latency

# Magnetic Tape

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

# Disk Structure

- Disk drives are addressed as large 1-dimensional arrays of **logical blocks**, where the logical block is the smallest unit of transfer
  - ☞ Low-level formatting creates **logical blocks** on physical media
- 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
    - 📄 Except for bad sectors
    - 📄 Non-constant # of sectors per track via constant angular velocity

# Disk Attachment

- Host-attached storage accessed through I/O ports talking to I/O busses
- 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
- I/O directed to bus ID, device ID, logical unit (LUN)

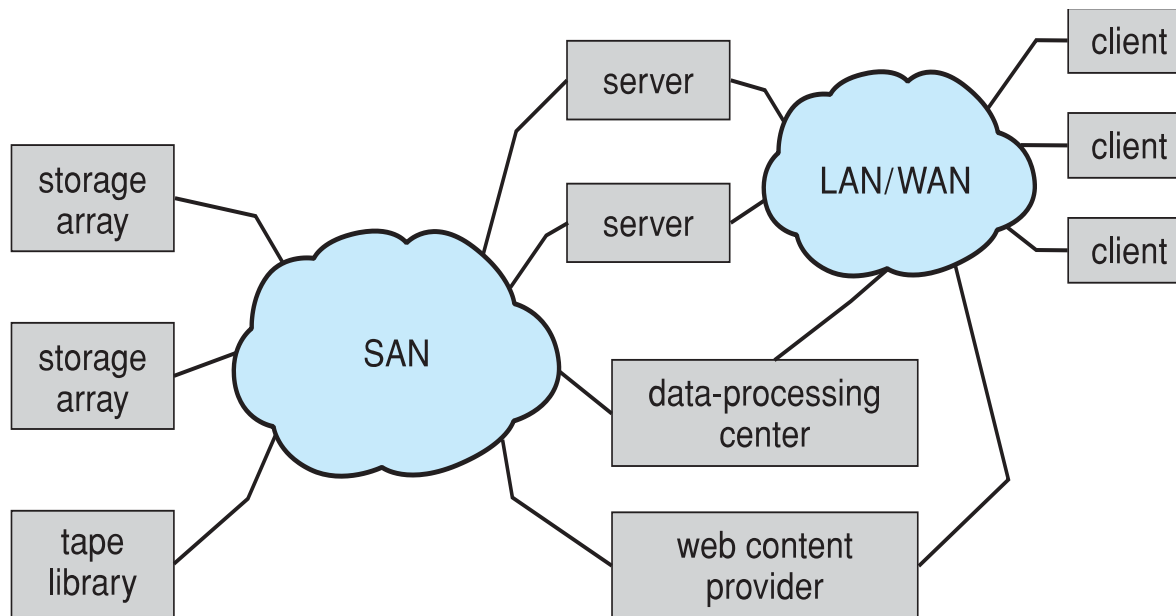


# Storage Array

- Can just attach disks, or arrays of disks
- Storage Array has controller(s), provides features to attached host(s)
  - ☞ Ports to connect hosts to array
  - ☞ Memory, controlling software (sometimes NVRAM, etc)
  - ☞ A few to thousands of disks
  - ☞ RAID, hot spares, hot swap (discussed later)
  - ☞ Shared storage -> more efficiency
  - ☞ Features found in some file systems
    - 📄 Snapshots, clones, thin provisioning, replication, deduplication, etc

# Storage Area Network

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

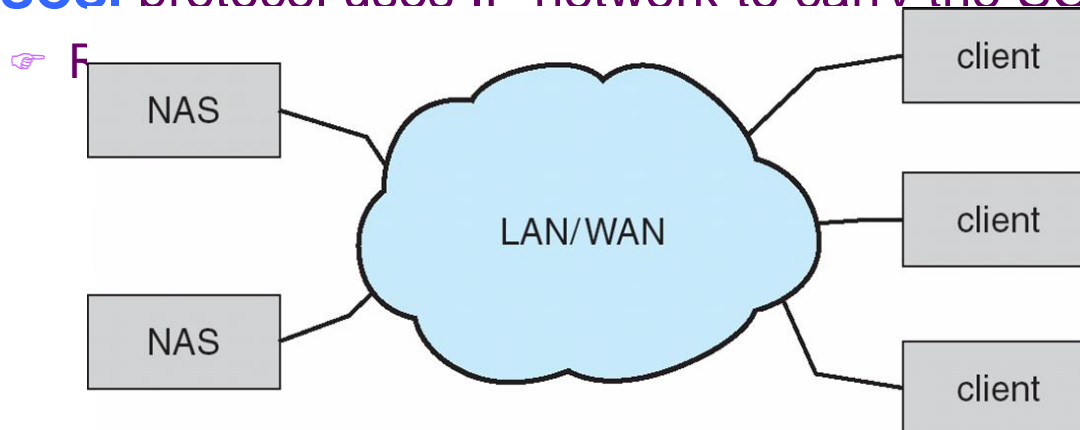


# Storage Area Network (Cont.)

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

# Network-Attached Storage

- 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
- NFS and CIFS are common protocols
- Implemented via remote procedure calls (RPCs) between host and storage over typically TCP or UDP on IP network
- **iSCSI** protocol uses IP network to carry the SCSI protocol



# Disk Scheduling

- The operating system is responsible for using hardware efficiently — for the disk drives, this means having a fast access time and disk bandwidth
- Minimize seek time
- Seek time  $\approx$  seek distance
- 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

# Disk Scheduling (Cont.)

- There are many sources of disk I/O request
  - ☞ OS
  - ☞ System processes
  - ☞ Users processes
- I/O request includes input or output mode, disk address, memory address, number of sectors to transfer
- OS maintains queue of requests, per disk or device
- Idle disk can immediately work on I/O request, busy disk means work must queue
  - ☞ Optimization algorithms only make sense when a queue exists

# Disk Scheduling (Cont.)

- Note that drive controllers have small buffers and can manage a queue of I/O requests (of varying “depth”)
- Several algorithms exist to schedule the servicing of disk I/O requests
- The analysis is true for one or many platters
- We illustrate scheduling algorithms with a request queue (0-199)

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

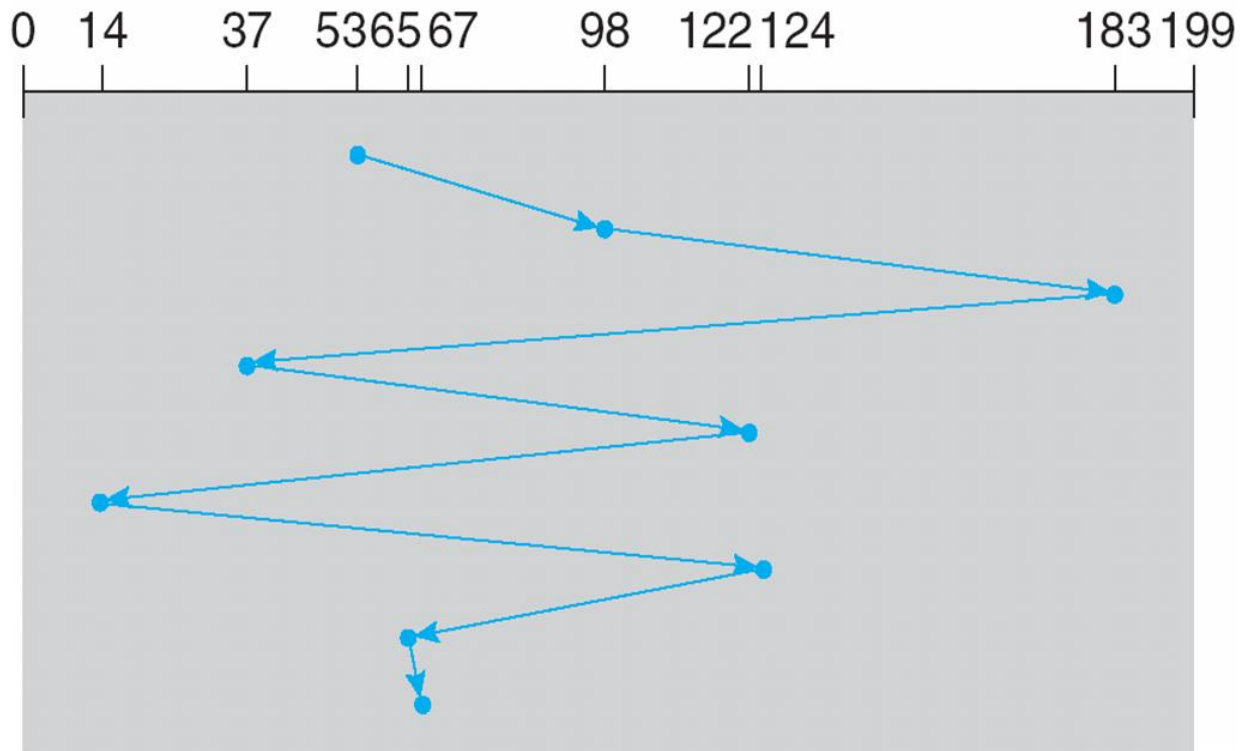
Head pointer 53

# FCFS

Illustration shows total head movement of 640 cylinders

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

head starts at 53

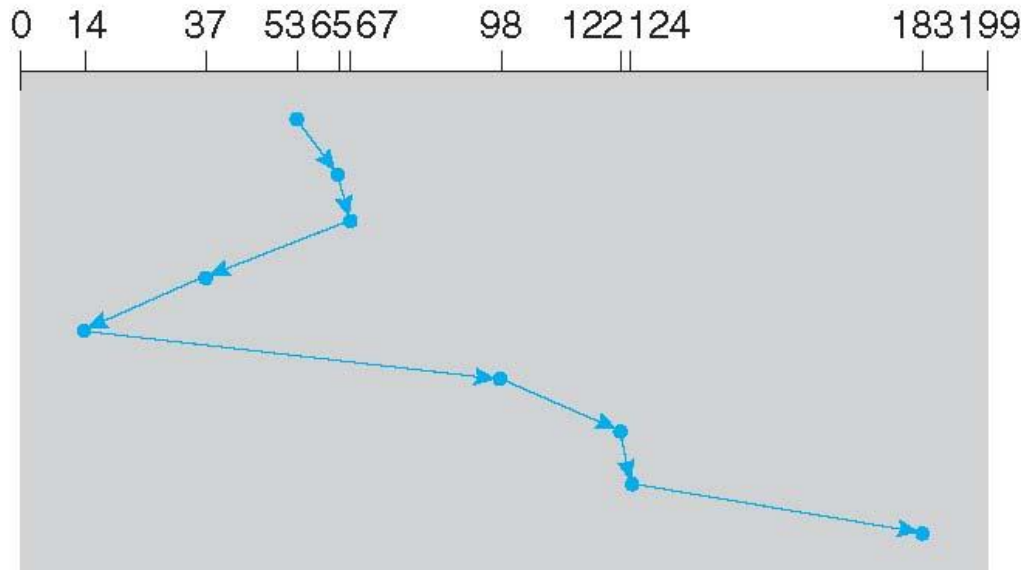




# SSTF

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

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

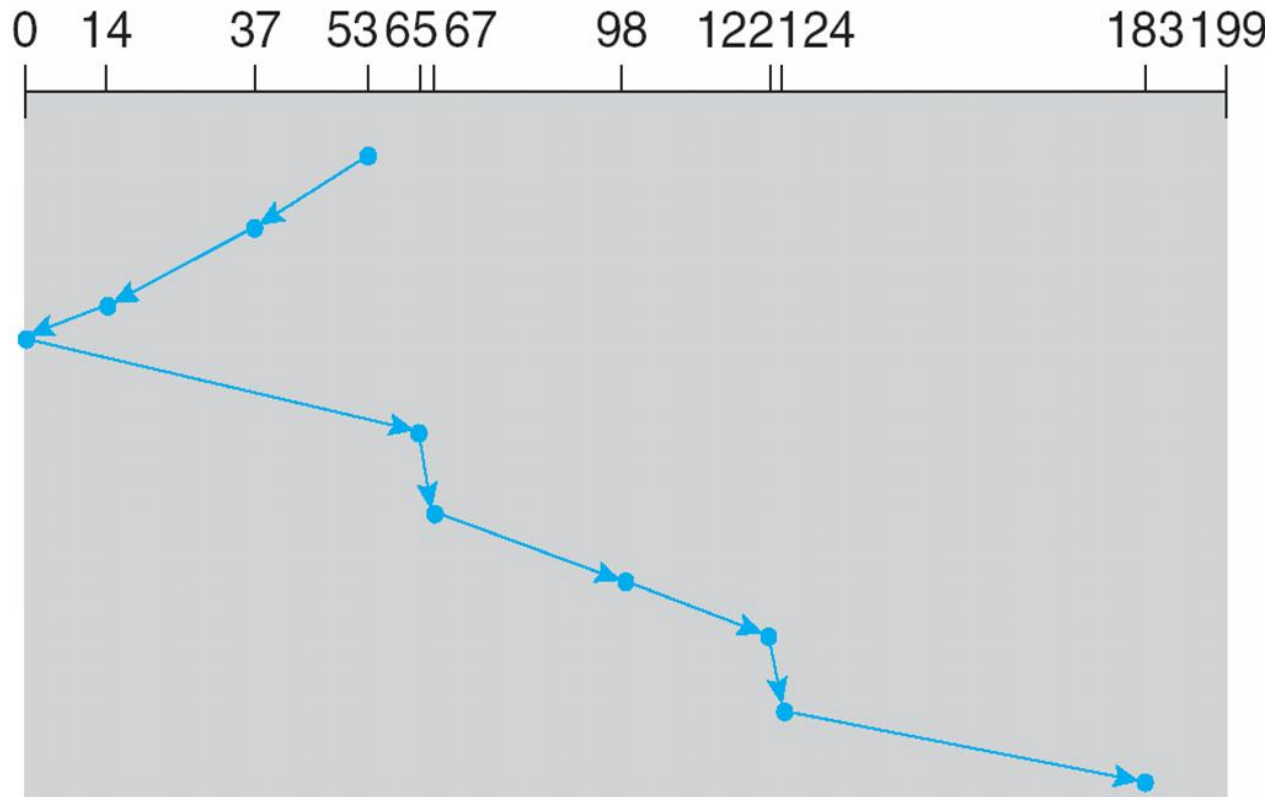


# SCAN

- 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.
- **SCAN algorithm** Sometimes called the **elevator algorithm**
- Illustration shows total head movement of 208 cylinders
- But note that if requests are uniformly dense, largest density at other end of disk and those wait the longest

# SCAN (Cont.)

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



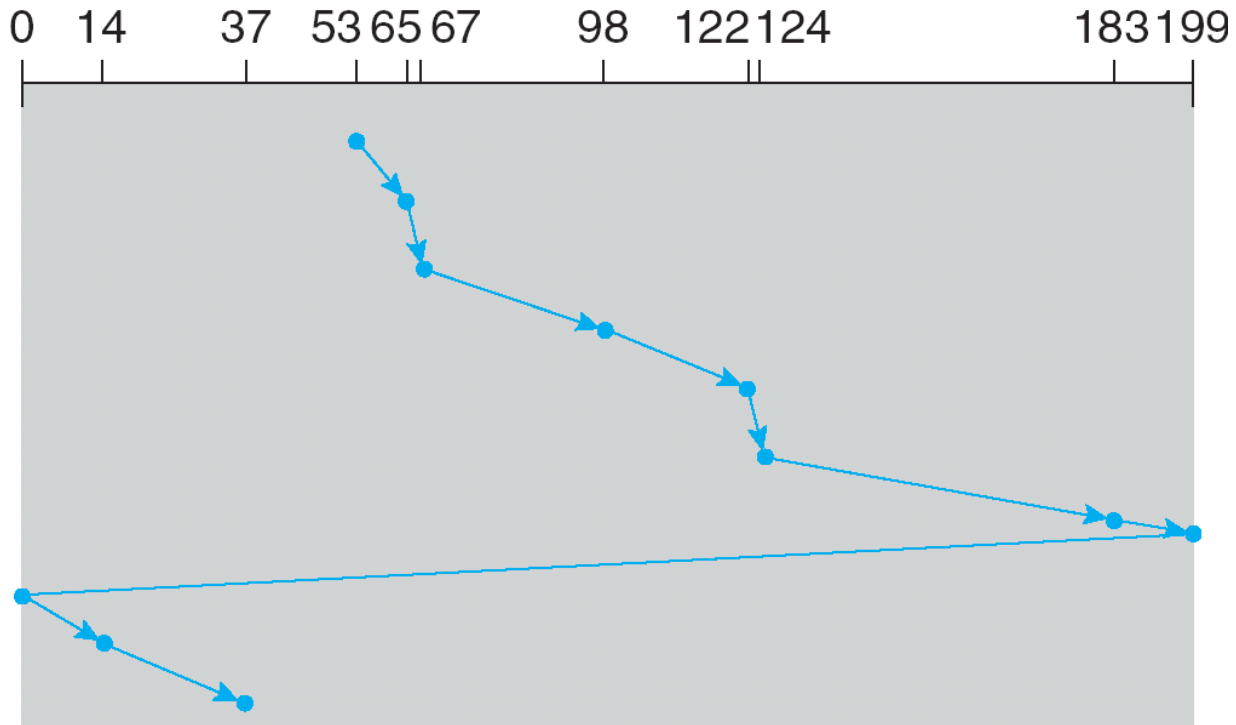
# C-SCAN

- Provides a more uniform wait time than SCAN
- 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
- Total number of cylinders?

# C-SCAN (Cont.)

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

head starts at 53



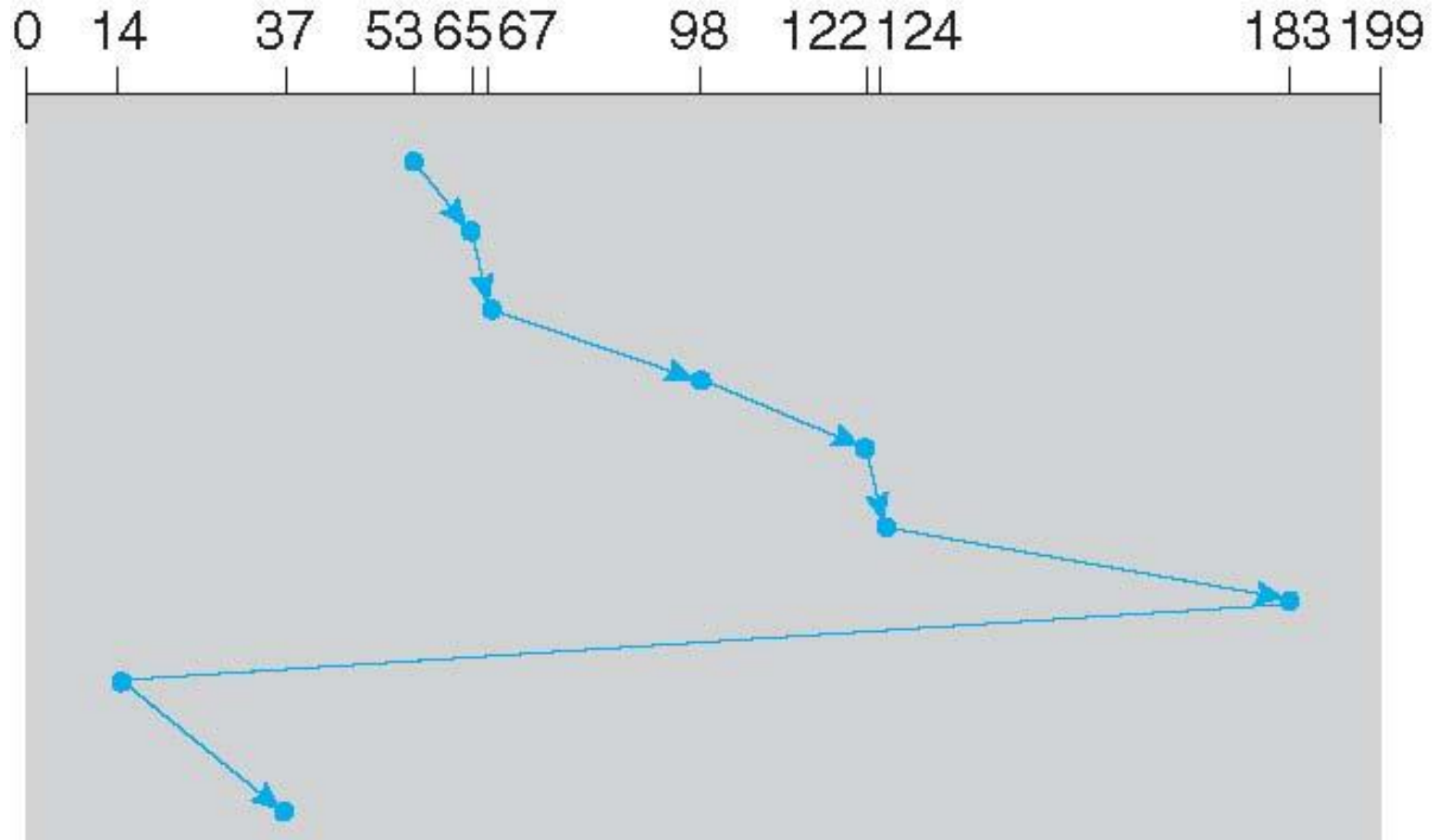
# C-LOOK

- LOOK a version of SCAN, C-LOOK a version of C-SCAN
- 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
- Total number of cylinders?

# C-LOOK (Cont.)

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

head starts at 53



# Selecting a Disk-Scheduling Algorithm

- SSTF is common and has a natural appeal
- SCAN and C-SCAN perform better for systems that place a heavy load on the disk
  - ☞ Less starvation
- Performance depends on the number and types of requests
- Requests for disk service can be influenced by the file-allocation method
  - ☞ And metadata layout
- The disk-scheduling algorithm should be written as a separate module of the operating system, allowing it to be replaced with a different algorithm if necessary
- Either SSTF or LOOK is a reasonable choice for the default algorithm
- What about rotational latency?
  - ☞ Difficult for OS to calculate
- How does disk-based queueing effect OS queue ordering efforts?



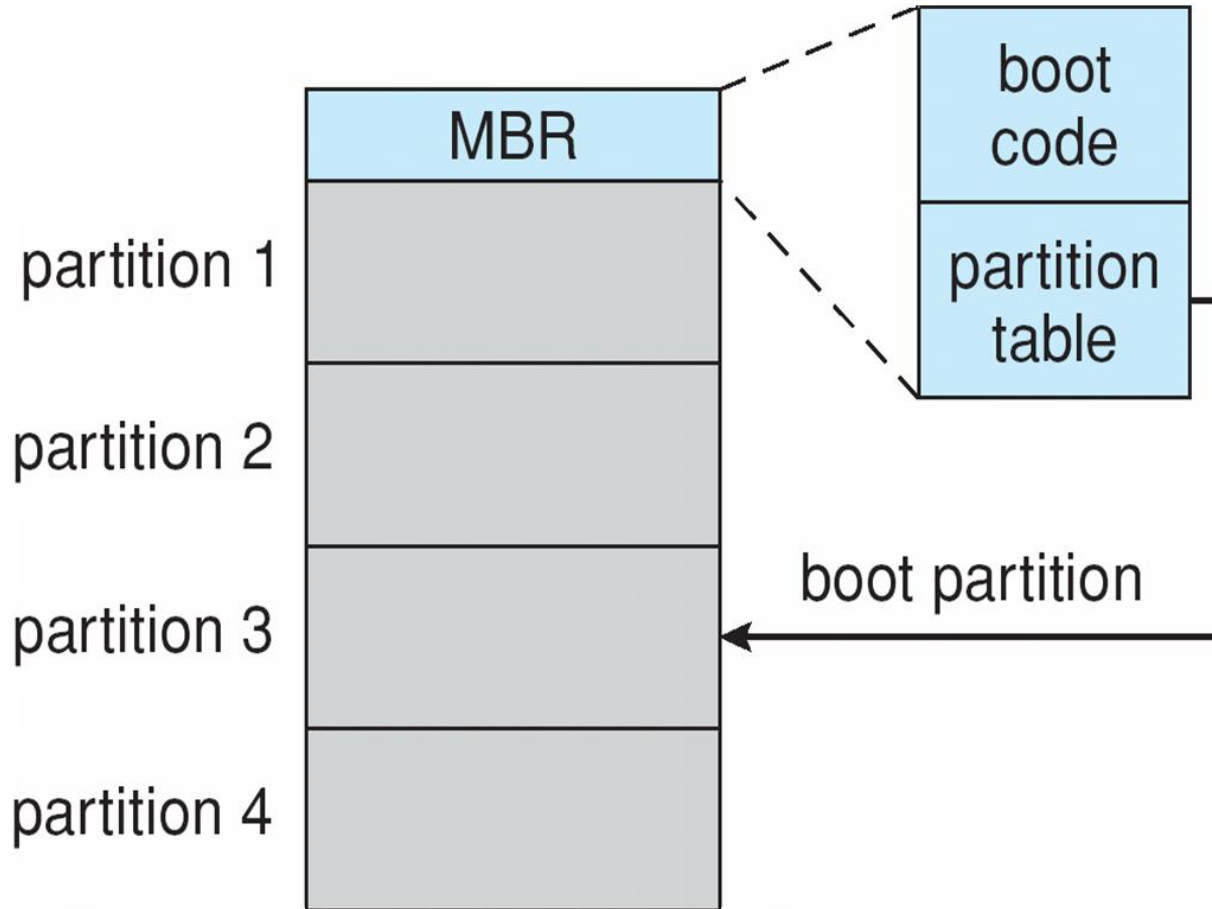
# Disk Management

- **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 file system”
  - ☞ To increase efficiency most file systems group blocks into **clusters**
    - 📄 Disk I/O done in blocks
    - 📄 File I/O done in clusters

# Disk Management (Cont.)

- Raw disk access for apps that want to do their own block management, keep OS out of the way (databases for example)
- 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

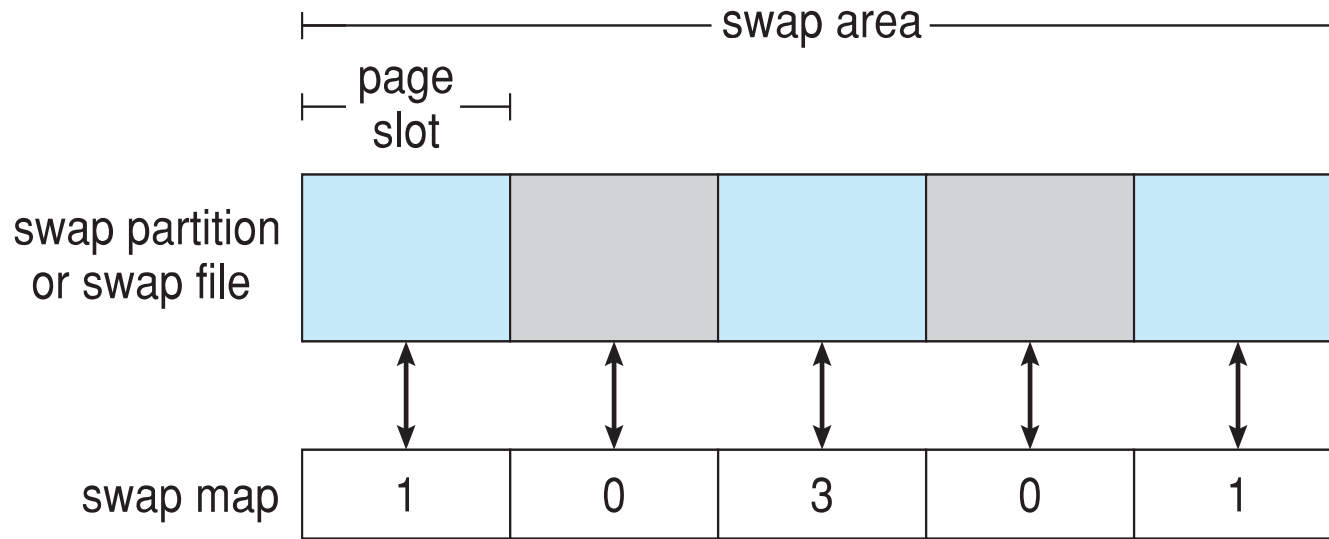
# Booting from a Disk in Windows



# Swap-Space Management

- Swap-space — Virtual memory uses disk space as an extension of main memory
  - ☞ Less common now due to memory capacity increases
- Swap-space can be carved out of the normal file system, or, more commonly, it can be in a separate disk partition (raw)
- 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
    - 📄 File data written to swap space until write to file system requested
    - 📄 Other dirty pages go to swap space due to no other home
    - 📄 Text segment pages thrown out and reread from the file system as needed
- What if a system runs out of swap space?
- Some systems allow multiple swap spaces

# Data Structures for Swapping on Linux Systems



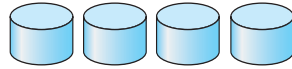
# RAID Structure

- RAID – redundant array of inexpensive disks
  - ☞ multiple disk drives provides reliability via **redundancy**
- Increases the **mean time to failure**
- **Mean time to repair** – exposure time when another failure could cause data loss
- **Mean time to data loss** based on above factors
- If mirrored disks fail independently, consider disk with 1300,000 mean time to failure and 10 hour mean time to repair
  - ☞ Mean time to data loss is  $100,000^2 / (2 * 10) = 500 * 10^6$  hours, or 57,000 years!
- Frequently combined with **NVRAM** to improve write performance
- Several improvements in disk-use techniques involve the use of multiple disks working cooperatively

# RAID (Cont.)

- Disk **striping** uses a group of disks as one storage unit
- RAID is arranged into six different levels
- RAID schemes improve performance and improve the reliability of the storage system by storing redundant data
  - ☞ **Mirroring** or **shadowing (RAID 1)** keeps duplicate of each disk
  - ☞ Striped mirrors (**RAID 1+0**) or mirrored stripes (**RAID 0+1**) provides high performance and high reliability
  - ☞ **Block interleaved parity (RAID 4, 5, 6)** uses much less redundancy
- RAID within a storage array can still fail if the array fails, so automatic **replication** of the data between arrays is common
- Frequently, a small number of **hot-spare** disks are left unallocated, automatically replacing a failed disk and having data rebuilt onto them

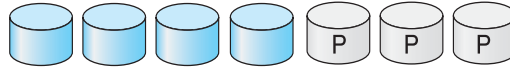
# RAID Levels



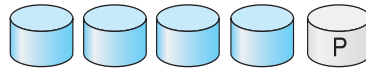
(a) RAID 0: non-redundant striping.



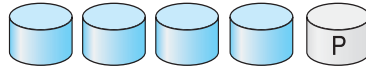
(b) RAID 1: mirrored disks.



(c) RAID 2: memory-style error-correcting codes.



(d) RAID 3: bit-interleaved parity.



(e) RAID 4: block-interleaved parity.



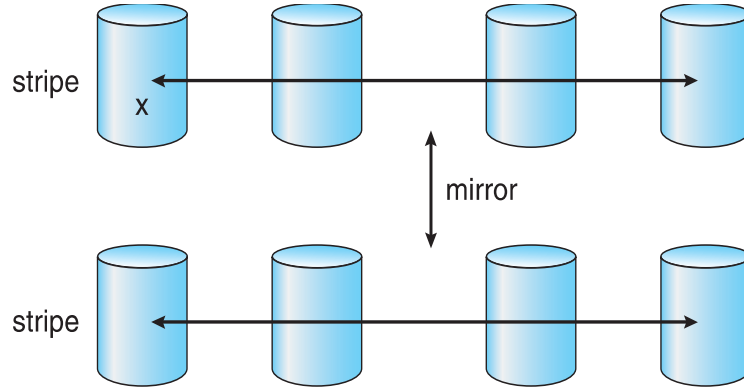
(f) RAID 5: block-interleaved distributed parity.



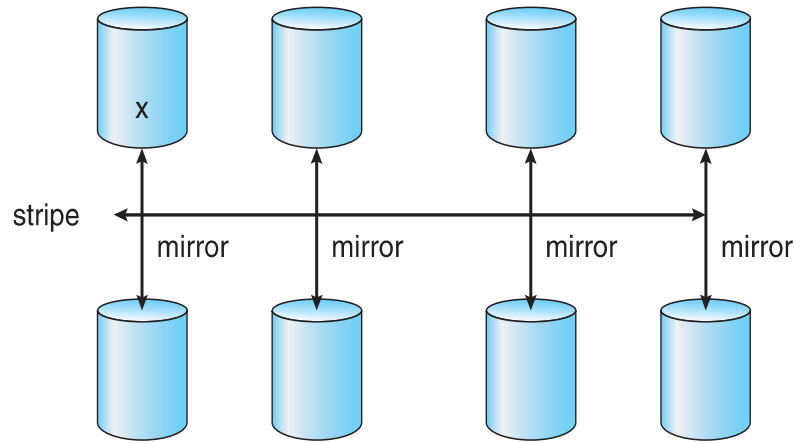
(g) RAID 6: P + Q redundancy.



# RAID (0 + 1) and (1 + 0)



a) RAID 0 + 1 with a single disk failure.



b) RAID 1 + 0 with a single disk failure.

# Other Features

- Regardless of where RAID implemented, other useful features can be added
- **Snapshot** is a view of file system before a set of changes take place (i.e. at a point in time)
  - ☞ More in Ch 12
- Replication is automatic duplication of writes between separate sites
  - ☞ For redundancy and disaster recovery
  - ☞ Can be synchronous or asynchronous
- Hot spare disk is unused, automatically used by RAID production if a disk fails to replace the failed disk and rebuild the RAID set if possible
  - ☞ Decreases mean time to repair

# Extensions

- RAID alone does not prevent or detect data corruption or other errors, just disk failures
- Solaris ZFS adds **checksums** of all data and metadata
- Checksums kept with pointer to object, to detect if object is the right one and whether it changed
- Can detect and correct data and metadata corruption
- ZFS also removes volumes, partitions
  - ☞ Disks allocated in **pools**
  - ☞ Filesystems with a pool share that pool, use and release space like `malloc()` and `free()` memory allocate / release calls

# UNIT V: Deadlocks

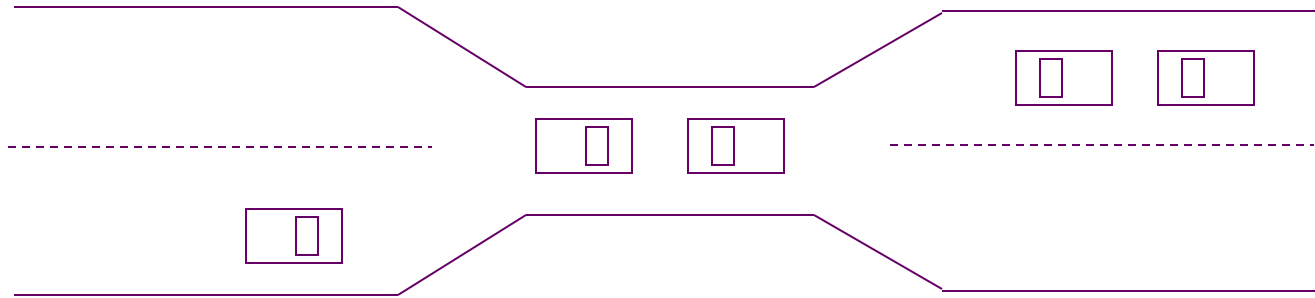
- System Model
- Deadlock Characterization
- Methods for Handling Deadlocks
- Deadlock Prevention
- Deadlock Avoidance
- Deadlock Detection
- Recovery from Deadlock
- Combined Approach to Deadlock Handling

# The Deadlock Problem

- A set of blocked processes each holding a resource and waiting to acquire a resource held by another process in the set.
- Example
  - ☞ System has 2 tape drives.
  - ☞  $P_1$  and  $P_2$  each hold one tape drive and each needs another one.
- Example
  - ☞ semaphores  $A$  and  $B$ , initialized to 1

$P_0$	$P_1$
<i>wait (A);</i>	<i>wait(B)</i>
<i>wait (B);</i>	<i>wait(A)</i>

# Bridge Crossing Example



- Traffic only in one direction.
- Each section of a bridge can be viewed as a resource.
- If a deadlock occurs, it can be resolved if one car backs up (preempt resources and rollback).
- Several cars may have to be backed up if a deadlock occurs.
- Starvation is possible.

# System Model

- Resource types  $R_1, R_2, \dots, R_m$   
*CPU cycles, memory space, I/O devices*
- Each resource type  $R_i$  has  $W_i$  instances.
- Each process utilizes a resource as follows:
  - ☞ request
  - ☞ use
  - ☞ release

# Deadlock Characterization

Deadlock can arise if four conditions hold simultaneously.

- **Mutual exclusion:** only one process at a time can use a resource.
- **Hold and wait:** a process holding at least one resource is waiting to acquire additional resources held by other processes.
- **No preemption:** a resource can be released only voluntarily by the process holding it, after that process has completed its task.
- **Circular wait:** there exists a set  $\{P_0, P_1, \dots, P_{n-1}\}$  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  $P_2$ , ...,  $P_{n-1}$  is waiting for a resource that is held by  $P_0$ .



# Resource-Allocation Graph

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

- $V$  is partitioned into two types:

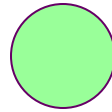
- ☞  $P = \{P_1, P_2, \dots, P_n\}$ , the set consisting of all the processes in the system.

- ☞  $R = \{R_1, R_2, \dots, R_m\}$ , the set consisting of all resource types in the system.

- request edge – directed edge  $P_i \rightarrow R_j$
- assignment edge – directed edge  $R_j \rightarrow P_i$

# Resource-Allocation Graph (Cont.)

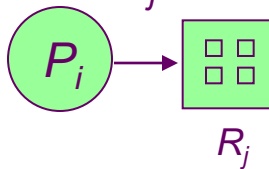
- Process



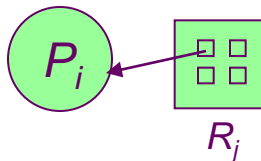
- Resource Type with 4 instances



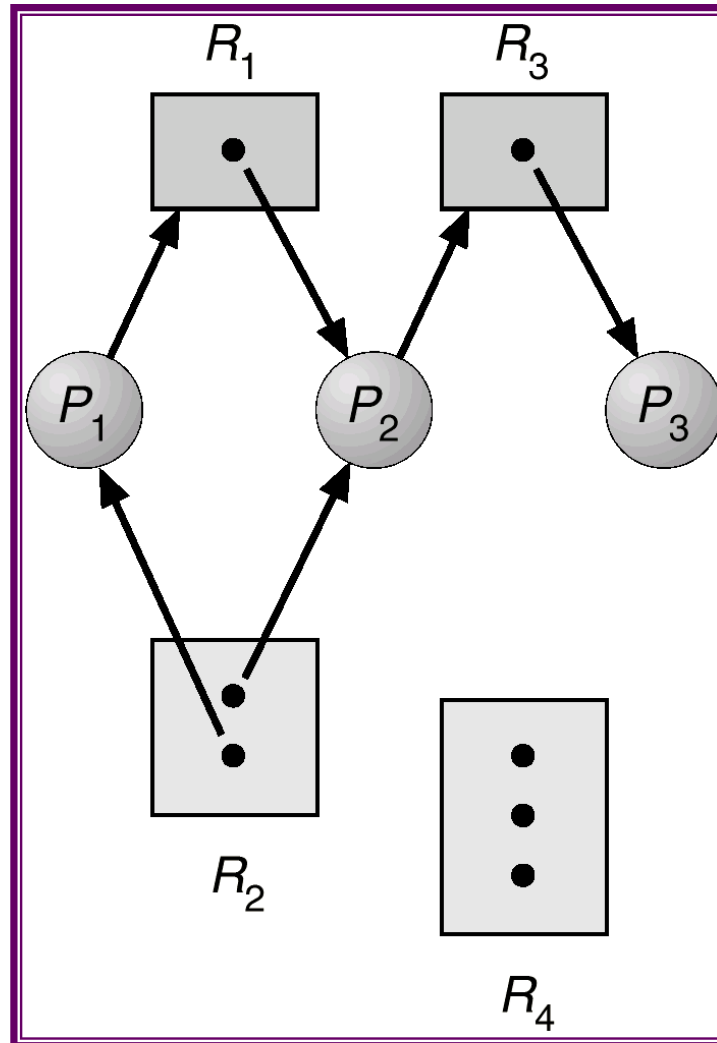
- $P_i$  requests instance of  $R_j$



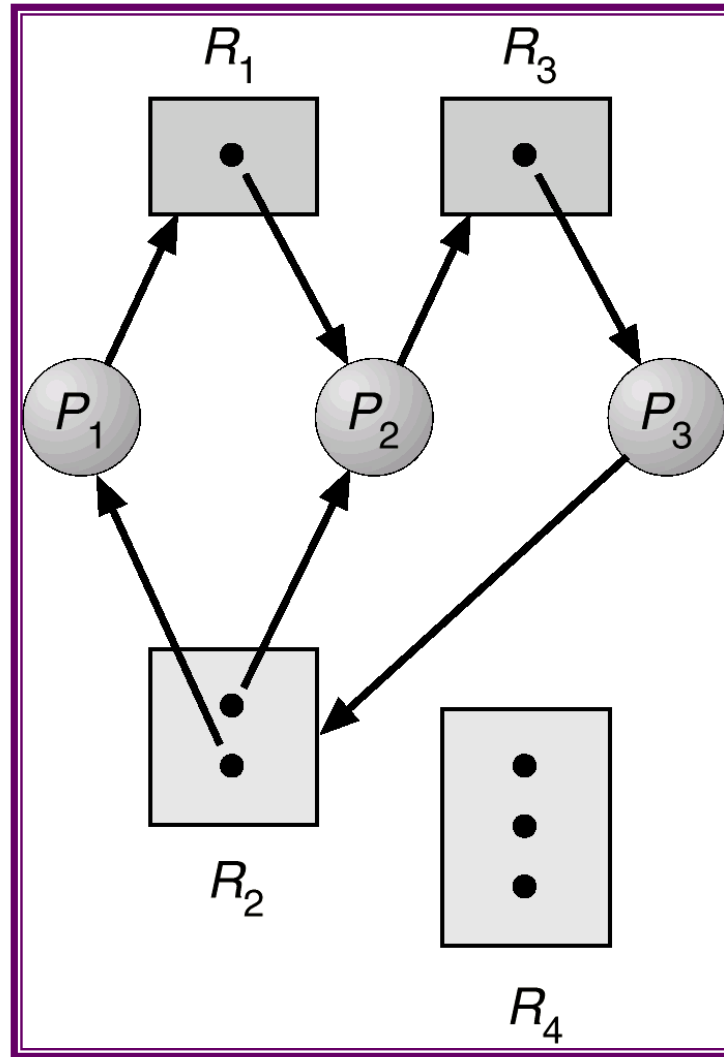
- $P_i$  is holding an instance of  $R_j$



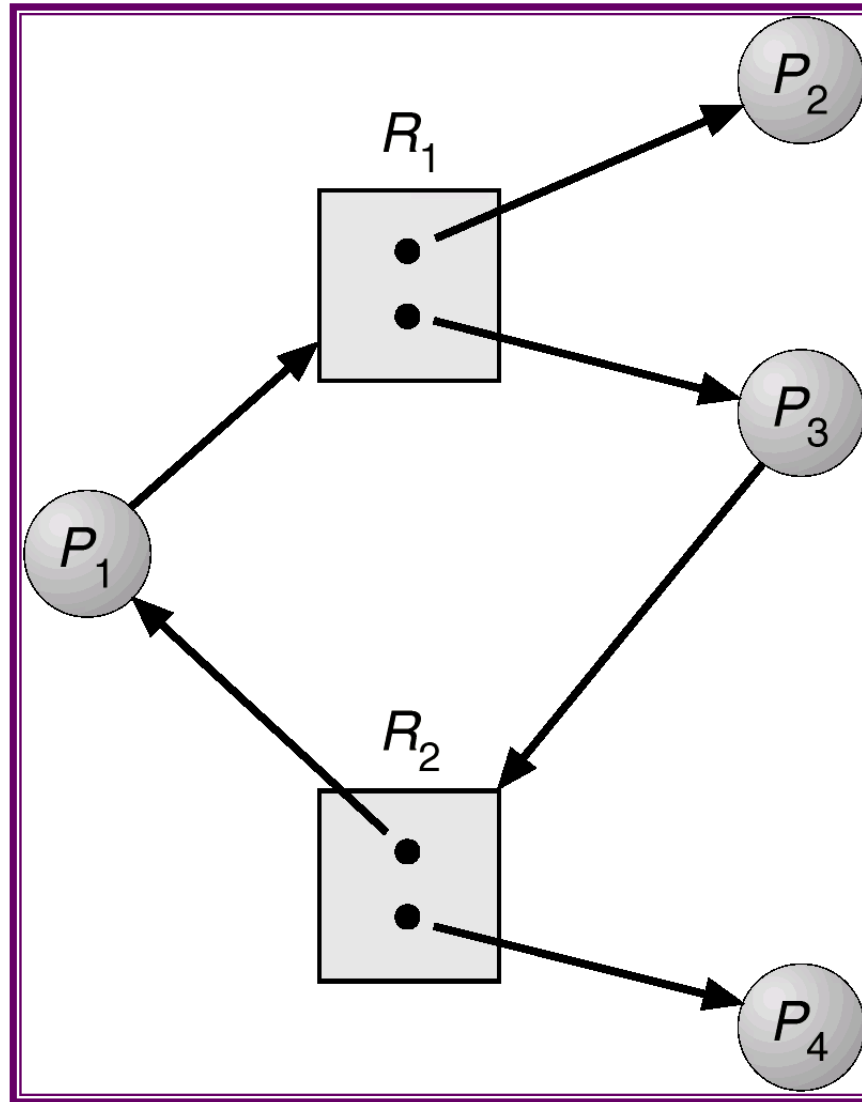
# Example of a Resource Allocation Graph



# Resource Allocation Graph With A Deadlock



# Resource Allocation Graph With A Cycle But No Deadlock



# Basic Facts

- If graph contains no cycles  $\Rightarrow$  no deadlock.
- 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

- Ensure that the system will *never* enter a deadlock state.
- Allow the system to enter a deadlock state and then recover.
- 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.

- **Mutual Exclusion** – not required for sharable resources; must hold for nonsharable resources.
- **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.



# Deadlock Prevention (Cont.)

## ■ 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.

## ■ 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.

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

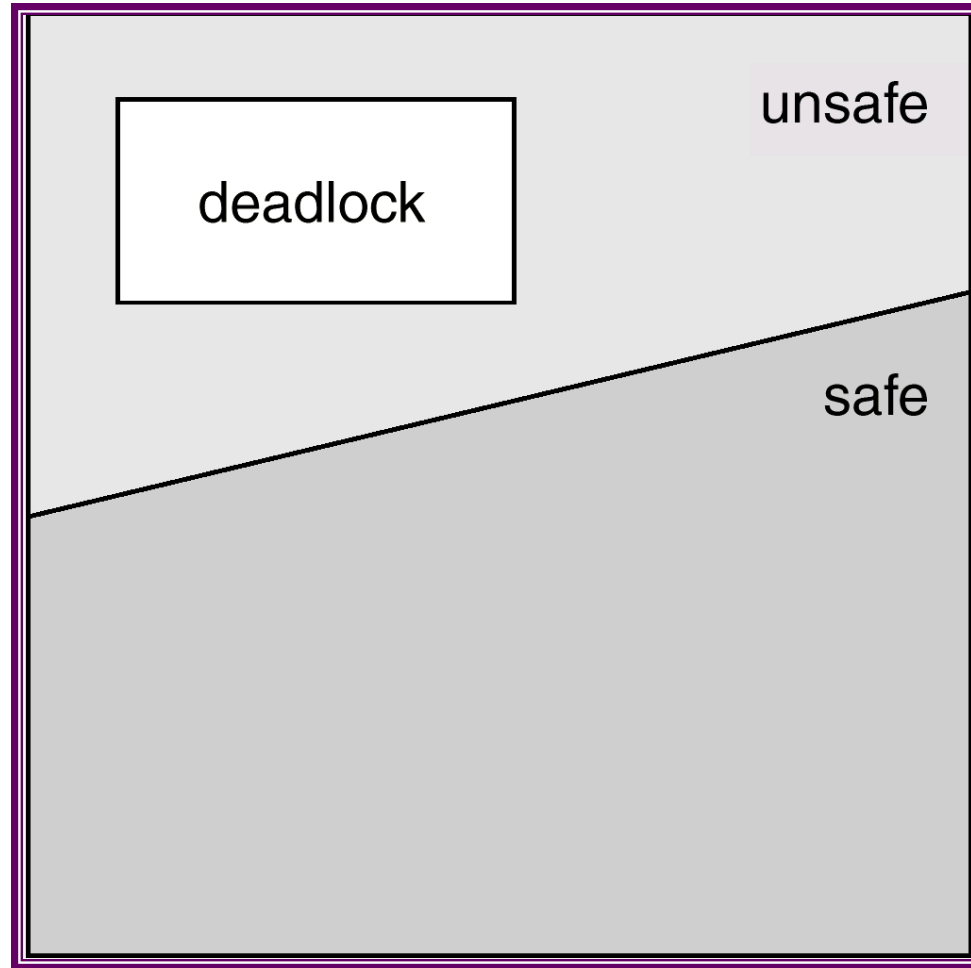
# Safe State

- When a process requests an available resource, system must decide if immediate allocation leaves the system in a *safe state*.
- System is in safe state if there exists a safe sequence of all processes.
- Sequence  $\langle P_1, P_2, \dots, P_n \rangle$  is safe if 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  $j < i$ .
  - ☞ 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 so on.

# Basic Facts

- If a system is in safe state  $\Rightarrow$  no deadlocks.
- If a system is in unsafe state  $\Rightarrow$  possibility of deadlock.
- Avoidance  $\Rightarrow$  ensure that a system will never enter an unsafe state.

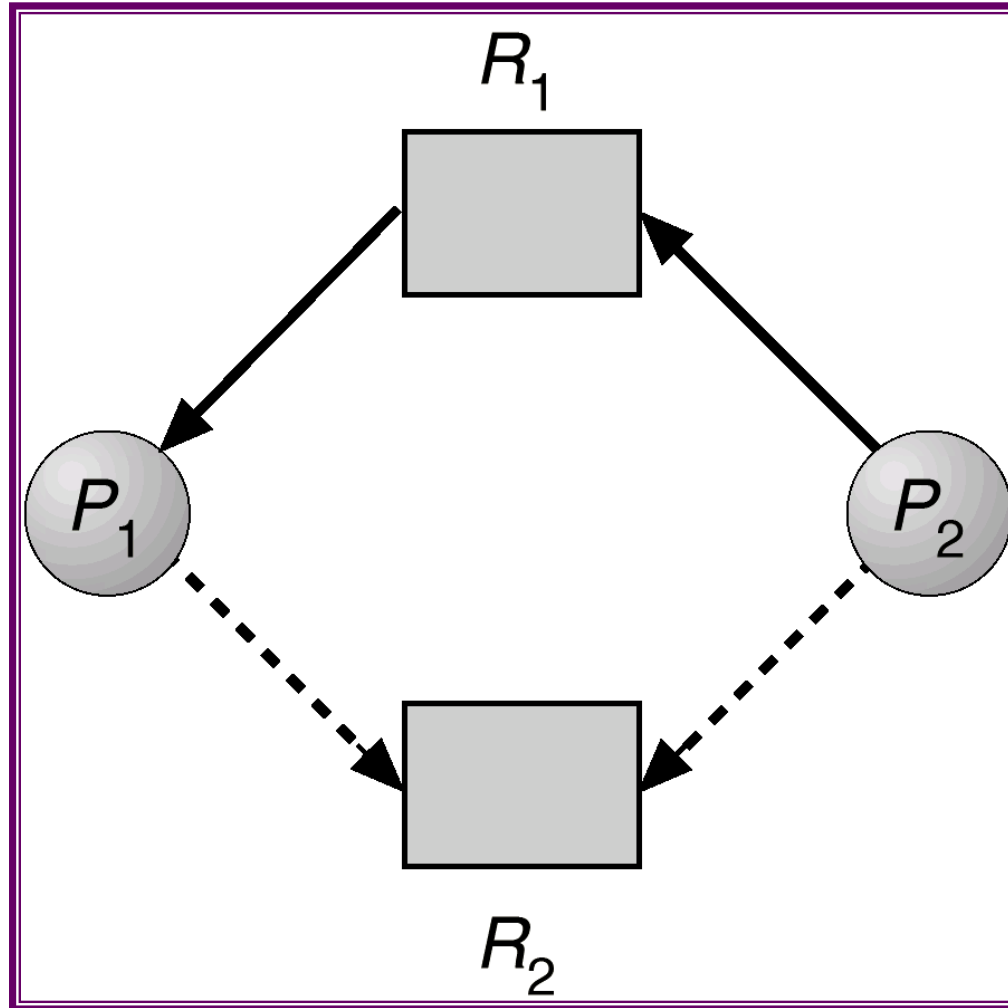
# Safe, Unsafe , Deadlock State



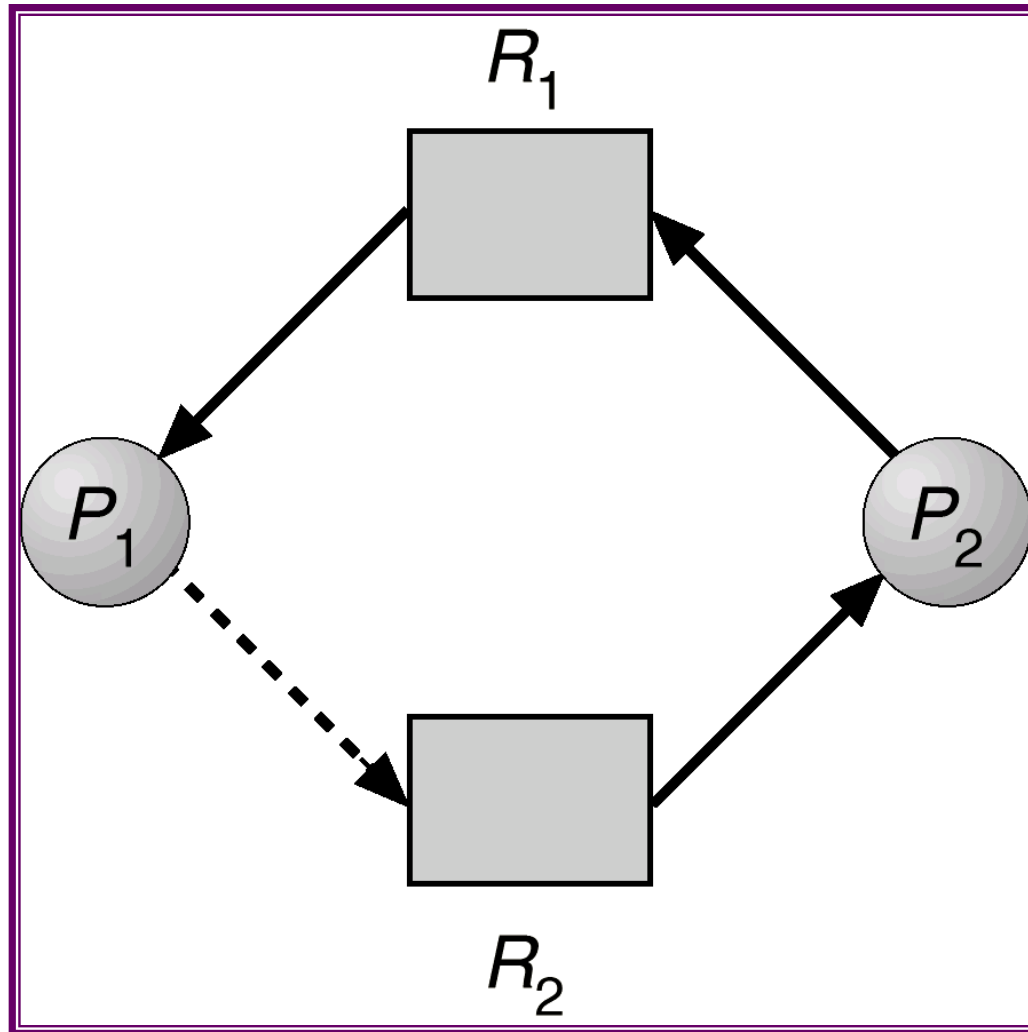
# Resource-Allocation Graph Algorithm

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

# Resource-Allocation Graph For Deadlock Avoidance



# Unsafe State In Resource-Allocation Graph





# Banker's Algorithm

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

# Data Structures for the Banker's Algorithm

Let  $n$  = number of processes, and  $m$  = number of resources types.

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

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

# Safety Algorithm

1. Let *Work* and *Finish* be vectors of length *m* and *n*, respectively. Initialize:  
 $Work = Available$   
 $Finish[i] = false$  for  $i = 1, 3, \dots, n$ .
2. Find and *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:

$$Available = Available - Request_i;$$

$$Allocation_i = Allocation_i + Request_i;$$

$$Need_i = Need_i - Request_i;$$

- 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

- 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>
	$A\ B\ C$	$A\ B\ C$	$A\ B\ C$
$P_0$	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	

# Example (Cont.)

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

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

- 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.

# Example $P_1$ Request (1,0,2) (Cont.)

- Check that Request  $\leq$  Available (that is,  $(1,0,2) \leq (3,3,2) \Rightarrow$  true.

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

- Executing safety algorithm shows that sequence  $\langle P_1, P_3, P_4, P_0, P_2 \rangle$  satisfies safety requirement.
- Can request for (3,3,0) by  $P_4$  be granted?
- Can request for (0,2,0) by  $P_0$  be granted?

# Deadlock Detection

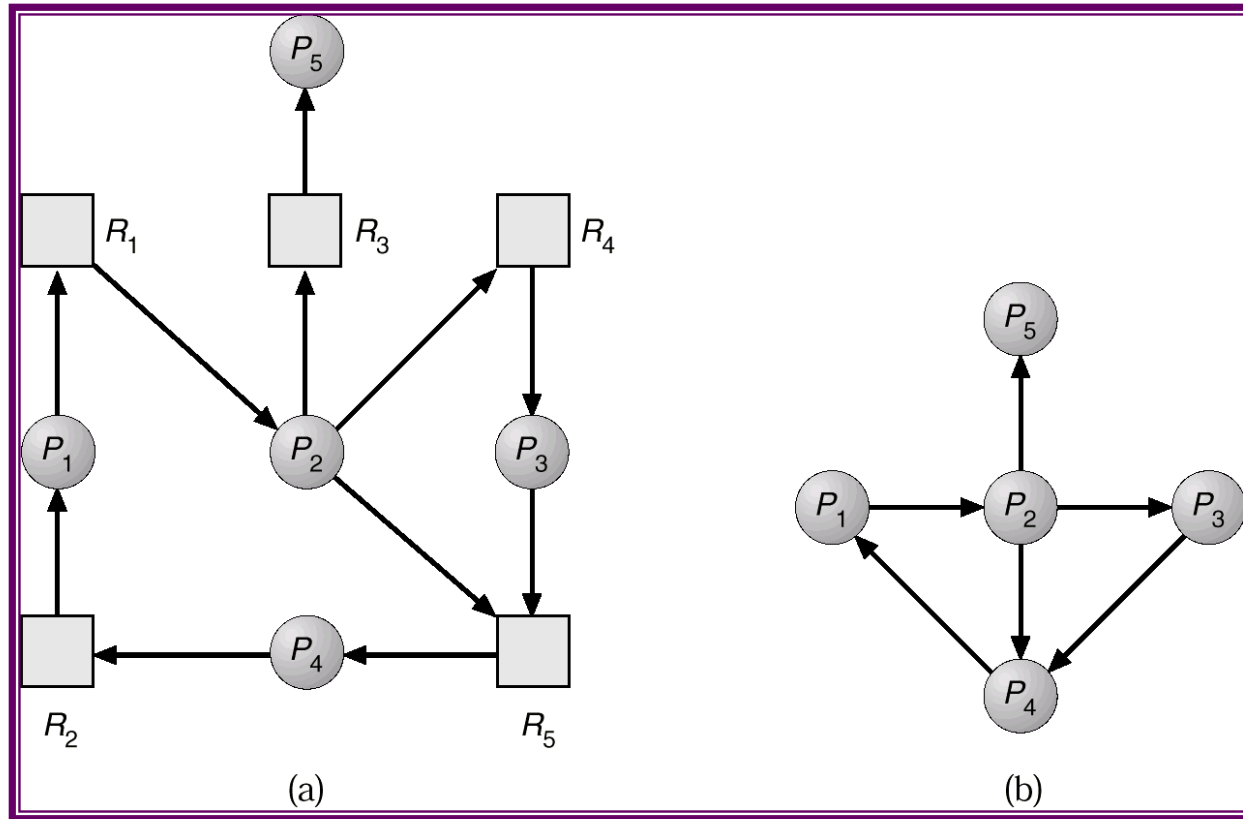
- Allow system to enter deadlock state
- Detection algorithm
- Recovery scheme



# Single Instance of Each Resource Type

- Maintain *wait-for* graph
  - ☞ Nodes are processes.
  - ☞  $P_i \rightarrow P_j$  if  $P_i$  is waiting for  $P_j$ .
- Periodically invoke an algorithm that searches for a cycle in the graph.
- An algorithm to detect a cycle in a graph requires an order of  $n^2$  operations, where  $n$  is the number of vertices in the graph.

# Resource-Allocation Graph and Wait-for Graph



Resource-Allocation Graph

Corresponding wait-for graph

# Several Instances of a Resource Type

- *Available:* A vector of length  $m$  indicates the number of available resources of each type.
- *Allocation:* An  $n \times m$  matrix defines the number of resources of each type currently allocated to each process.
- *Request:* An  $n \times m$  matrix indicates the current request of each process. If  $Request[i_j] = k$ , then process  $P_i$  is requesting  $k$  more instances of resource type  $R_j$ .

# Detection Algorithm

1. 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.

# Detection Algorithm (Cont.)

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.

Algorithm requires an order of  $O(m \times n^2)$  operations to detect whether the system is in deadlocked state.

# Example of Detection Algorithm

- Five processes  $P_0$  through  $P_4$ ; three resource types A (7 instances), B (2 instances), and C (6 instances).
- Snapshot at time  $T_0$ :

	<u>Allocation</u>			<u>Request</u>			<u>Available</u>		
	A	B	C	A	B	C	A	B	C
$P_0$	0	1	0	0	0	0	0	0	0
$P_1$	2	0	0	2	0	2			
$P_2$	3	0	3	0	0	0			
$P_3$	2	1	1	1	0	0			
$P_4$	0	0	2	0	0	2			

- Sequence  $\langle P_0, P_2, P_3, P_1, P_4 \rangle$  will result in  $Finish[i] = \text{true}$  for all  $i$ .

# Example (Cont.)

- $P_2$  requests an additional instance of type C.

	<u>Request</u>		
	A	B	C
$P_0$	0	0	0
$P_1$	2	0	1
$P_2$	0	0	1
$P_3$	1	0	0
$P_4$	0	0	2

- State of system?

- ☞ Can reclaim resources held by process  $P_0$ , but insufficient resources to fulfill other processes; requests.
- ☞ Deadlock exists, consisting of processes  $P_1$ ,  $P_2$ ,  $P_3$ , and  $P_4$ .

# Detection-Algorithm Usage

- When, and how often, to invoke depends on:
  - ☞ How often a deadlock is likely to occur?
  - ☞ How many processes will need to be rolled back?
    - 📄 one for each disjoint cycle
- If detection algorithm is invoked arbitrarily, there may be many cycles in the resource graph and so we would not be able to tell which of the many deadlocked processes “caused” the deadlock.



# Recovery from Deadlock: Process Termination

- Abort all deadlocked processes.
- Abort one process at a time until the deadlock cycle is eliminated.
- 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?

# Recovery from Deadlock: Resource Preemption

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

# Combined Approach to Deadlock Handling

- Combine the three basic approaches

- ☞ prevention

- ☞ avoidance

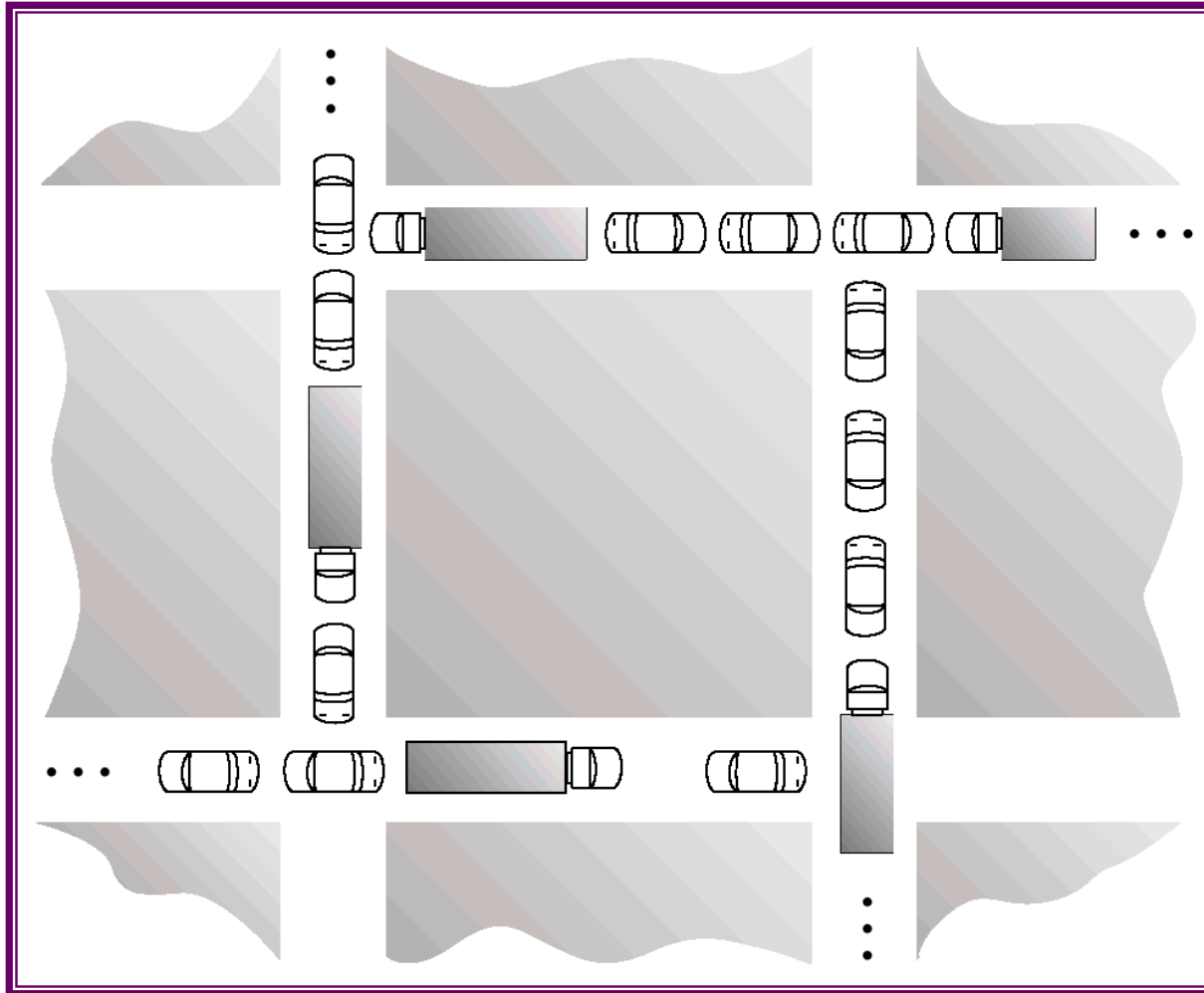
- ☞ detection

allowing the use of the optimal approach for each of resources in the system.

- Partition resources into hierarchically ordered classes.

- Use most appropriate technique for handling deadlocks within each class.

# Traffic Deadlock for Exercise 8.4



# UNIT V: Protection

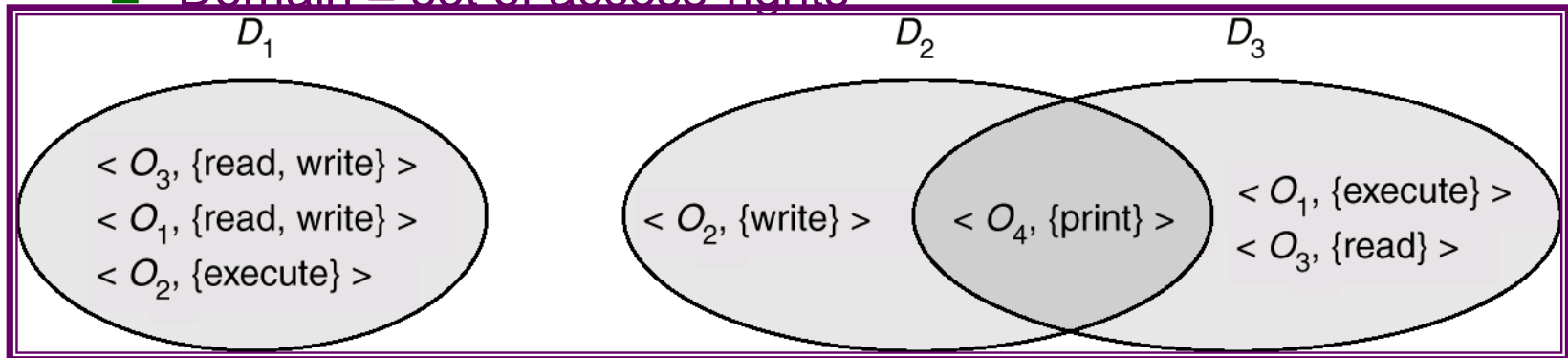
- Goals of Protection
- Domain of Protection
- Access Matrix
- Implementation of Access Matrix
- Revocation of Access Rights
- Capability-Based Systems
- Language-Based Protection

# Protection

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

# Domain Structure

- 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.
- Domain = set of access-rights



# Domain Implementation (UNIX)

- System consists of 2 domains:

- ☞ User
- ☞ Supervisor

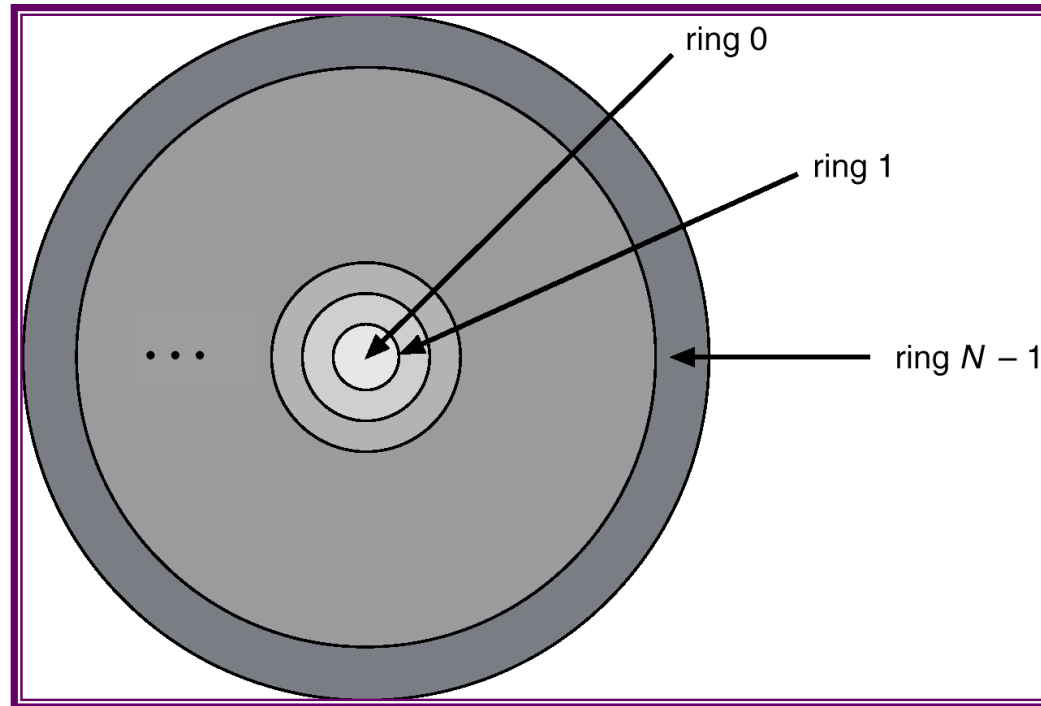
- UNIX

- ☞ Domain = user-id
- ☞ Domain switch accomplished via file system.
  - 📄 Each file has associated with it a domain bit (setuid bit).
  - 📄 When file is executed and setuid = on, then user-id is set to owner of the file being executed. When execution completes user-id is reset.



# Domain Implementation (Multics)

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



Multics Rings

# Access Matrix

- View protection as a matrix (*access matrix*)
- Rows represent domains
- Columns represent objects
- $Access(i, j)$  is the set of operations that a process executing in Domain<sub>*i*</sub> can invoke on Object<sub>*j*</sub>

# Access Matrix

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

Figure A

# Use of Access Matrix

- If a process in Domain  $D_i$  tries to do “op” on object  $O_j$ , then “op” must be in the access matrix.
- Can be expanded to dynamic protection.
  - ☞ Operations to add, delete access rights.
  - ☞ Special access rights:
    - 📄 *owner of  $O_i$*
    - 📄 *copy op from  $O_i$  to  $O_j$*
    - 📄 *control –  $D_i$  can modify  $D_j$  access rights*
    - 📄 *transfer – switch from domain  $D_i$  to  $D_j$*

# Use of Access Matrix (Cont.)

- Access matrix design separates mechanism from policy.

- ☞ Mechanism

- 📄 Operating system provides access-matrix + rules.
    - 📄 If ensures that the matrix is only manipulated by authorized agents and that rules are strictly enforced.

- ☞ Policy

- 📄 User dictates policy.
    - 📄 Who can access what object and in what mode.

# Implementation of Access Matrix

- Each column = Access-control list for one object  
Defines who can perform what operation.

Domain 1 = Read, Write

Domain 2 = Read

Domain 3 = Read

⋮

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

Object 1 – Read

Object 4 – Read, Write, Execute

Object 5 – Read, Write, Delete, Copy

# Access Matrix of Figure A With Domains as Objects

domain \ object	$F_1$	$F_2$	$F_3$	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			

**Figure B**

# 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		
$D_2$		owner read* write*	read* owner write*
$D_3$		write	write

(b)

# Modified Access Matrix of Figure B

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 control
$D_3$		read	execute					
$D_4$	write		write		switch			

# Revocation of Access Rights

- *Access List* – Delete access rights from access list.
  - ☞ Simple
  - ☞ Immediate
- *Capability List* – Scheme required to locate capability in the system before capability can be revoked.
  - ☞ Reacquisition
  - ☞ Back-pointers
  - ☞ Indirection
  - ☞ Keys

# Capability-Based Systems

## ■ Hydra

- ☞ Fixed set of access rights known to and interpreted by the system.
- ☞ Interpretation of user-defined rights performed solely by user's program; system provides access protection for use of these rights.

## ■ Cambridge CAP System

- ☞ Data capability - provides standard read, write, execute of individual storage segments associated with object.
- ☞ Software capability - interpretation left to the subsystem, through its protected procedures.

# Language-Based Protection

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

# Protection in Java 2

- Protection is handled by the Java Virtual Machine (JVM)
- A class is assigned a protection domain when it is loaded by the JVM.
- The protection domain indicates what operations the class can (and cannot) perform.
- If a library method is invoked that performs a privileged operation, the stack is inspected to ensure the operation can be performed by the library.

# Stack Inspection

protection domain:	untrusted applet	URL loader	networking
socket permission:	none	*.lucent.com:80, connect	any
class:	gui: ... get(url); open(addr); ...	get(URL u): ... doPrivileged { open('proxy.lucent.com:80'); } <request u from proxy> ...	open(Addr a): ... checkPermission(a, connect); connect (a); ...