

# Processes

- Process Concept
- Process Scheduling
- Operations on Processes
- Cooperating Processes
- Inter process Communication
- Communication in Client-Server Systems

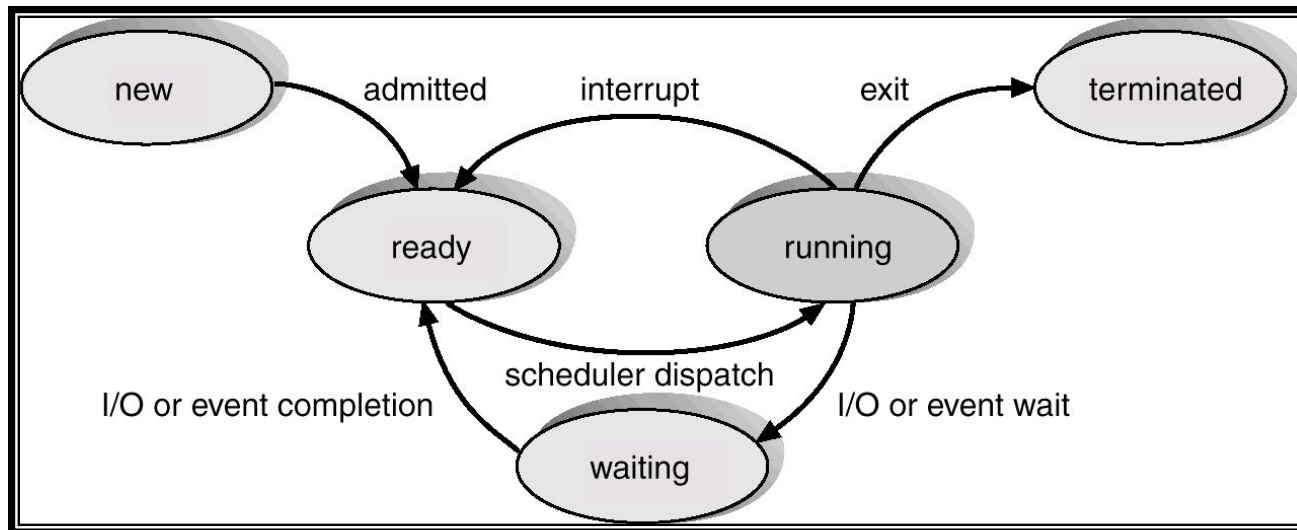
# Process Concept

- 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

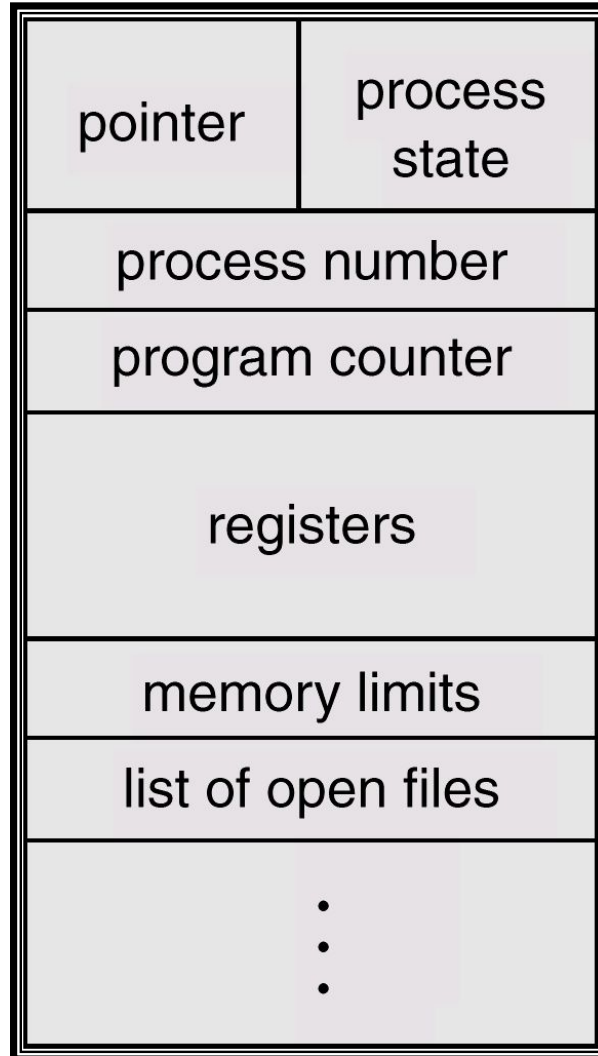


# Process Control Block (PCB)

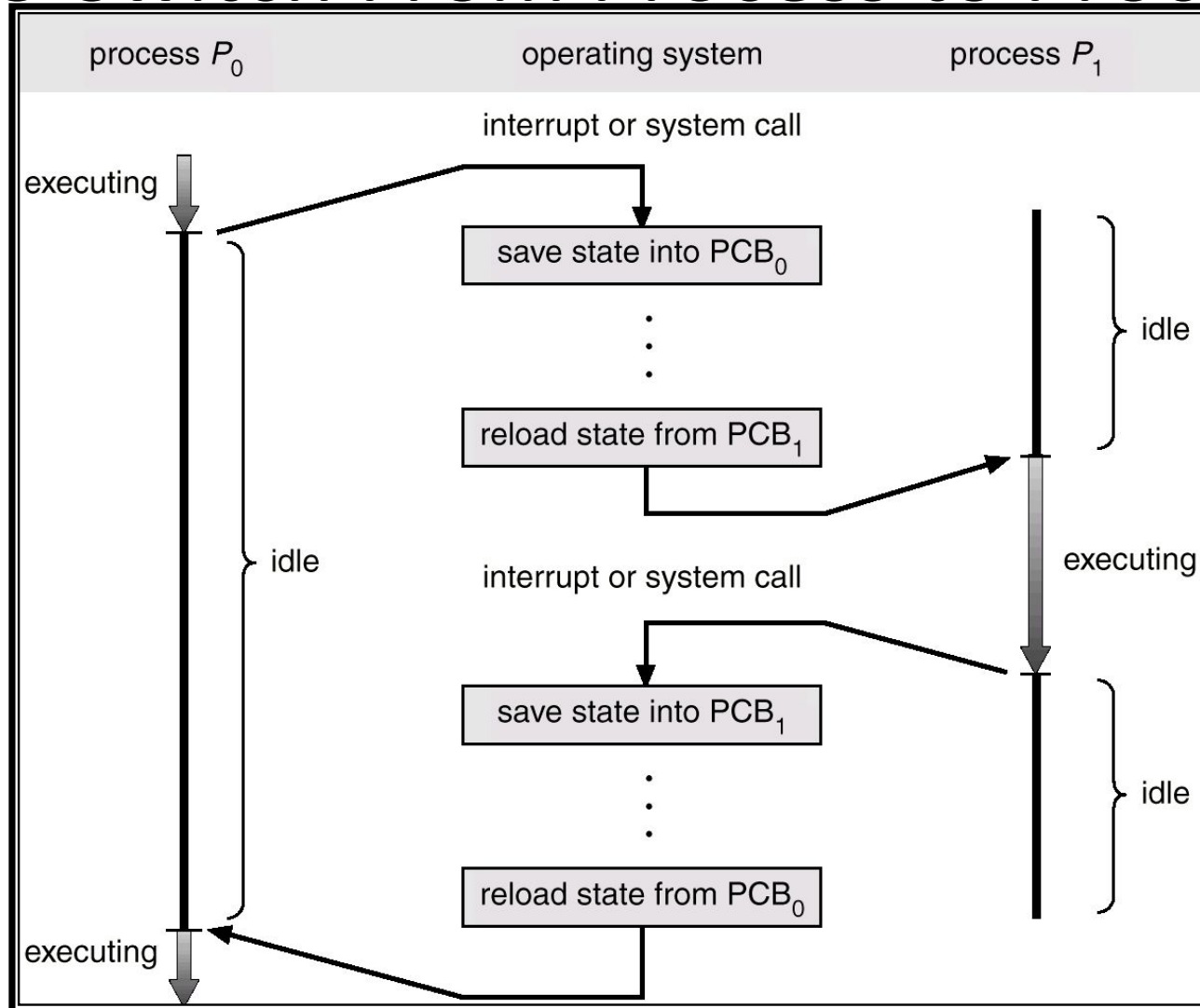
Information associated with each process.

- 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

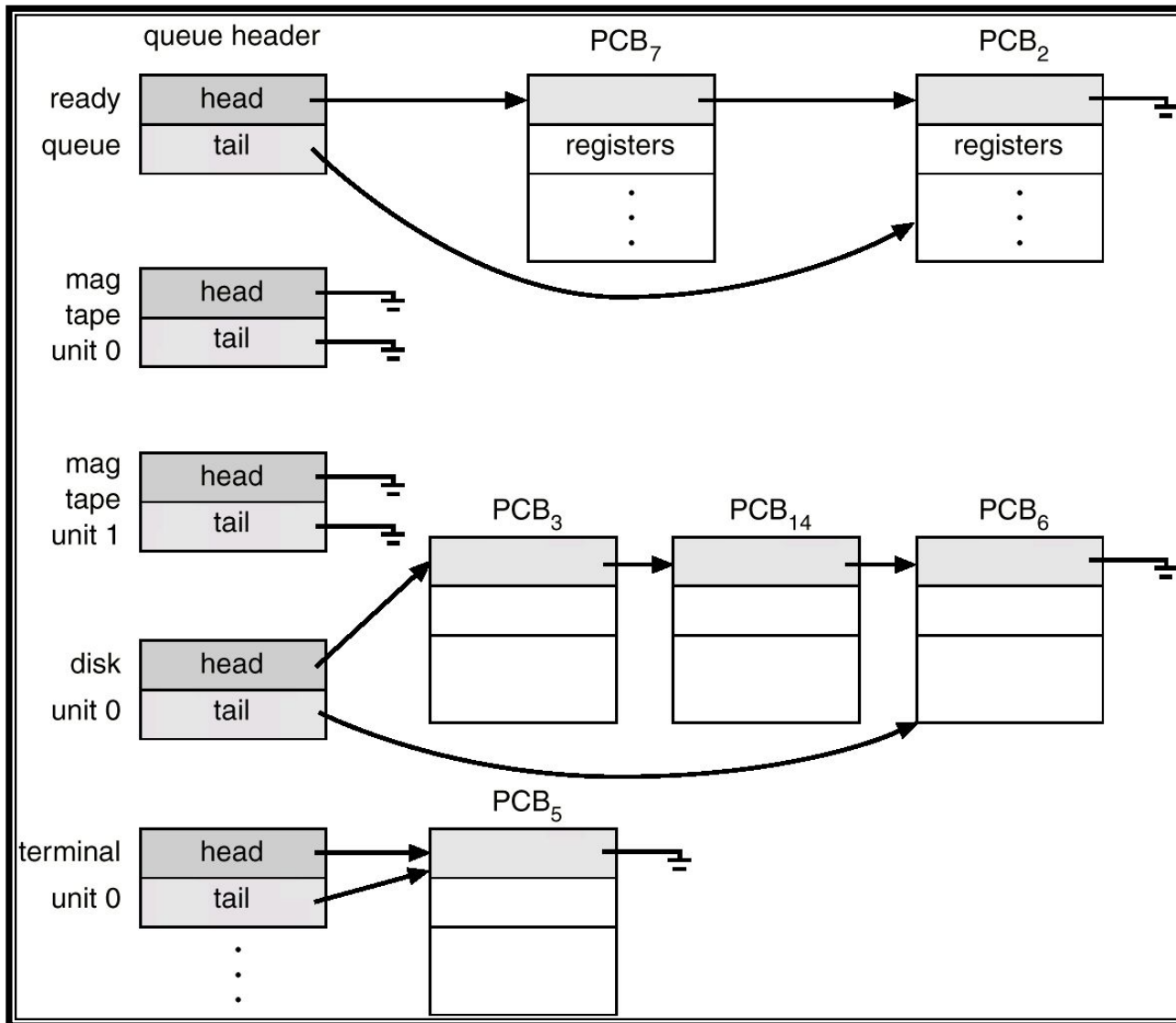


# Process Scheduling Queues

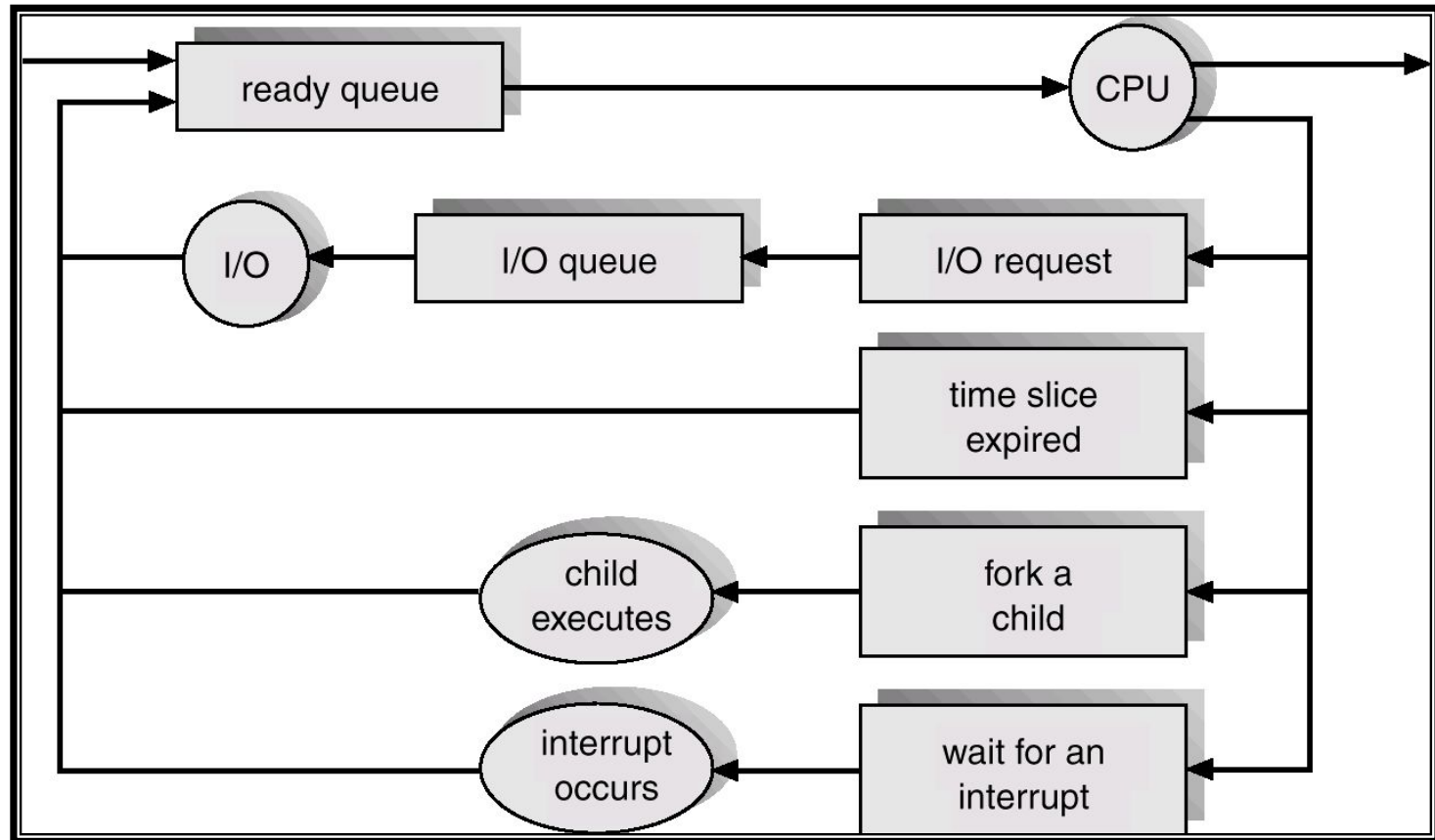
- Job queue – set of all processes in the system.
- Ready queue – set of all processes residing in main memory, ready and waiting to execute.
- Device queues – set of processes waiting for an I/O device.
- Process migration between the various queues.



# Ready Queue And Various I/O Device Queues



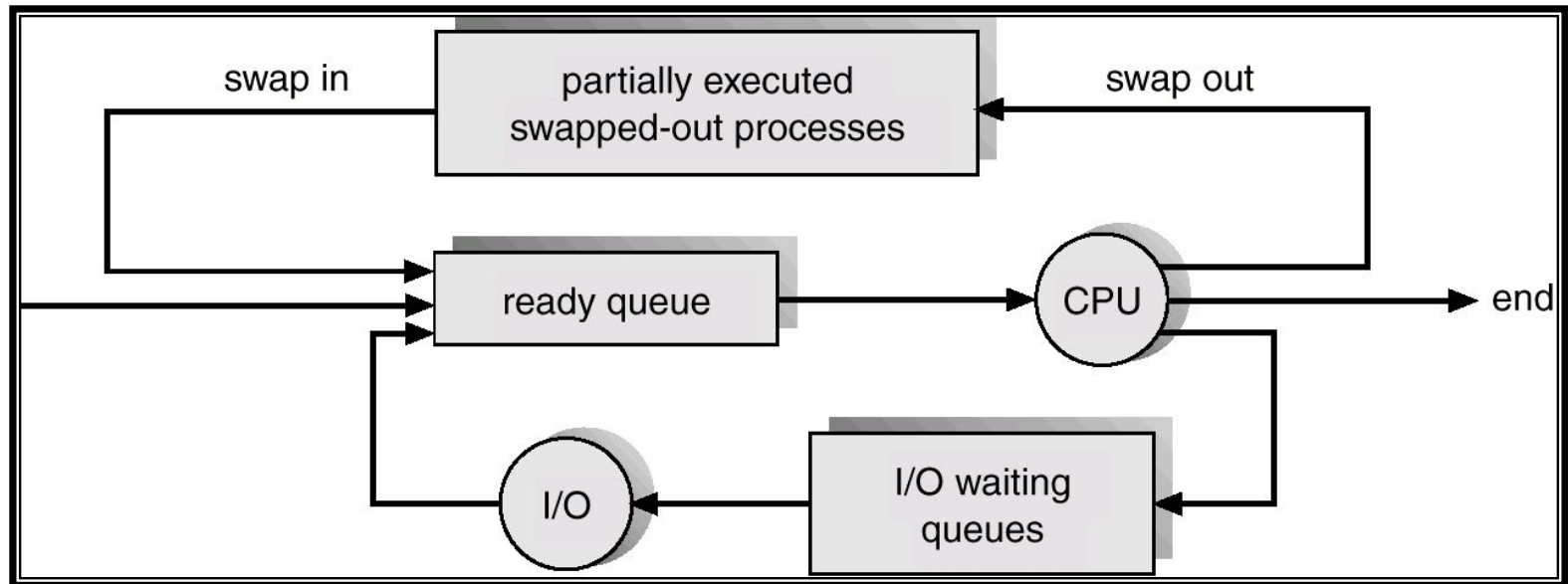
# Representation of Process Scheduling



# 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

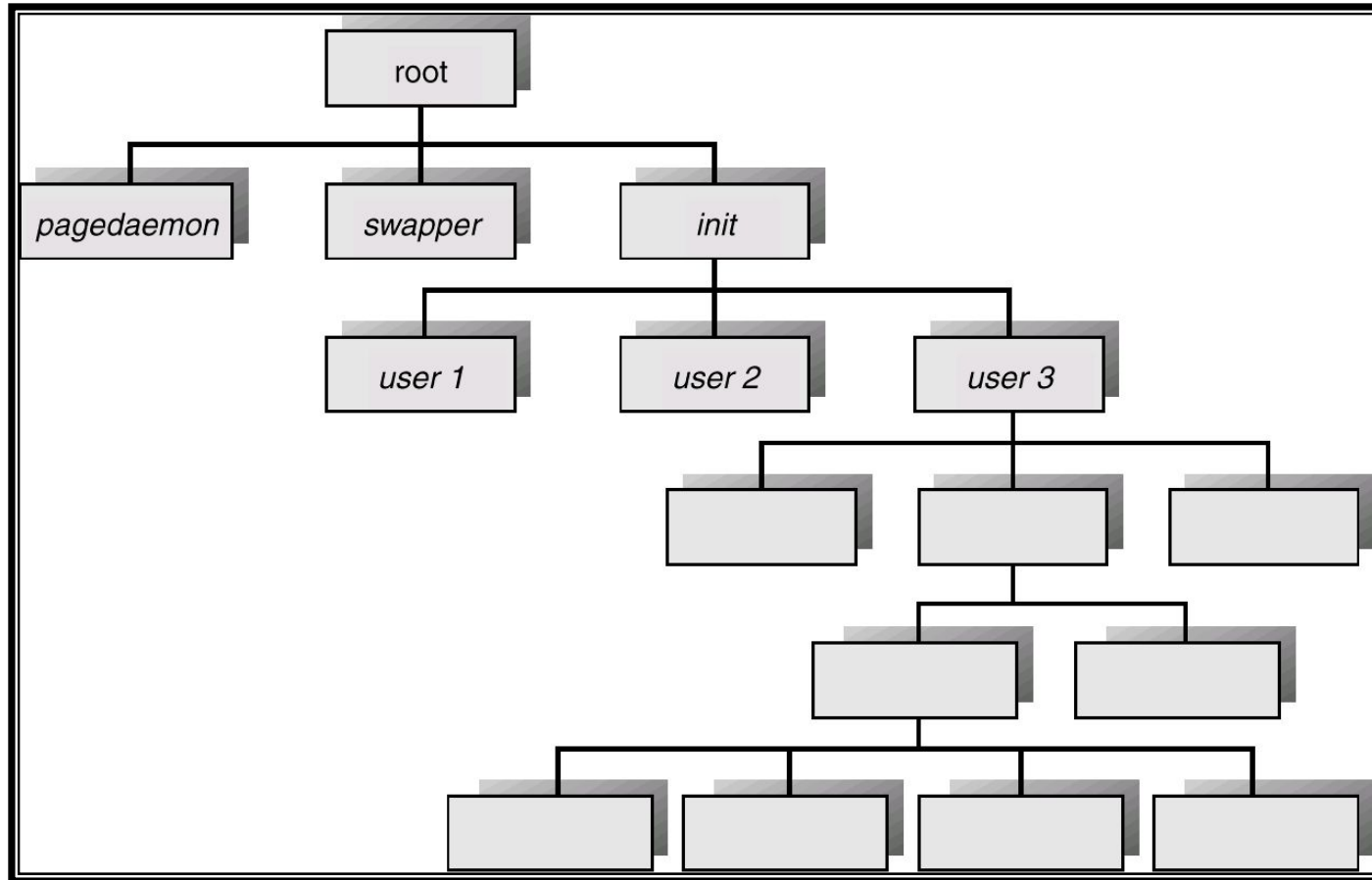
- Parent process create children processes, which, in turn create other processes, forming a tree of processes.
- Resource sharing
  - Parent and children share all resources.
  - Children share subset of parent's resources.
  - Parent and child share no resources.
- Execution
  - Parent and children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

# Buffering

- Queue of messages attached to the link; implemented in one of three ways.
  1. Zero capacity – 0 messages  
Sender must wait for receiver (rendezvous).
  2. Bounded capacity – finite length of  $n$  messages  
Sender must wait if link full.
  3. Unbounded capacity – infinite length  
Sender never waits.

# Client-Server Communication

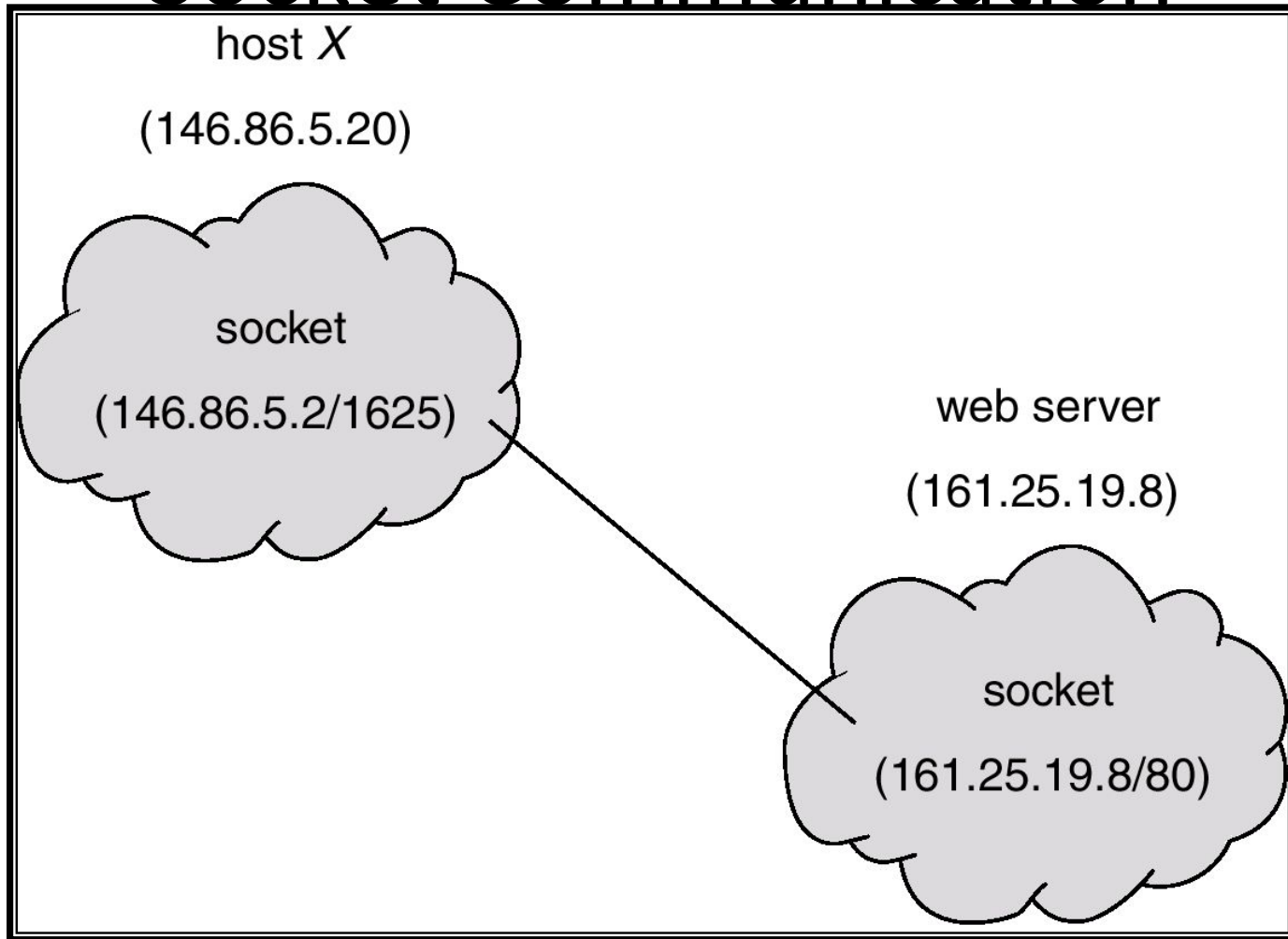
- Sockets
- Remote Procedure Calls
- Remote Method Invocation (Java)

# Sockets

- A socket is defined as an *endpoint for communication*.
- Concatenation of IP address and port
- The socket **161.25.19.8:1625** refers to port **1625** on host **161.25.19.8**
- Communication consists between a pair of sockets.



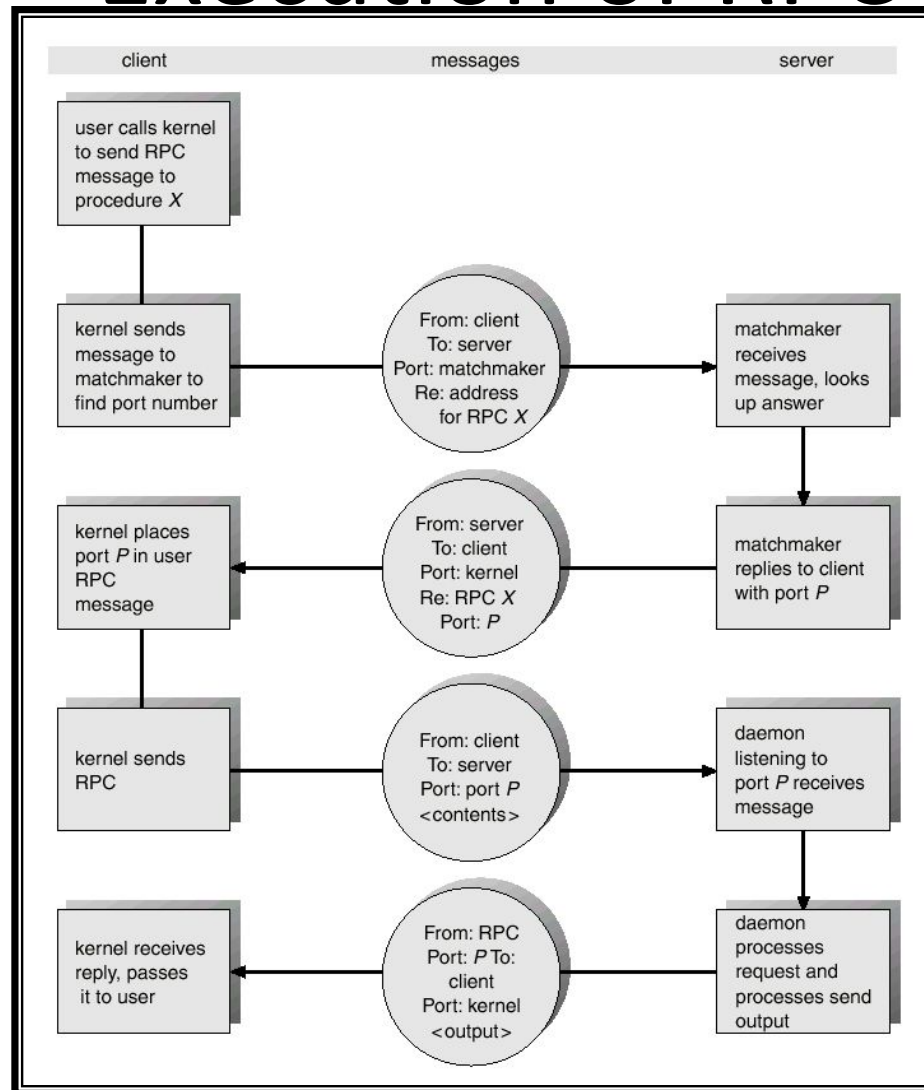
# Socket Communication



# Remote Procedure Calls

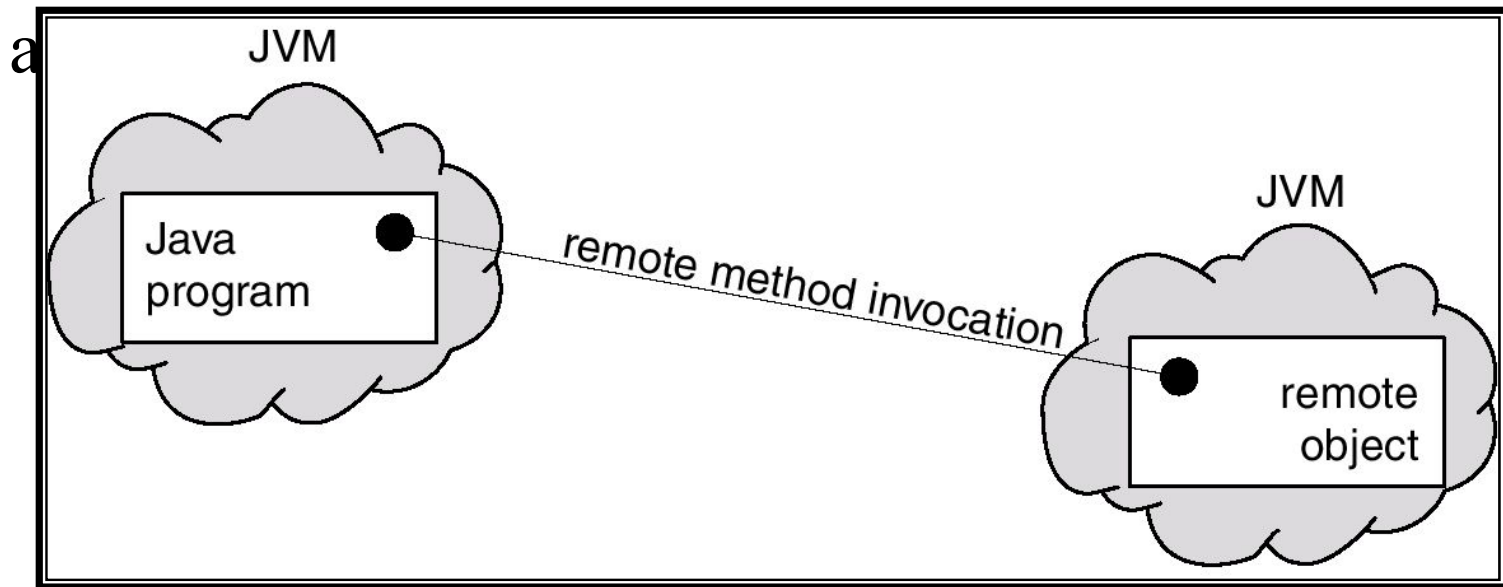
- Remote procedure call (RPC) abstracts procedure calls between processes on networked systems.
- **Stubs** – client-side proxy for the actual procedure on the server.
- The client-side stub locates the server and *marshalls* the parameters.
- The server-side stub receives this message, unpacks the marshalled parameters, and performs the procedure on the server.

# Execution of RPC

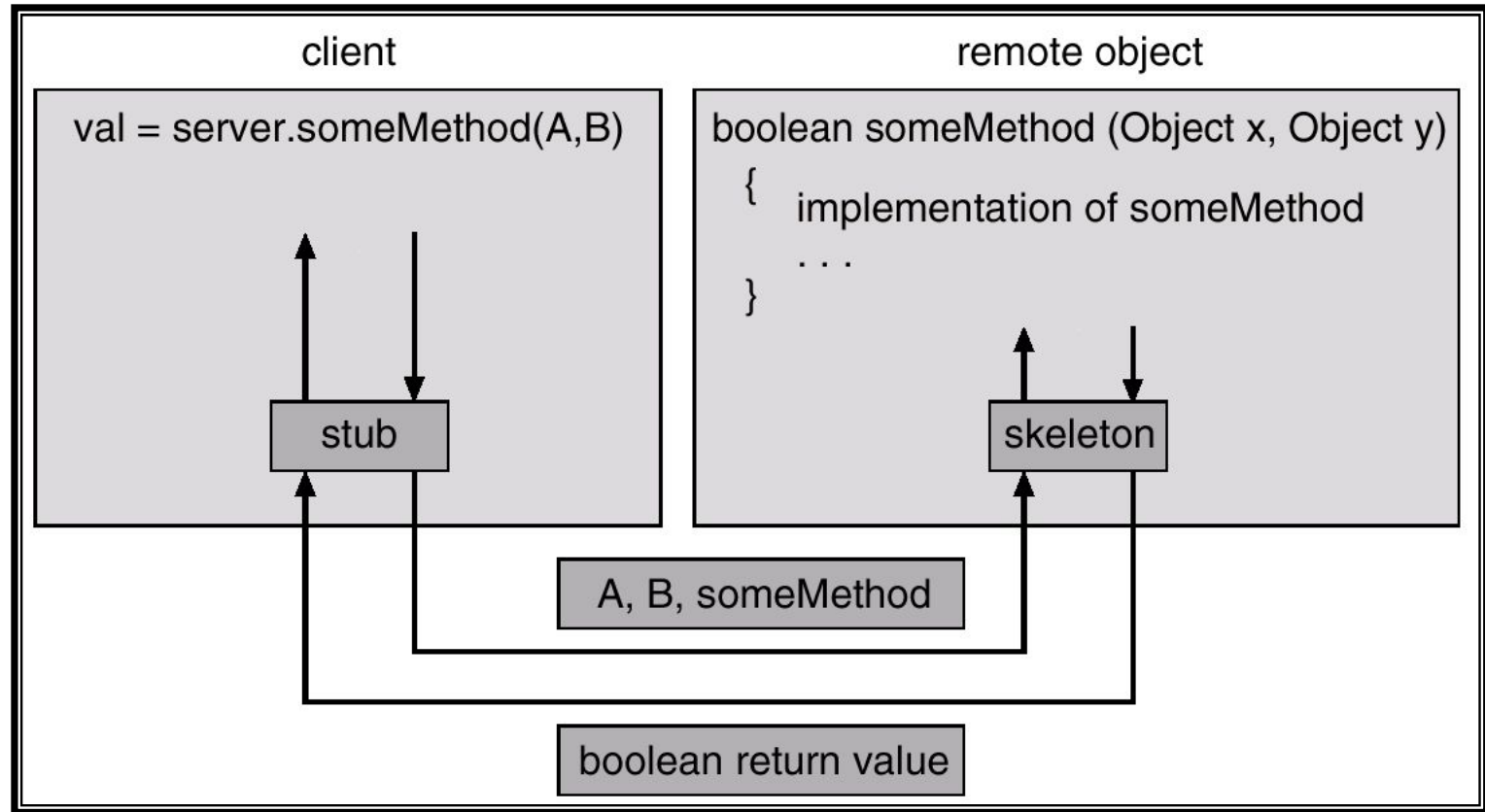


# Remote Method Invocation

- Remote Method Invocation (RMI) is a Java mechanism similar to RPCs.
- RMI allows a Java program on one machine to invoke



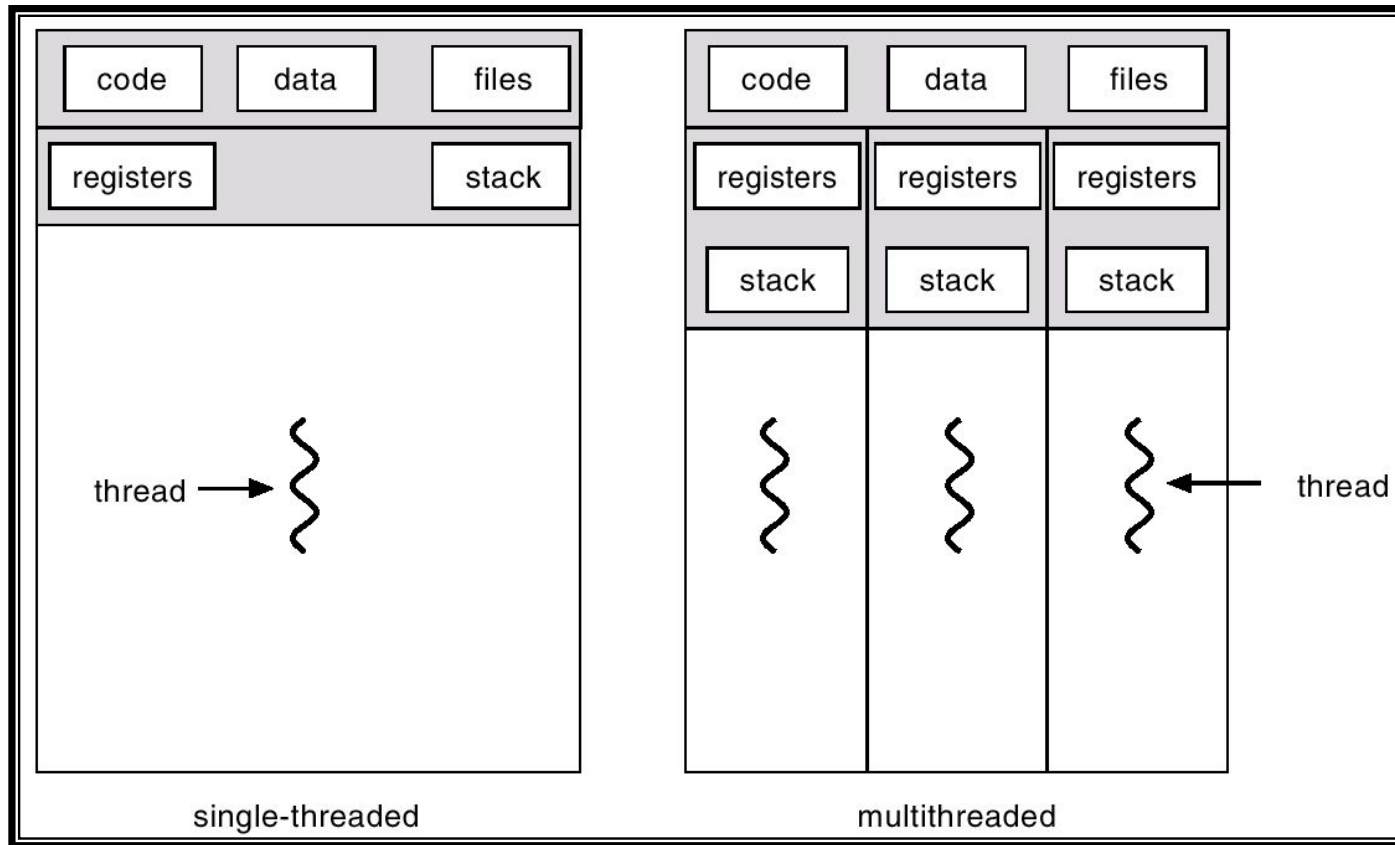
# Marshalling Parameters



# Chapter 5: Threads

- Overview
- Multithreading Models
- Threading Issues
- Pthreads
- Solaris 2 Threads
- Windows 2000 Threads
- Linux Threads
- Java Threads

# 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

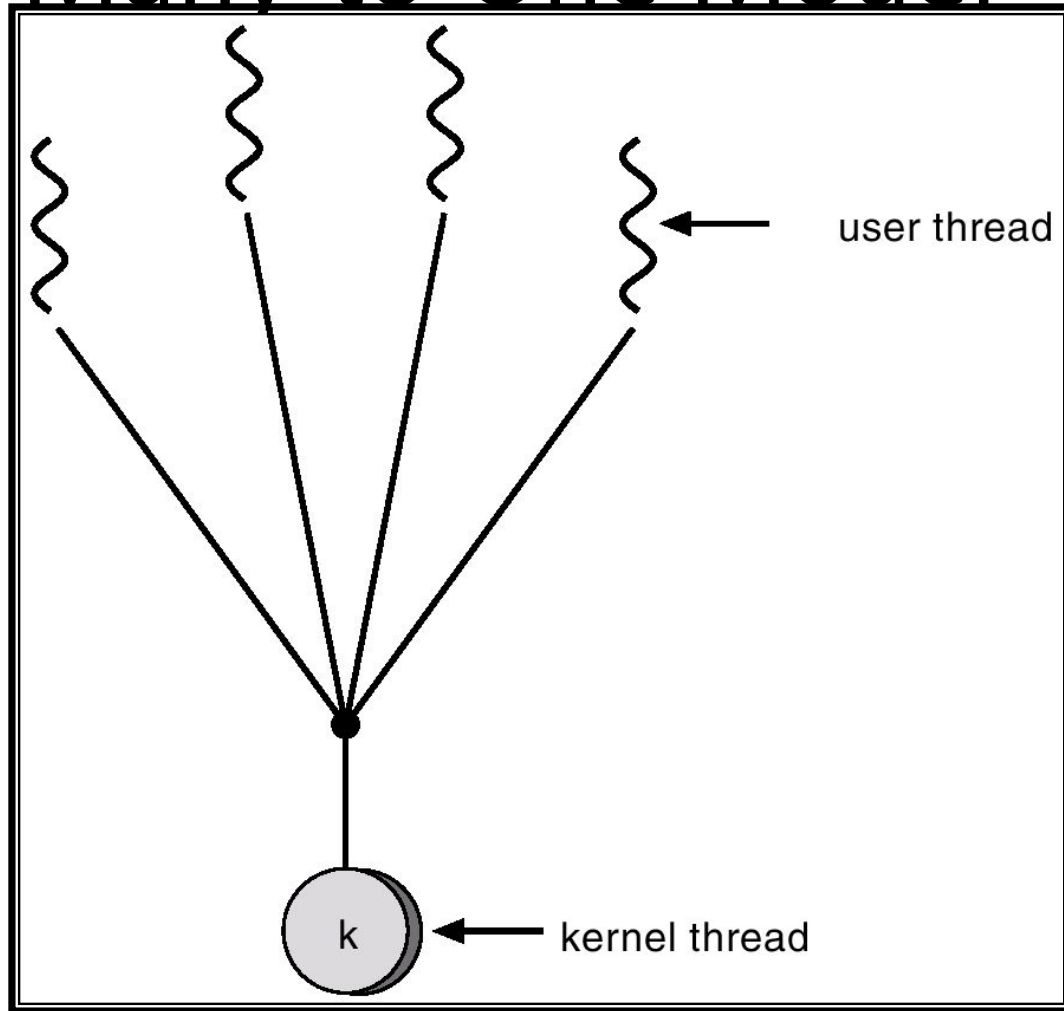
# Multithreading Models

- 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.

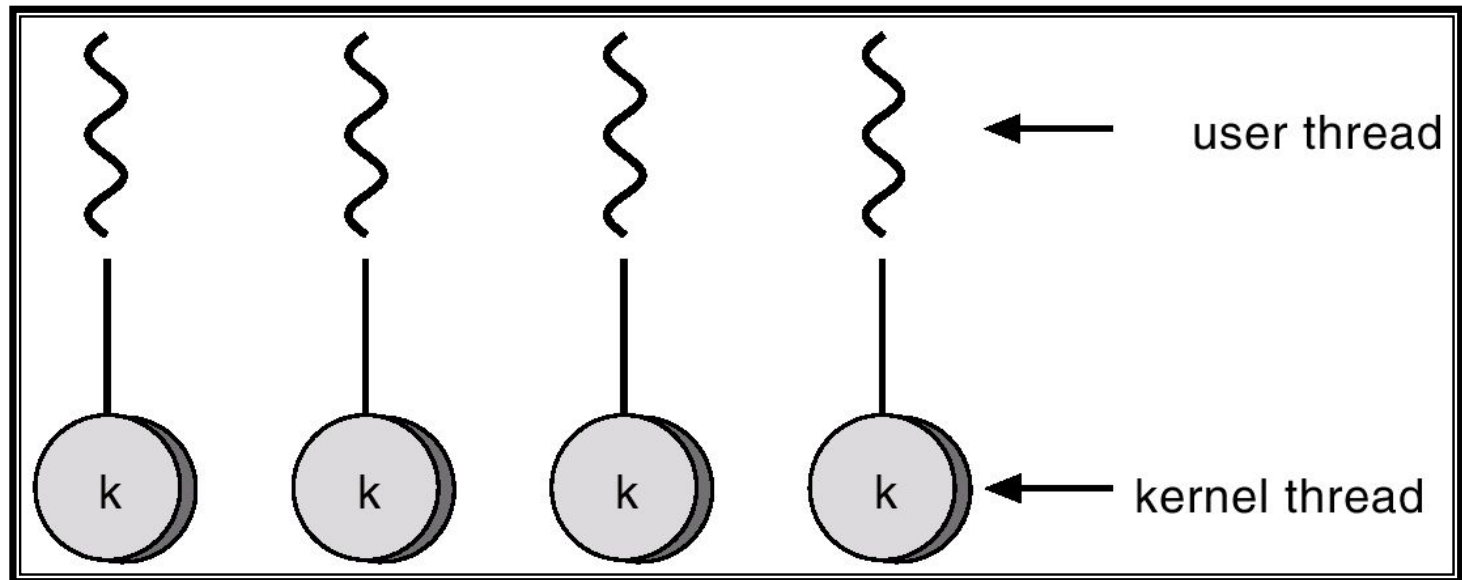
# Many-to-One Model



# One-to-One

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

# One-to-one Model

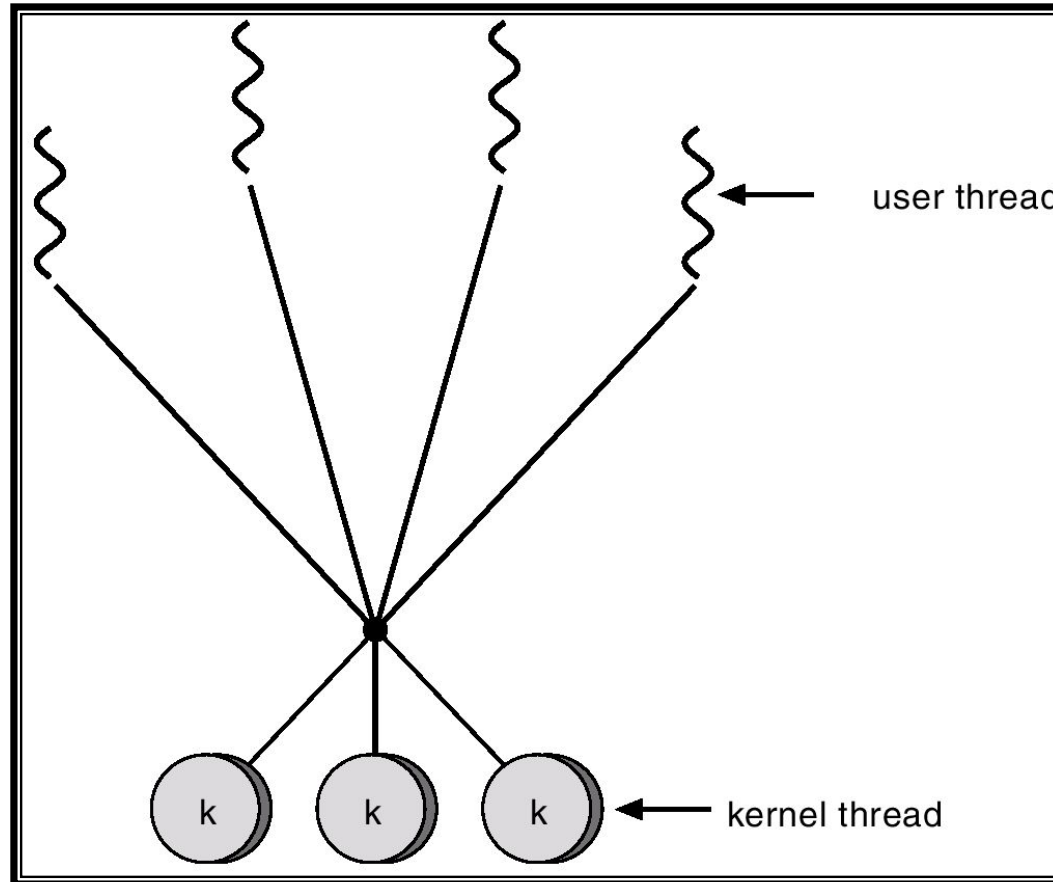


# Many-to-Many Model

- 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



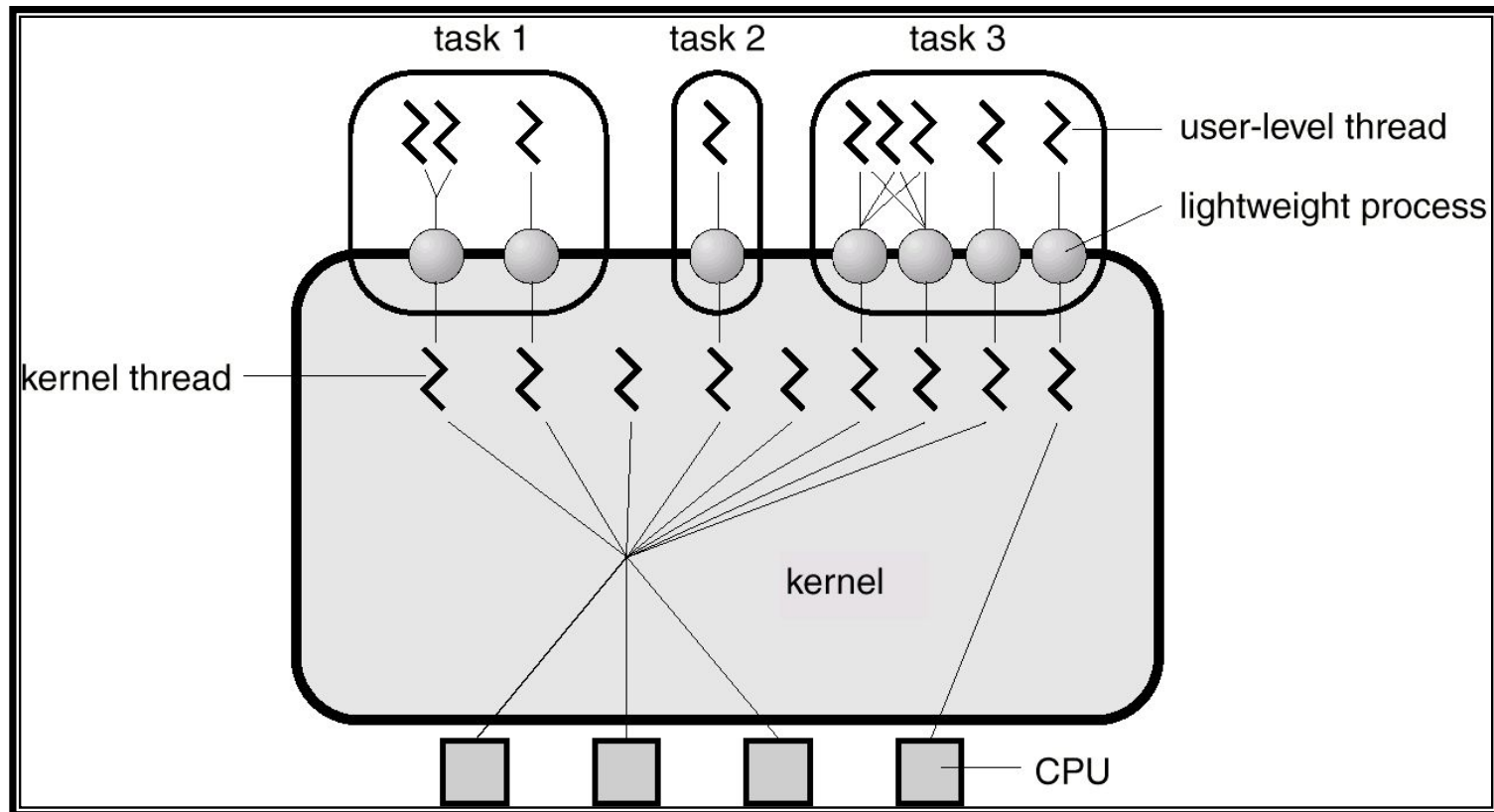
# Threading Issues

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

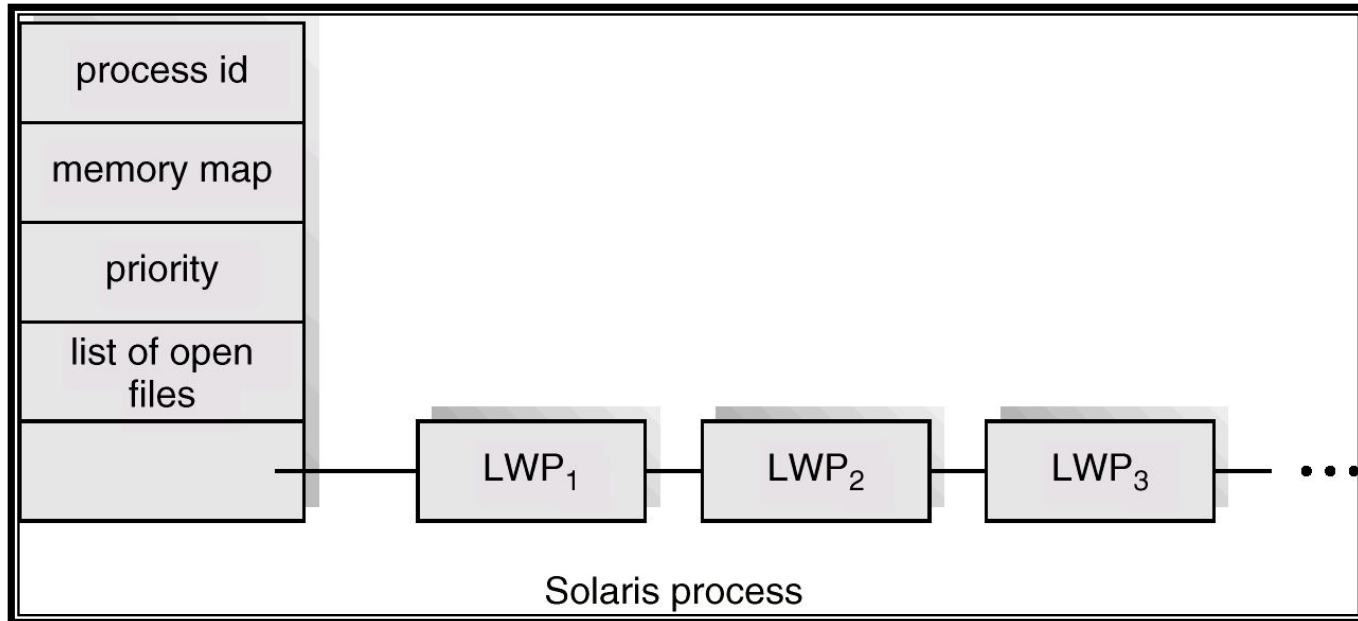
# Pthreads

- 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

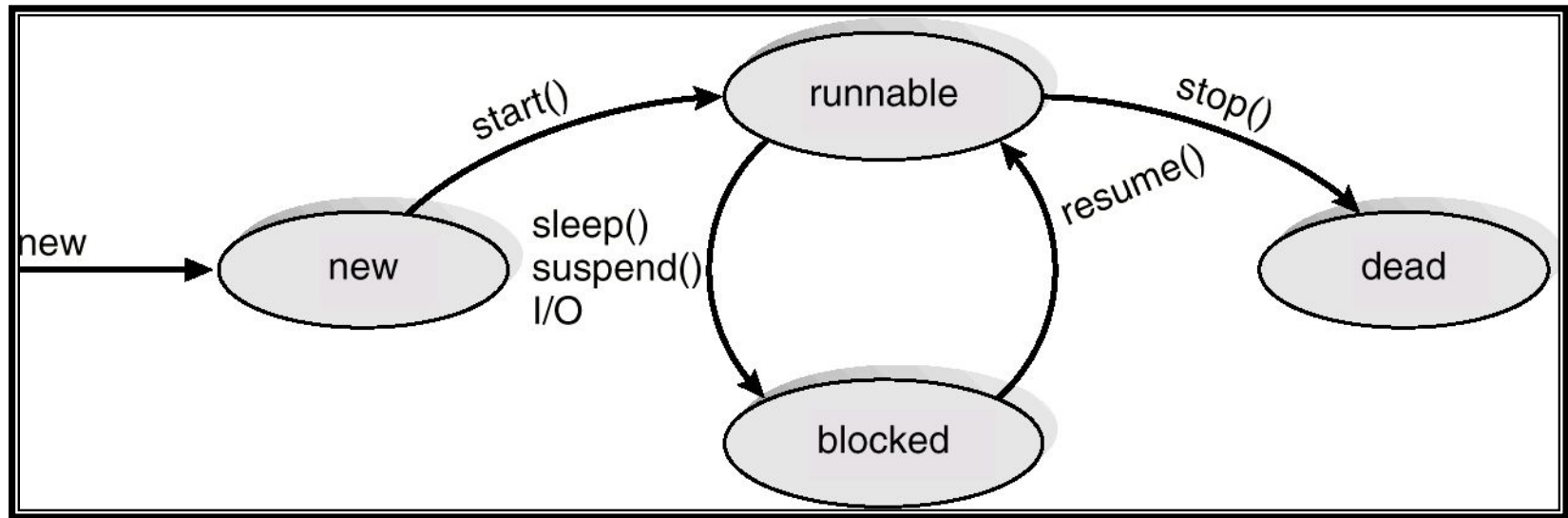
- Linux refers to them as *tasks* rather than *threads*.
- Thread creation is done through clone() system call.
- Clone() allows a child task to share the address space of the parent task (process)

# 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



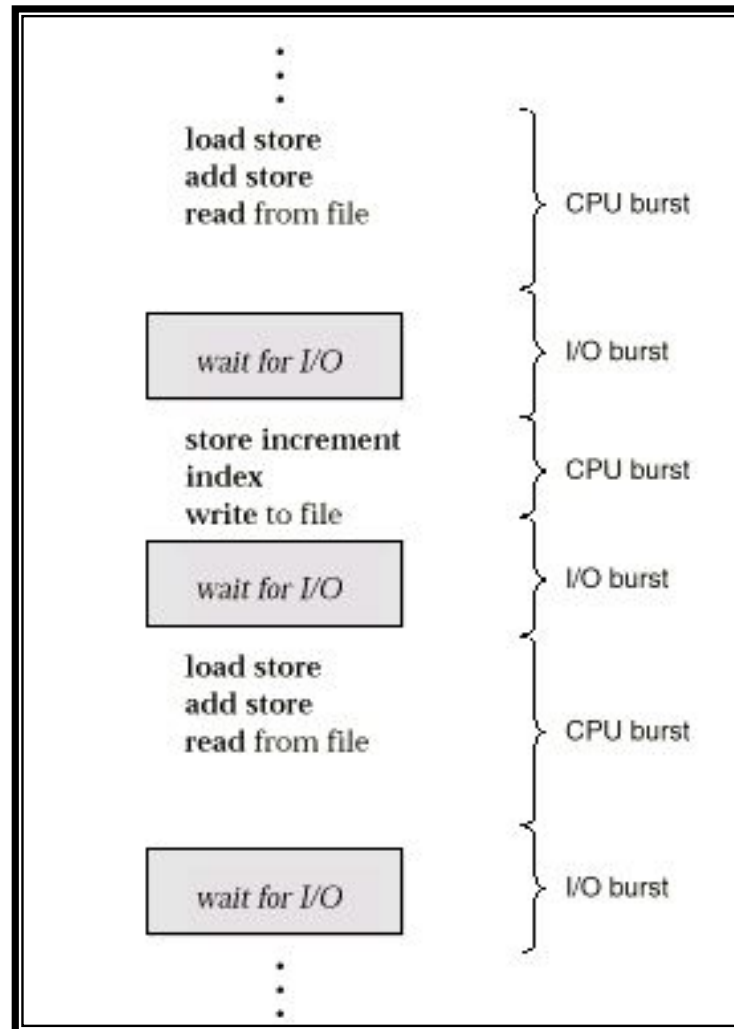
# Chapter 6: CPU Scheduling

- Basic Concepts
- Scheduling Criteria
- Scheduling Algorithms
- Multiple-Processor Scheduling
- Real-Time Scheduling
- Algorithm Evaluation

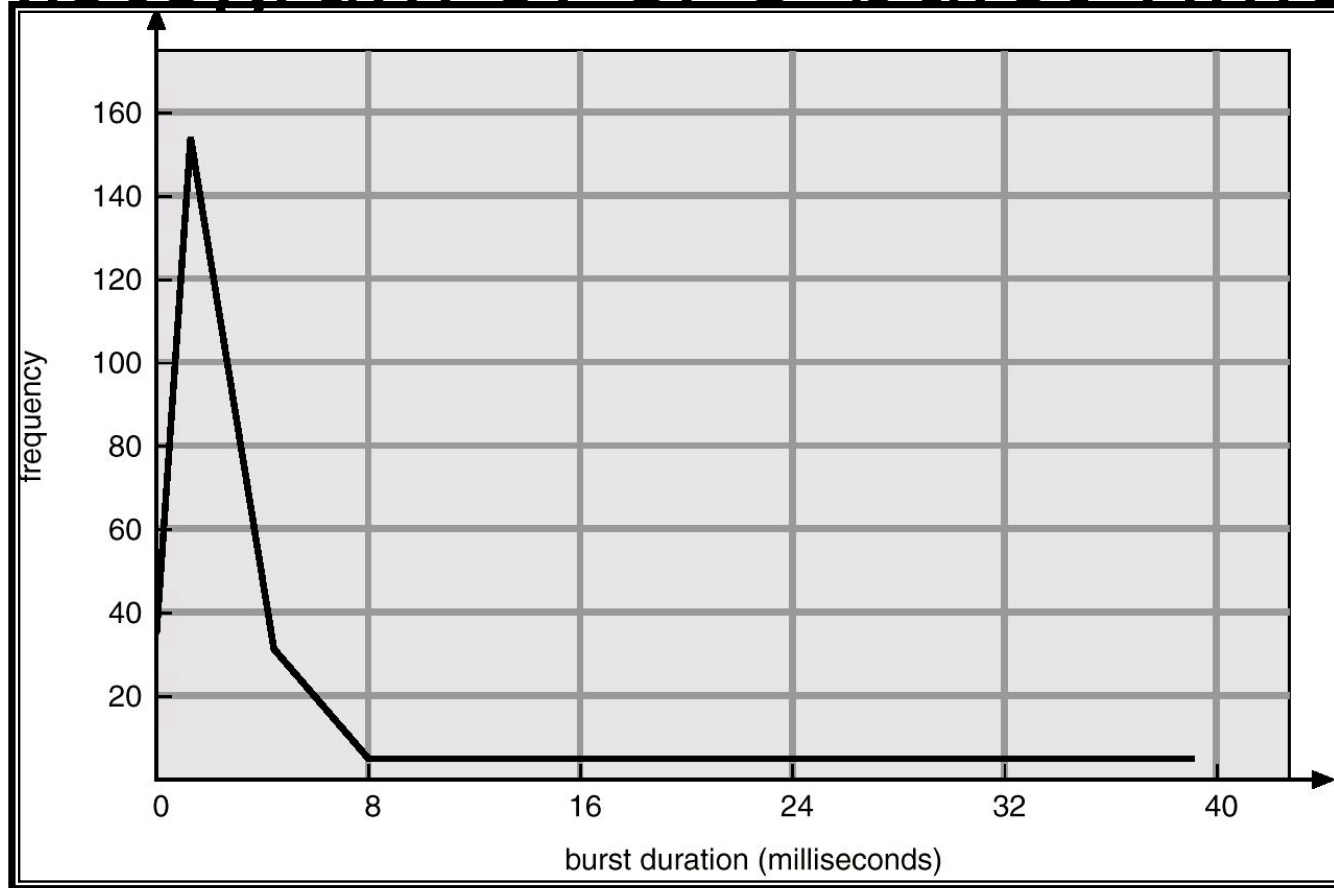
# Basic Concepts

- 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

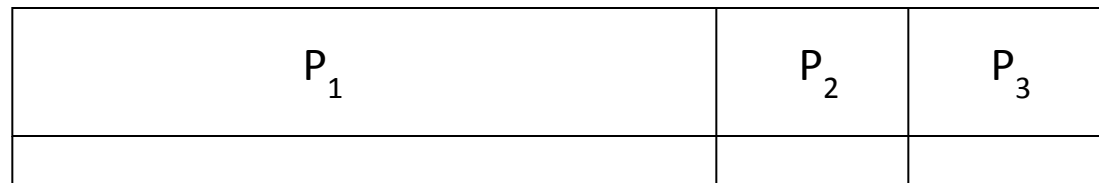
Process Burst Time

$P_1$  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:



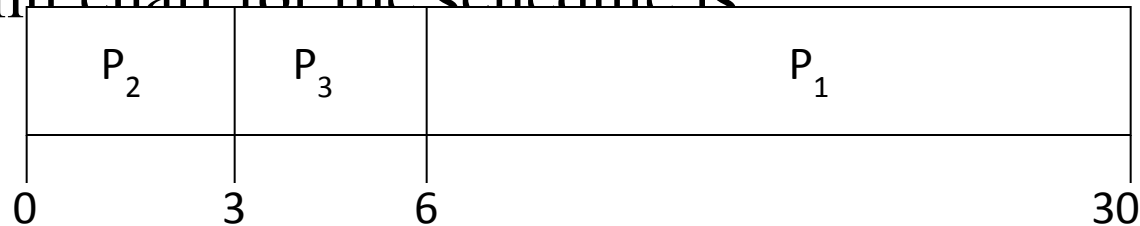
- Waiting time for  $P_1 = 0$ ;  $P_2 = 24$ ;  $P_3 = 27$
- Average waiting time:  $(0 + 24 + 27)/3 = 17$

# FCFS Scheduling (Cont.)

Suppose that the processes arrive in the order

$P_2, P_3, P_1$ .

- The Gantt chart for the schedule is:



- 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

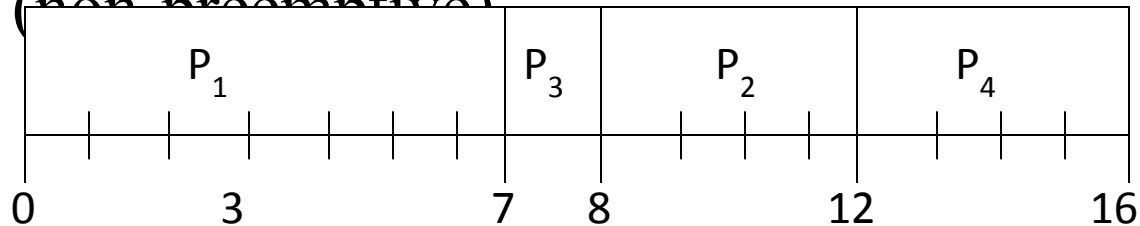
# Shortest-Job-First (SJR) Scheduling

- 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                 |

- SJF (non-preemptive)

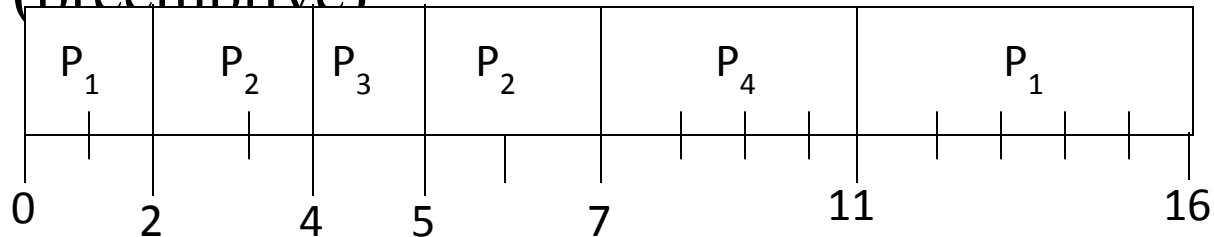


- Average waiting time =  $(0 + 6 + 3 + 7)/4 - 4$

# Example of Preemptive SJF

| <u>Process</u> | <u>Arrival Time</u> | <u>Burst Time</u> |
|----------------|---------------------|-------------------|
| $P_1$          | 0                   | 7                 |
| $P_2$          | 2                   | 4                 |
| $P_3$          | 4                   | 1                 |
| $P_4$          | 5                   | 4                 |

- SJF (preemptive)



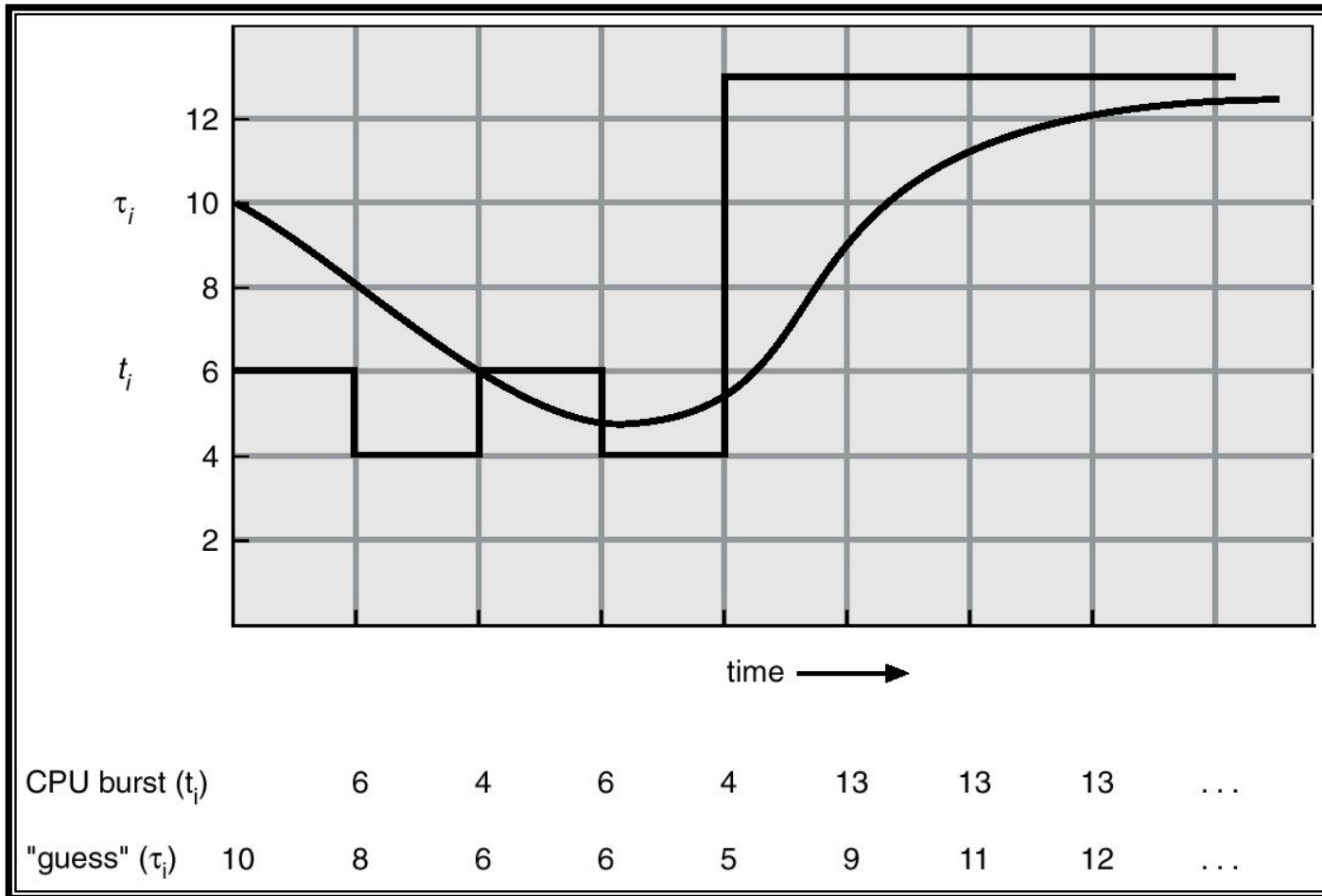
- Average waiting time =  $(9 + 1 + 0 + 2)/4 = 3$

# 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  $n^{th}$  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:
$$\tau_{n+1} = \alpha t_n + (1 - \alpha) \alpha t_n - 1 + \dots$$
$$+ (1 - \alpha)^j \alpha t_n - 1 + \dots$$
$$+ (1 - \alpha)^{n-1} t_n \tau_0$$
- 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

Process   Burst Time

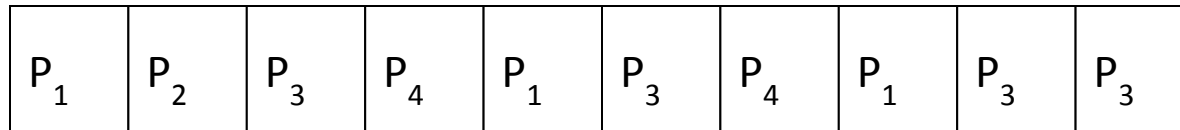
$P_1$    53

$P_2$    17

$P_3$    68

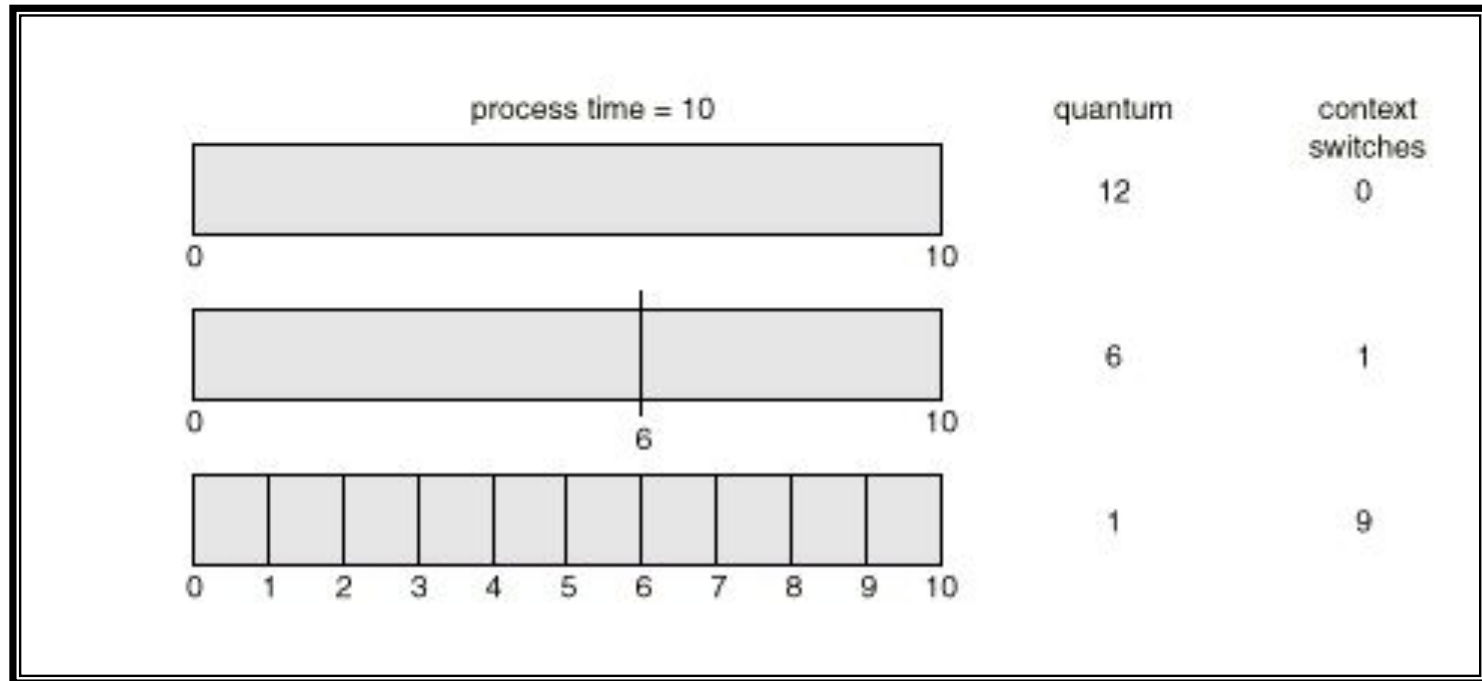
$P_4$    24

- The Gantt chart is:

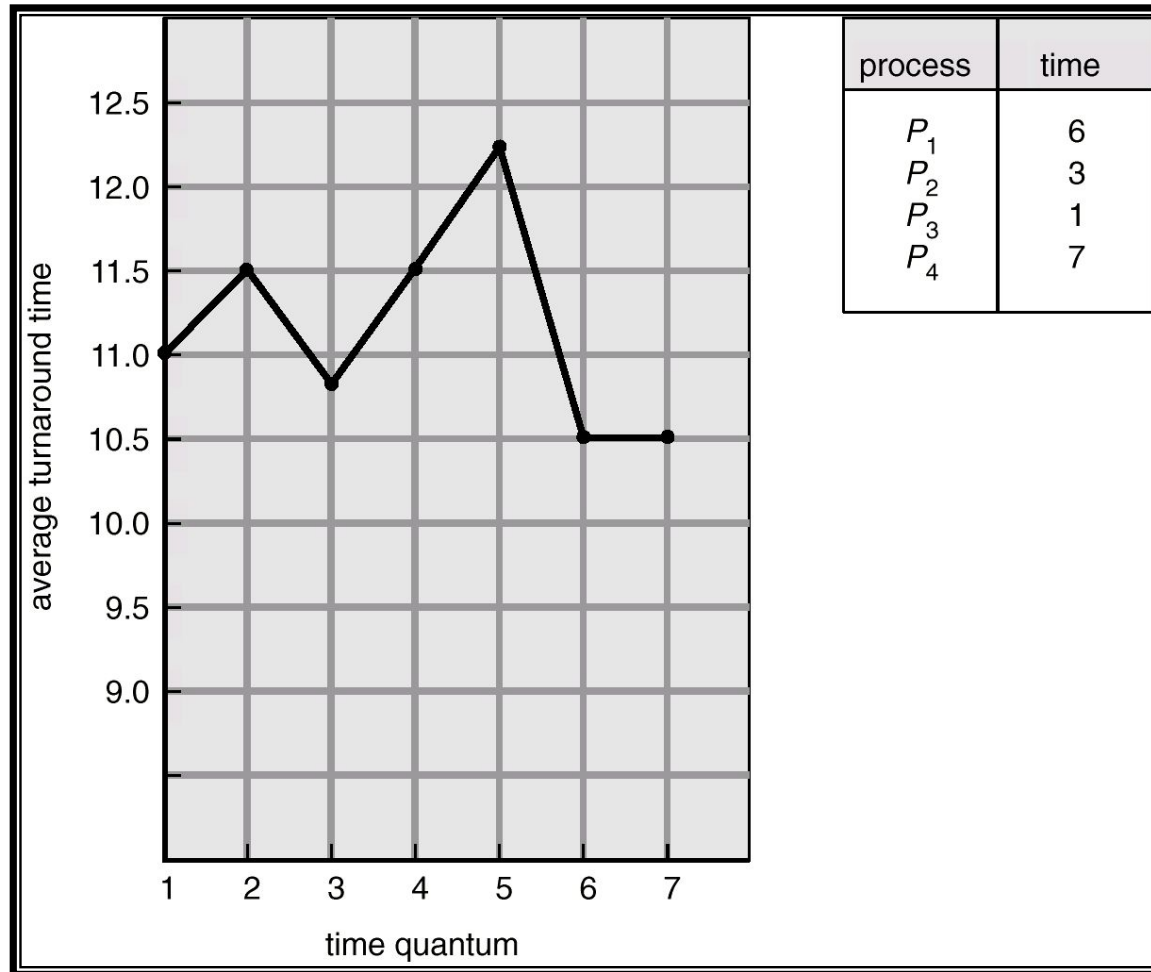


- Typically, higher average turnaround than SJF, but better response.

# Time Quantum and Context Switch Time



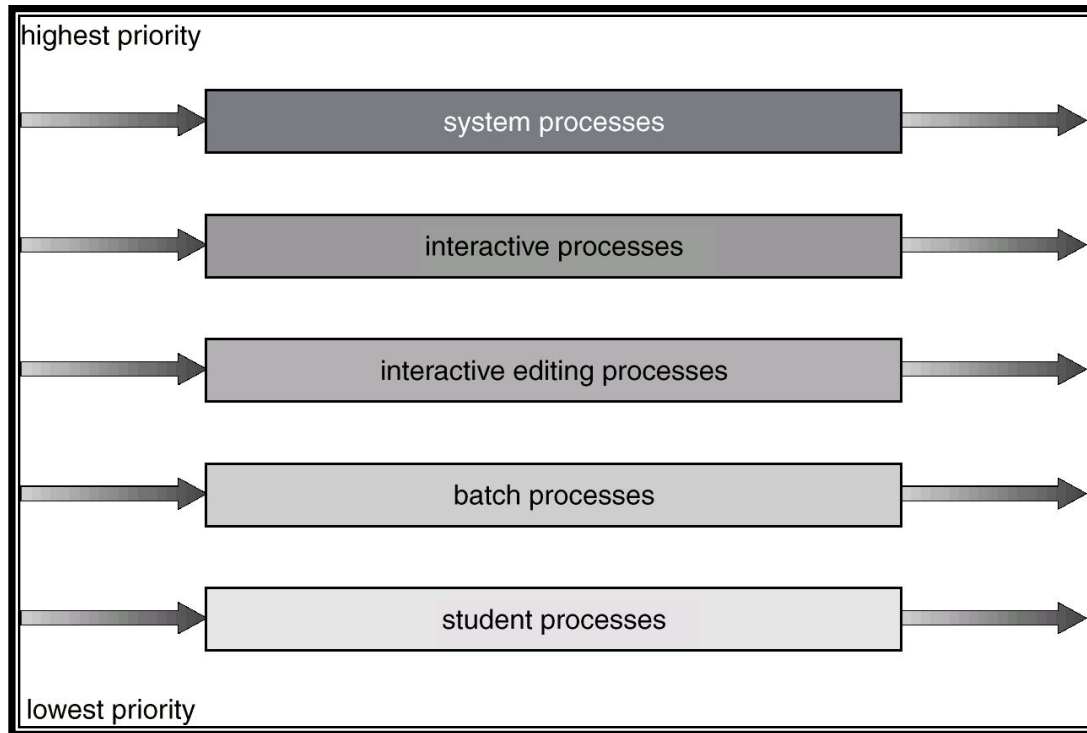
# Turnaround Time Varies With The Time Quantum



# 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





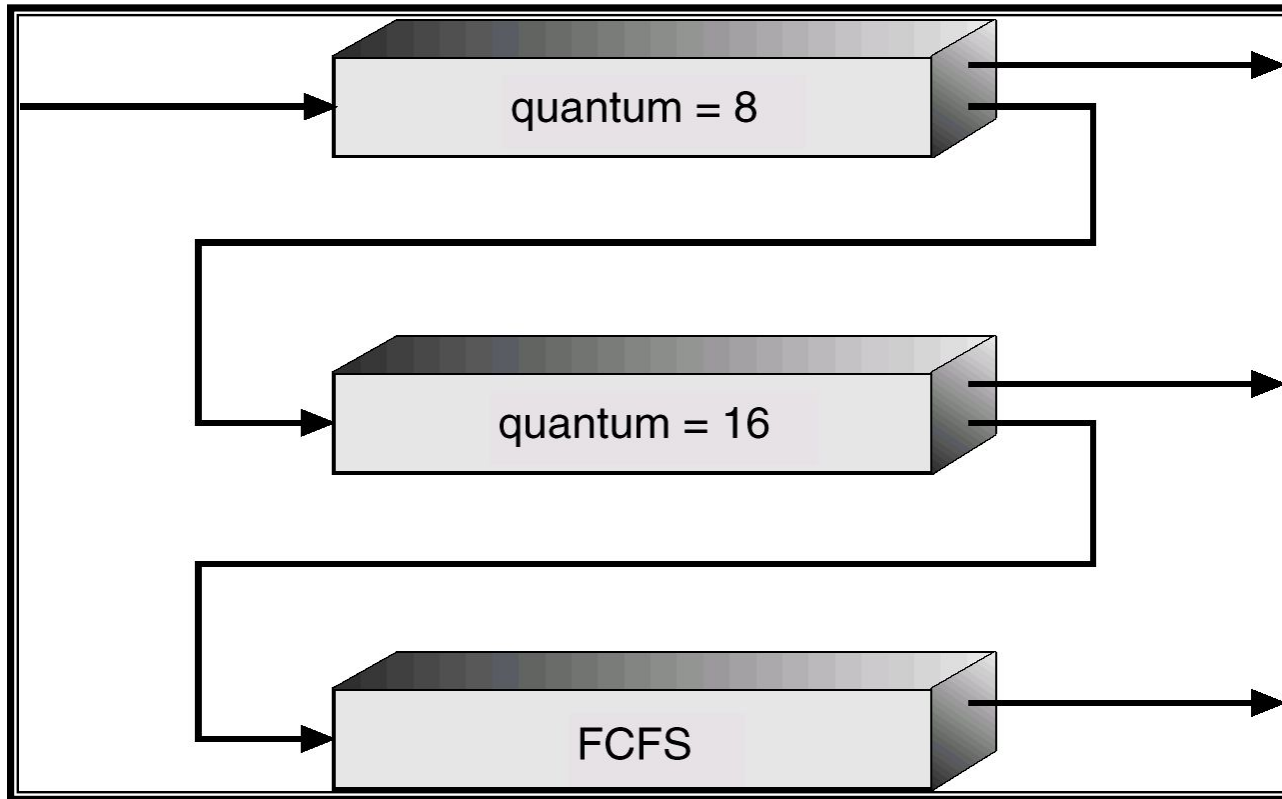
# Multilevel Feedback Queue

- 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



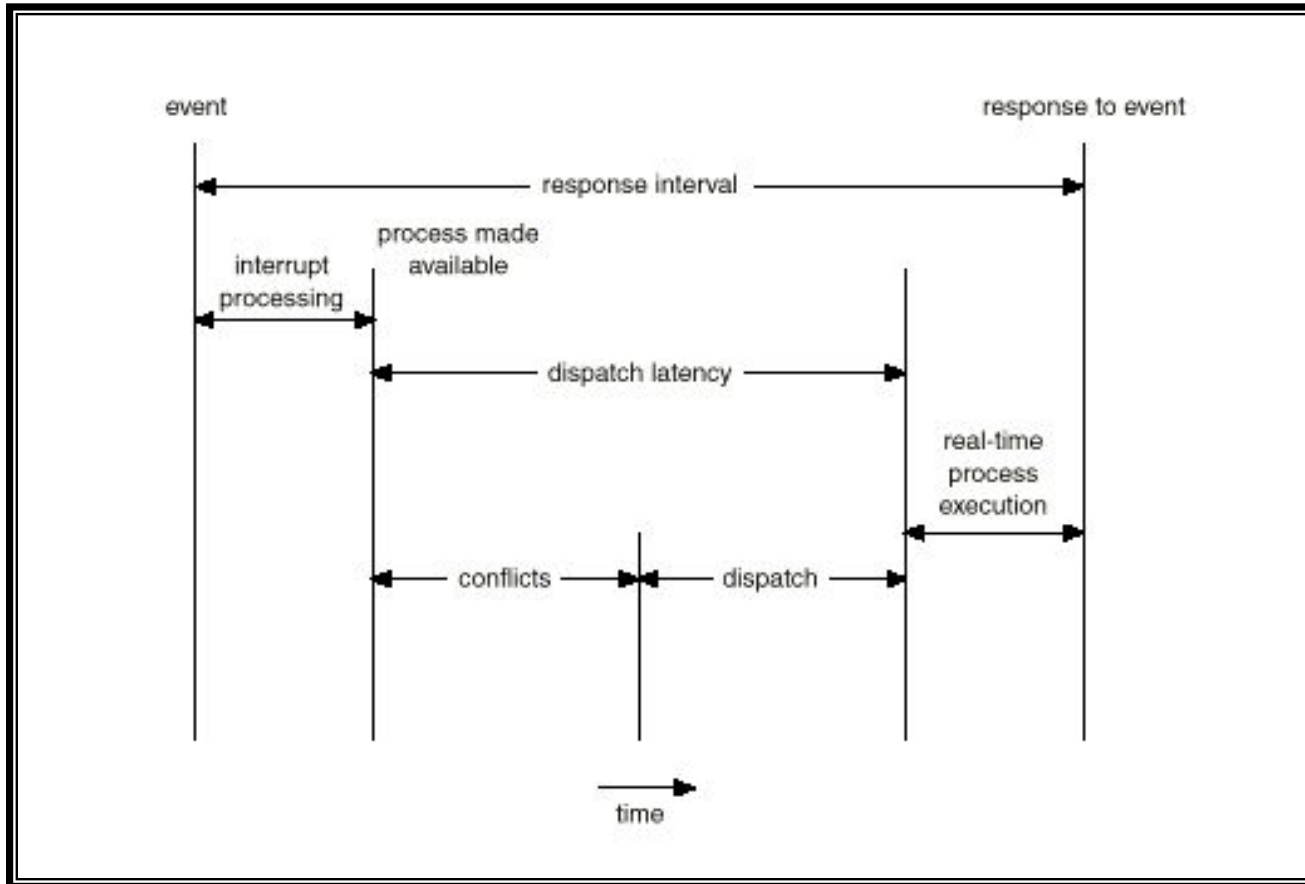
# Multiple-Processor Scheduling

- 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.

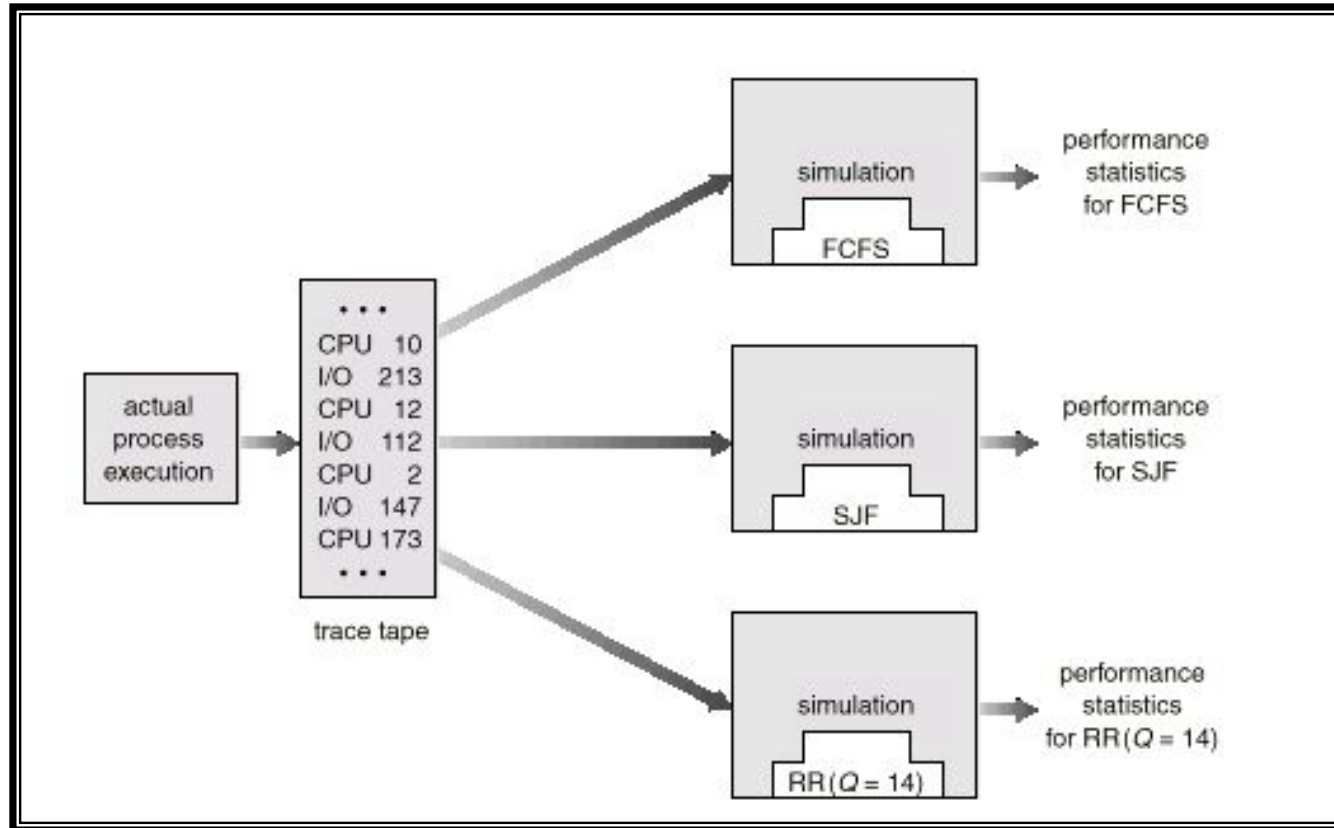
# Dispatch Latency



# Algorithm Evaluation

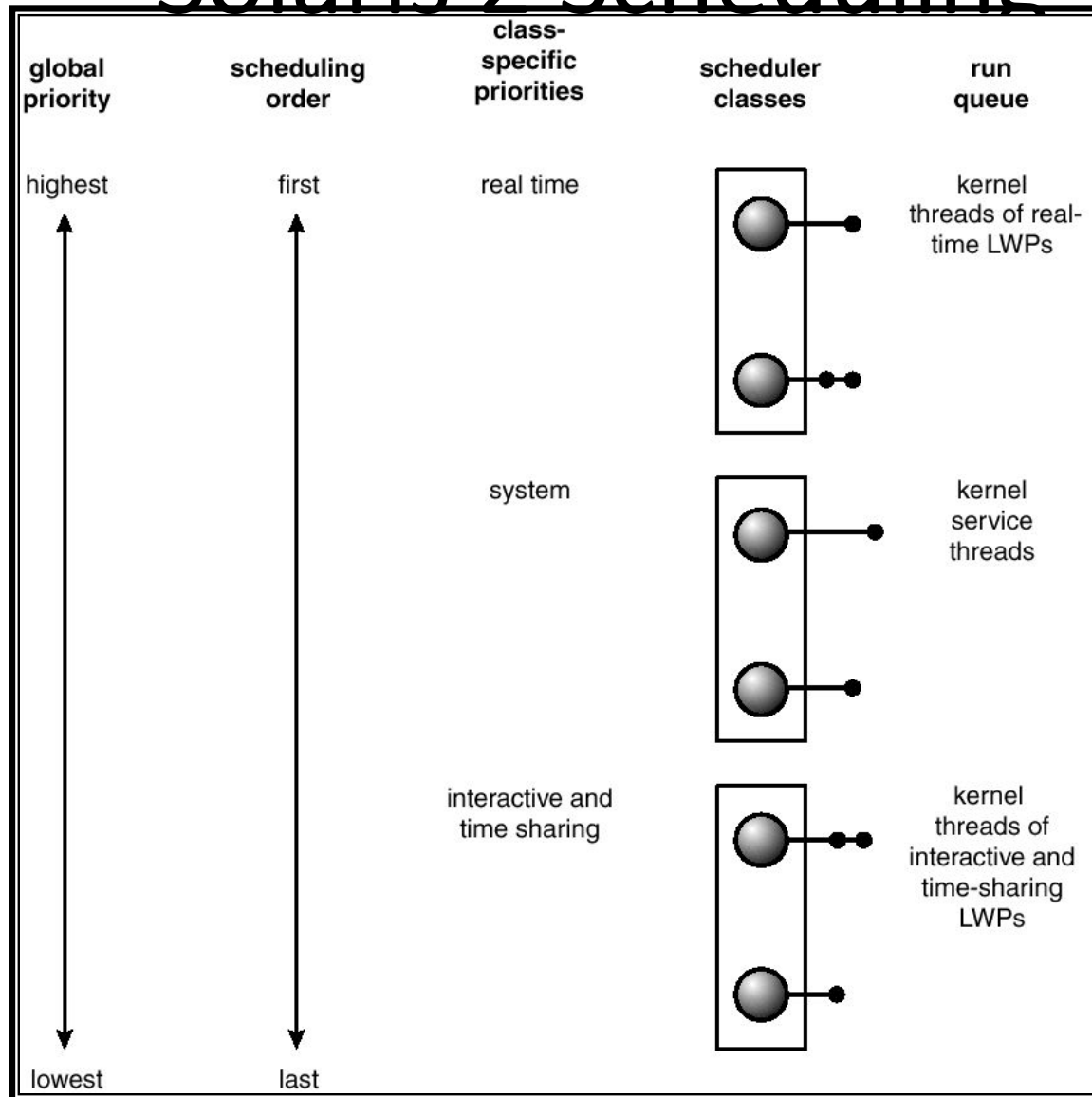
- Deterministic modeling – takes a particular predetermined workload and defines the performance of each algorithm for that workload.
- Queueing models
- Implementation

# Evaluation of CPU Schedulers by Simulation





# Solaris 2 Scheduling



# Windows 2000 Priorities

|               | real-time | high | above normal | normal | below normal | idle priority |
|---------------|-----------|------|--------------|--------|--------------|---------------|
| time-critical | 31        | 15   | 15           | 15     | 15           | 15            |
| highest       | 26        | 15   | 12           | 10     | 8            | 6             |
| above normal  | 25        | 14   | 11           | 9      | 7            | 5             |
| normal        | 24        | 13   | 10           | 8      | 6            | 4             |
| below normal  | 23        | 12   | 9            | 7      | 5            | 3             |
| lowest        | 22        | 11   | 8            | 6      | 4            | 2             |
| idle          | 16        | 1    | 1            | 1      | 1            | 1             |

# 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 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*  
        critical 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**  $\Rightarrow P_i$  can enter its critical section
- Process  $P_i$ 
  - do** {
    - while (turn != i) ;**  
critical section
    - turn = i;**  
remainder 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 = i;**
    - 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  $<\equiv$  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:

$P_i$   $P_j$   
 $\square$   $\square$

$A$   $wait(flag)$

$signal(flag)$   $B$

# 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

$P_0$     $P_1$

*wait(S); wait(Q);*

*wait(Q); wait(S);*

□   □

*signal(S); signal(Q);*

*signal(Q) signal(S);*

- **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

...

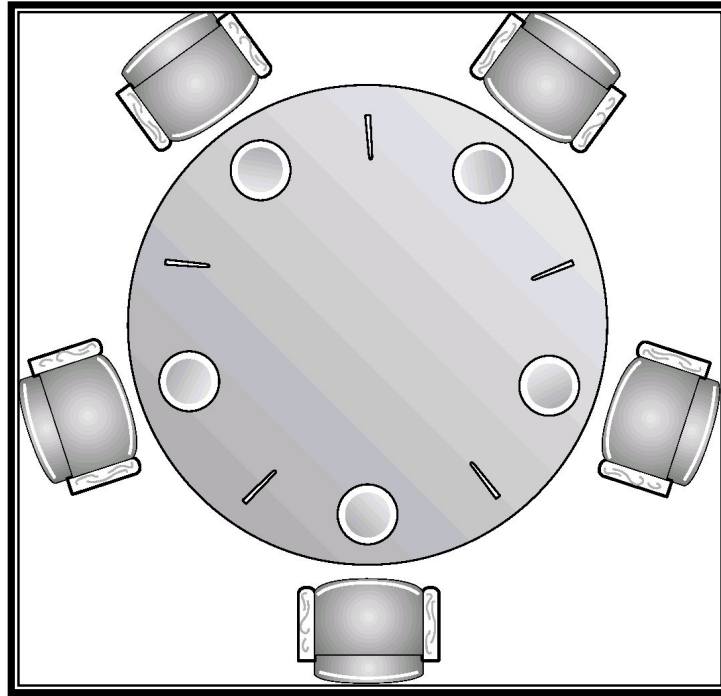
**signal(wrt);**



# 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

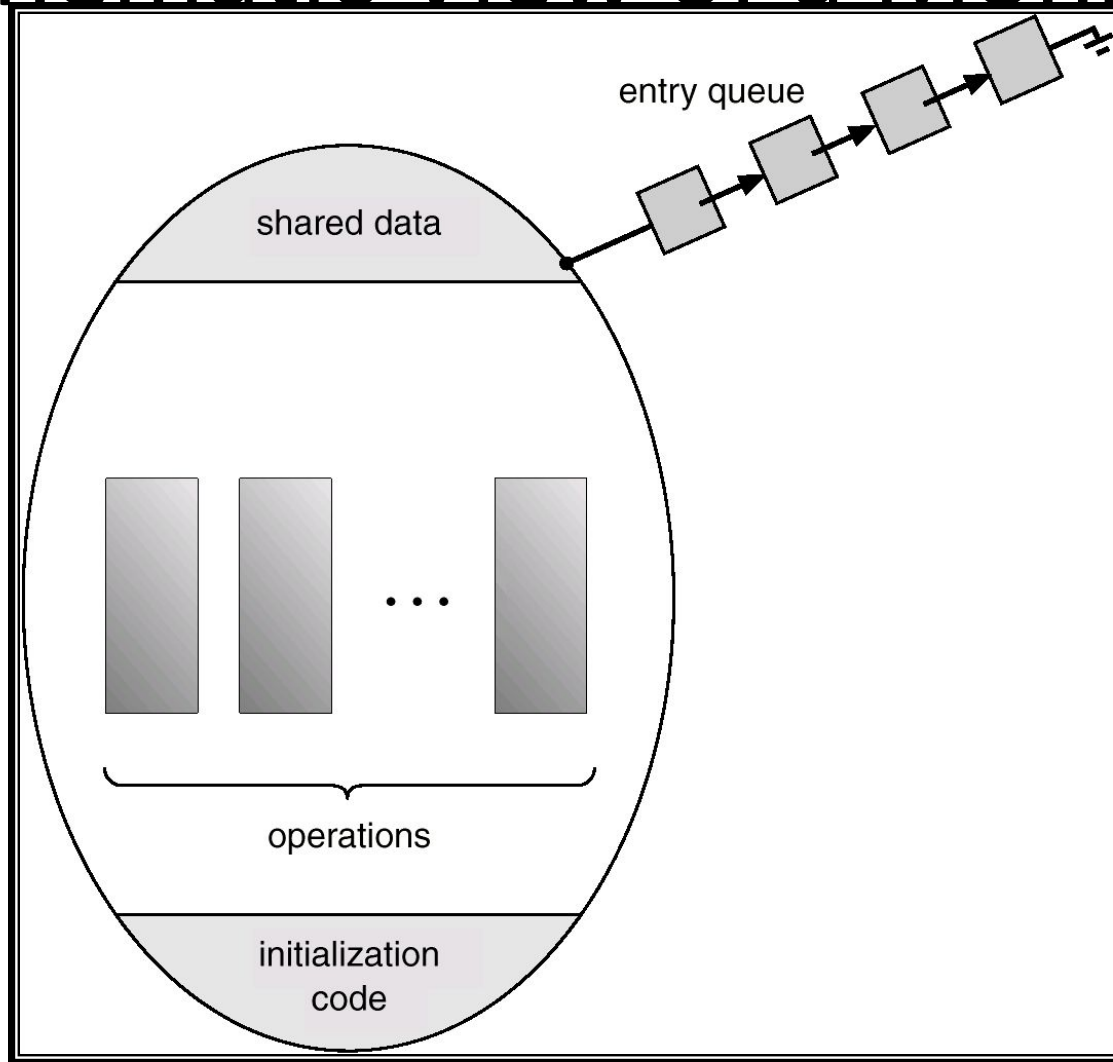
- High-level synchronization construct that allows the safe sharing of an abstract data type among concurrent processes.

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

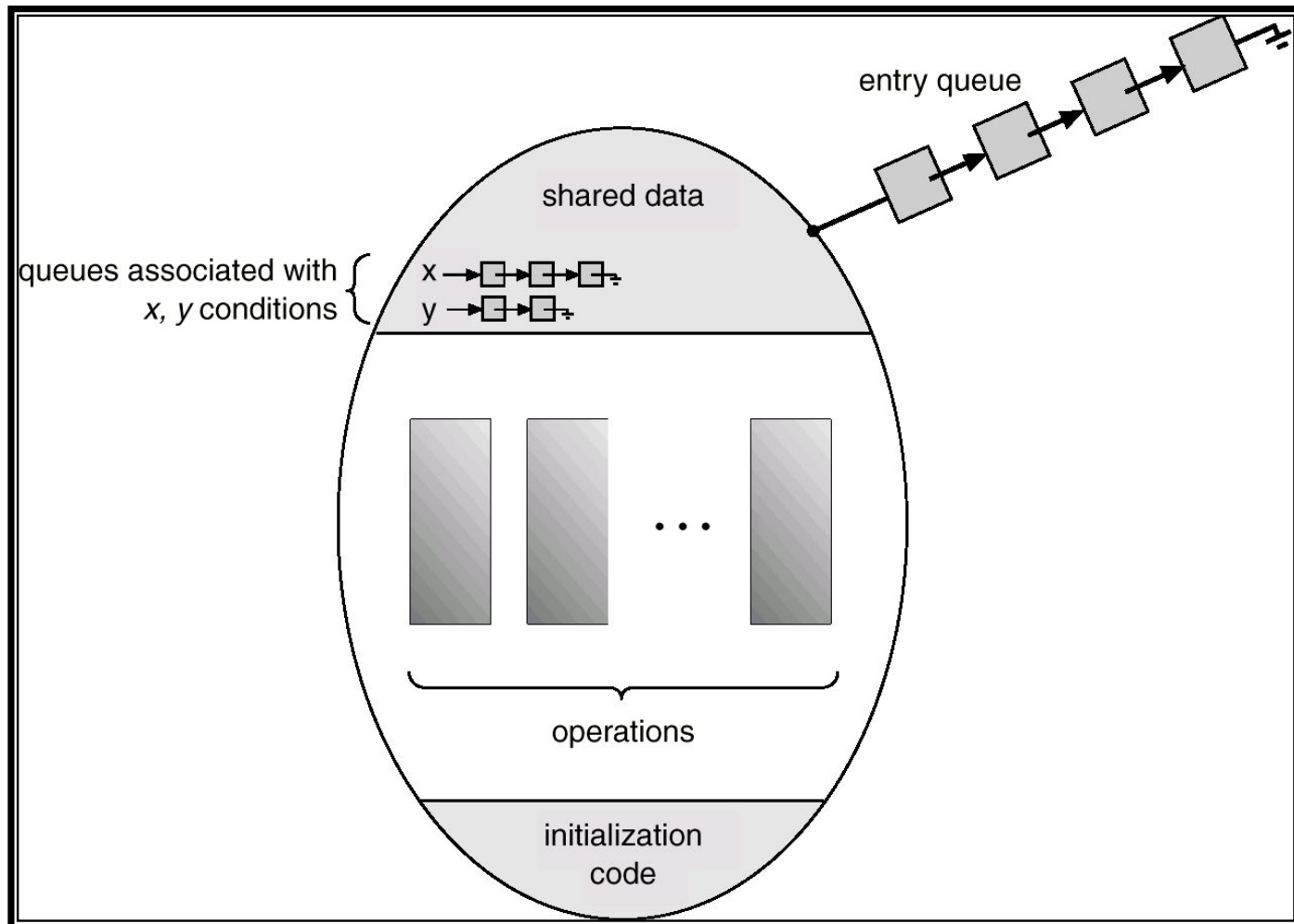
# Monitors

- 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  $x.\text{wait}$  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.

# Deadlocks

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

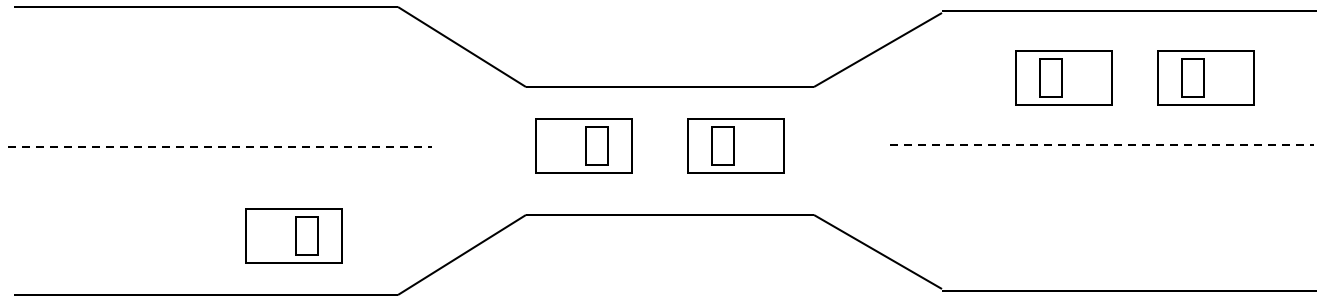


# The Deadlock Problem

- 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

|            |           |
|------------|-----------|
| $P_0$      | $P_1$     |
| $wait(A);$ | $wait(B)$ |
| $wait(B);$ | $wait(A)$ |

# 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  $P_0$ .

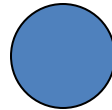
# Resource-Allocation Graph

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.)

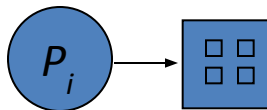
- Process



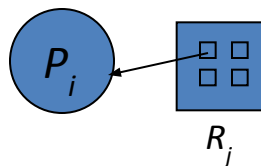
- Resource Type with 4 instances



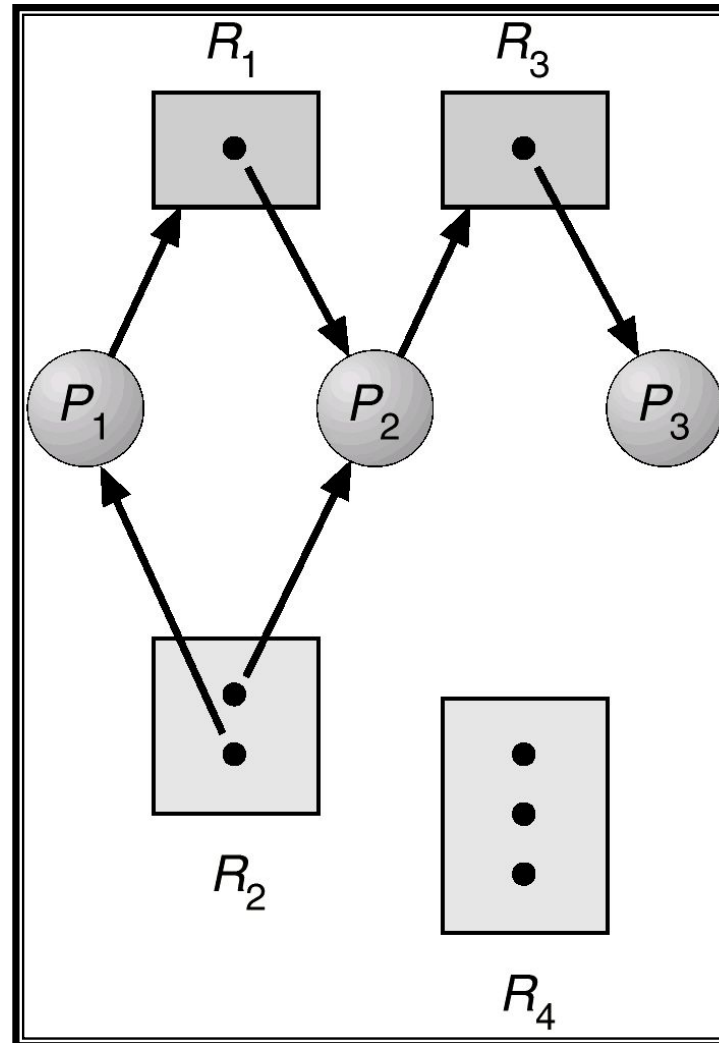
- $P_i$  requests instance of  $R_j$



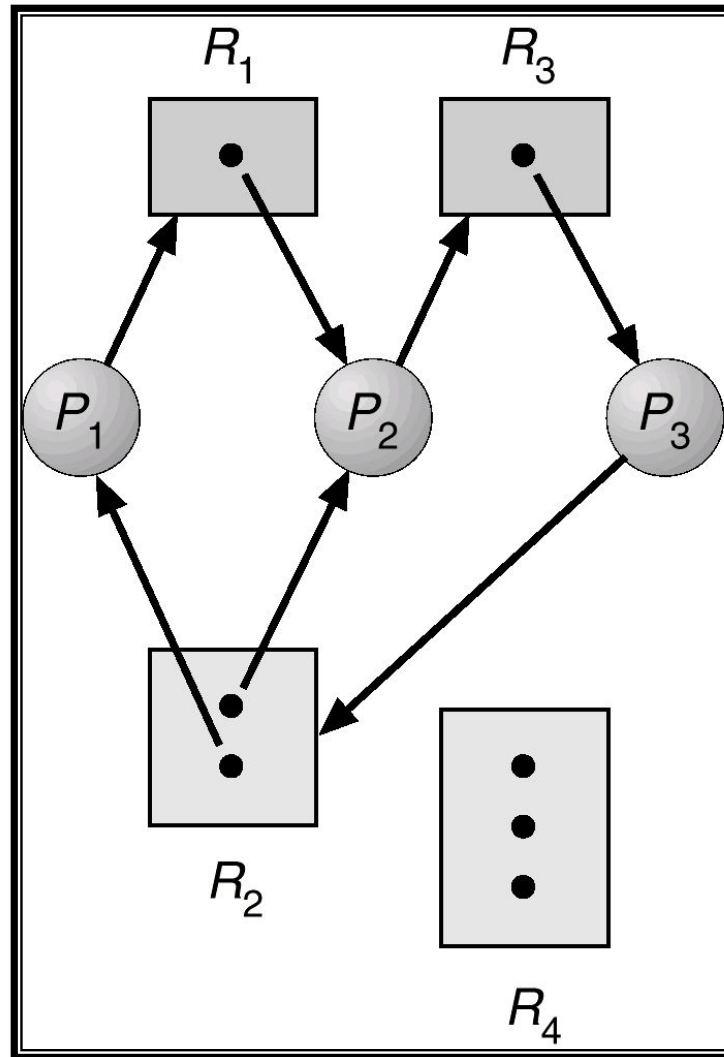
- $P_i$  is holding an instance of  $R_j$



# Example of a Resource Allocation Graph

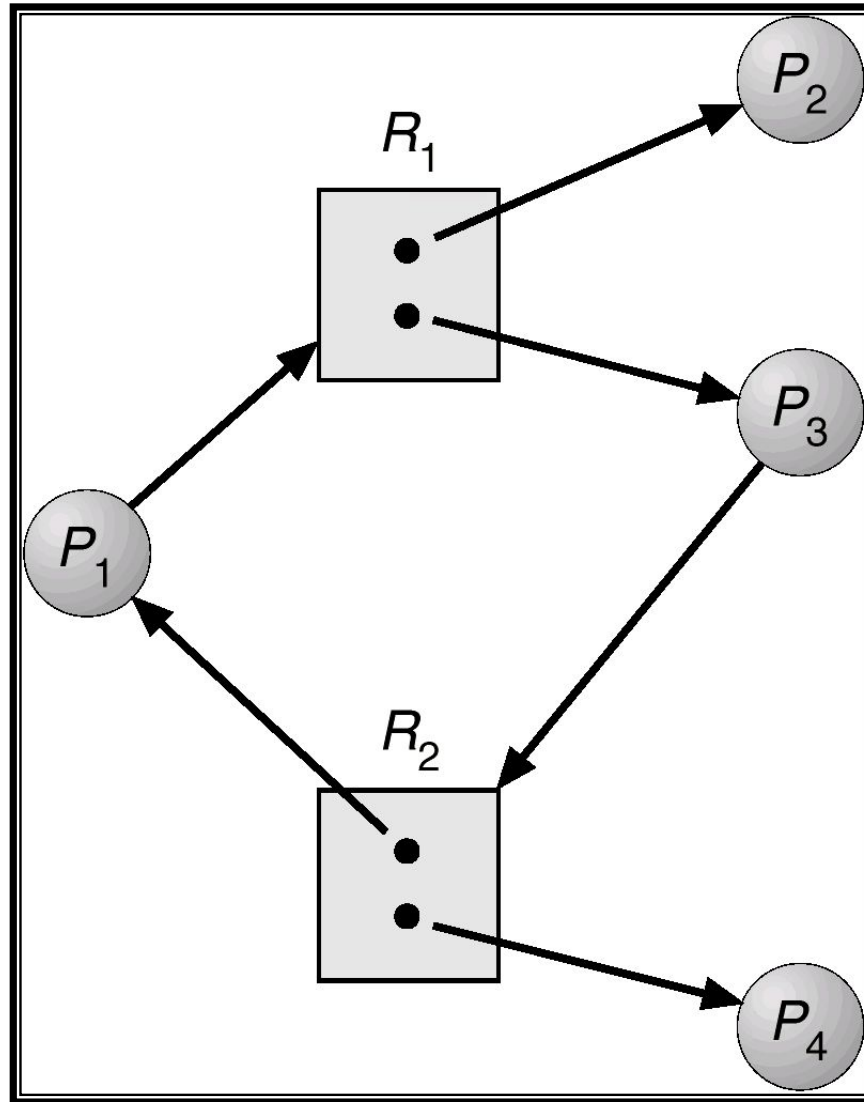


# Resource Allocation Graph With A Deadlock





## 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

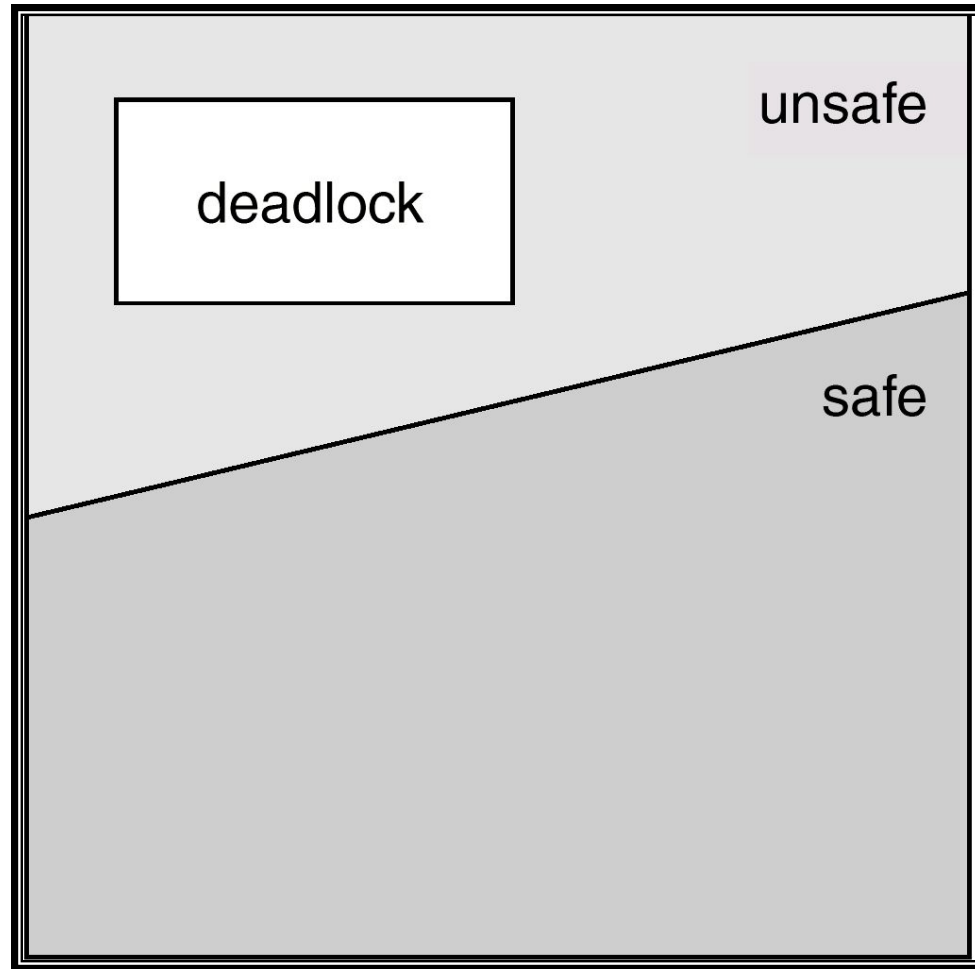
- When a process requests an available resource, system must decide if immediate allocation leaves the system in a *safe state*.
- System is in safe state if there exists a safe sequence of all processes.
- Sequence  $\langle P_1, P_2, \dots, P_n \rangle$  is safe if for each  $P_i$ , the resources that  $P_i$  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.



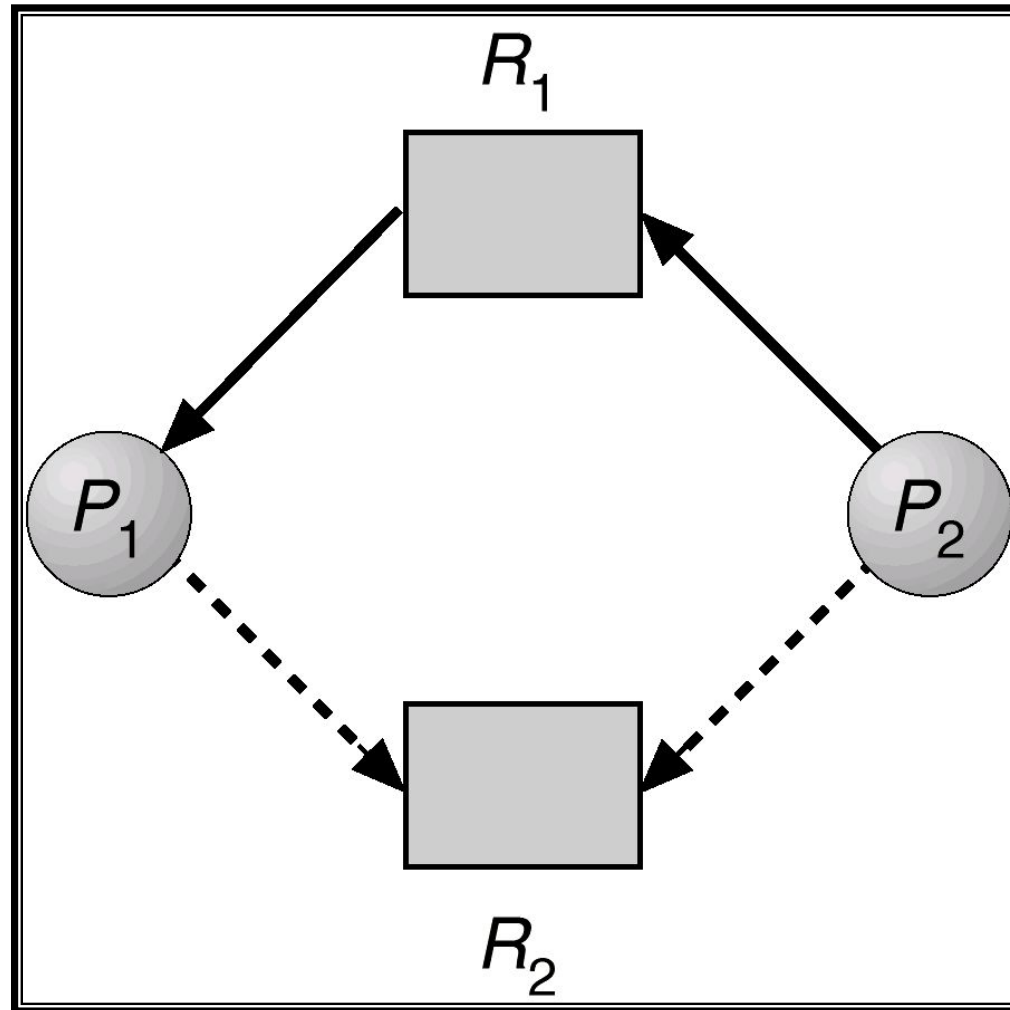
# Safe, Unsafe , Deadlock State



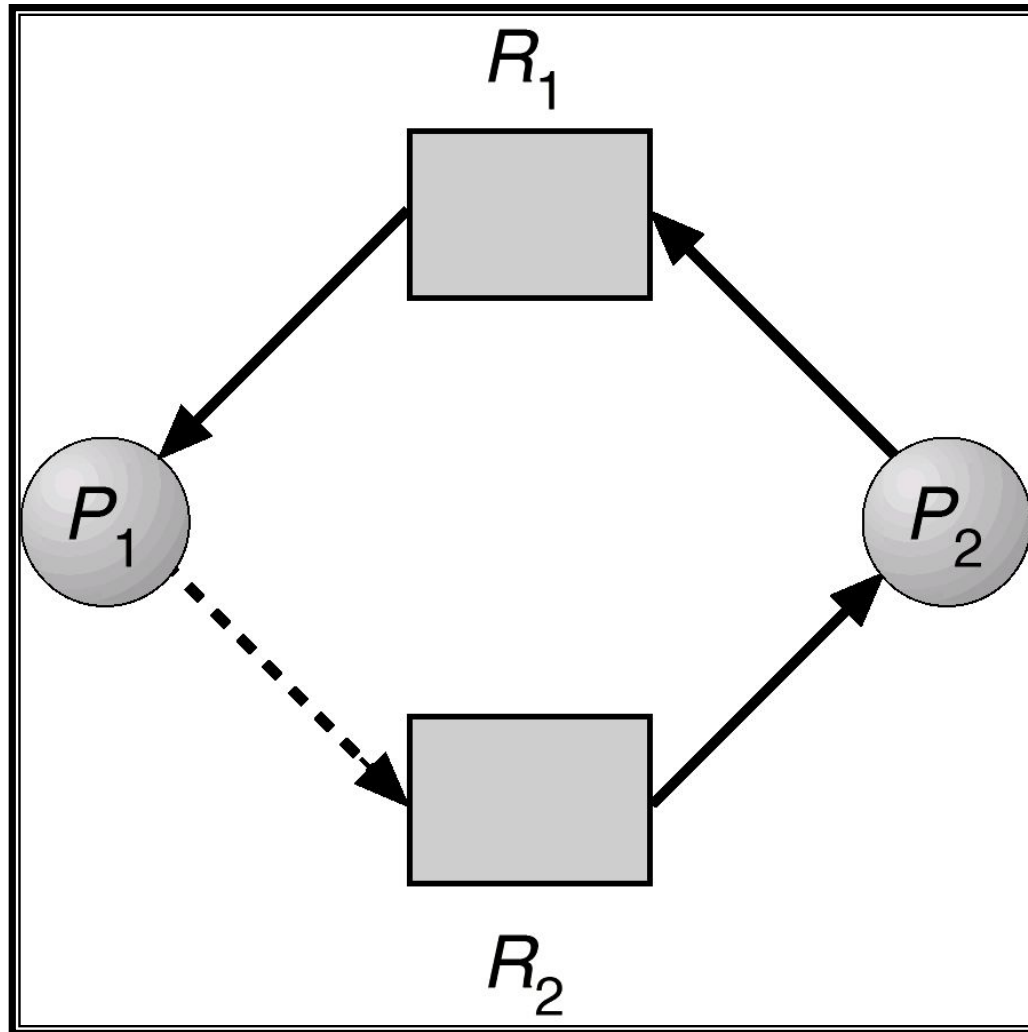
# Resource-Allocation Graph Algorithm

- *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  $i = 1, 2, \dots, n$ .
2. Find an *i* such that both:
  - (a)  $Finish[i] = false$
  - (b)  $Need_i \leq Work$If no such *i* exists, go to step 4.
3.  $Work = Work + Allocation_i$   
 $Finish[i] = true$   
go to step 2.
4. If  $Finish[i] == true$  for all *i*, then the system is in a safe state.

## Resource-Request Algorithm for Process $P_i$

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

1. If  $Request_i \leq Need_i$ , go to step 2. Otherwise, raise error condition, since process has exceeded its maximum claim.
2. If  $Request_i \leq Available$ , go to step 3. Otherwise  $P_i$  must wait, since resources are not available.
3. Pretend to allocate requested resources to  $P_i$  by modifying the state as follows:

$$Available = Available - Request_i;$$

$$Allocation_i = Allocation_i + Request_i;$$

$$Need_i = Need_i - Request_i;$$

- 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> |          |          |
|-------|-------------|----------|----------|
|       | <i>A</i>    | <i>B</i> | <i>C</i> |
| $P_0$ | 7           | 4        | 3        |
| $P_1$ | 1           | 2        | 2        |
| $P_2$ | 6           | 0        | 0        |
| $P_3$ | 0           | 1        | 1        |
| $P_4$ | 4           | 3        | 1        |

- 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,  $(1,0,2) \leq (3,3,2) \Rightarrow \text{true}$ ).

|       | <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?

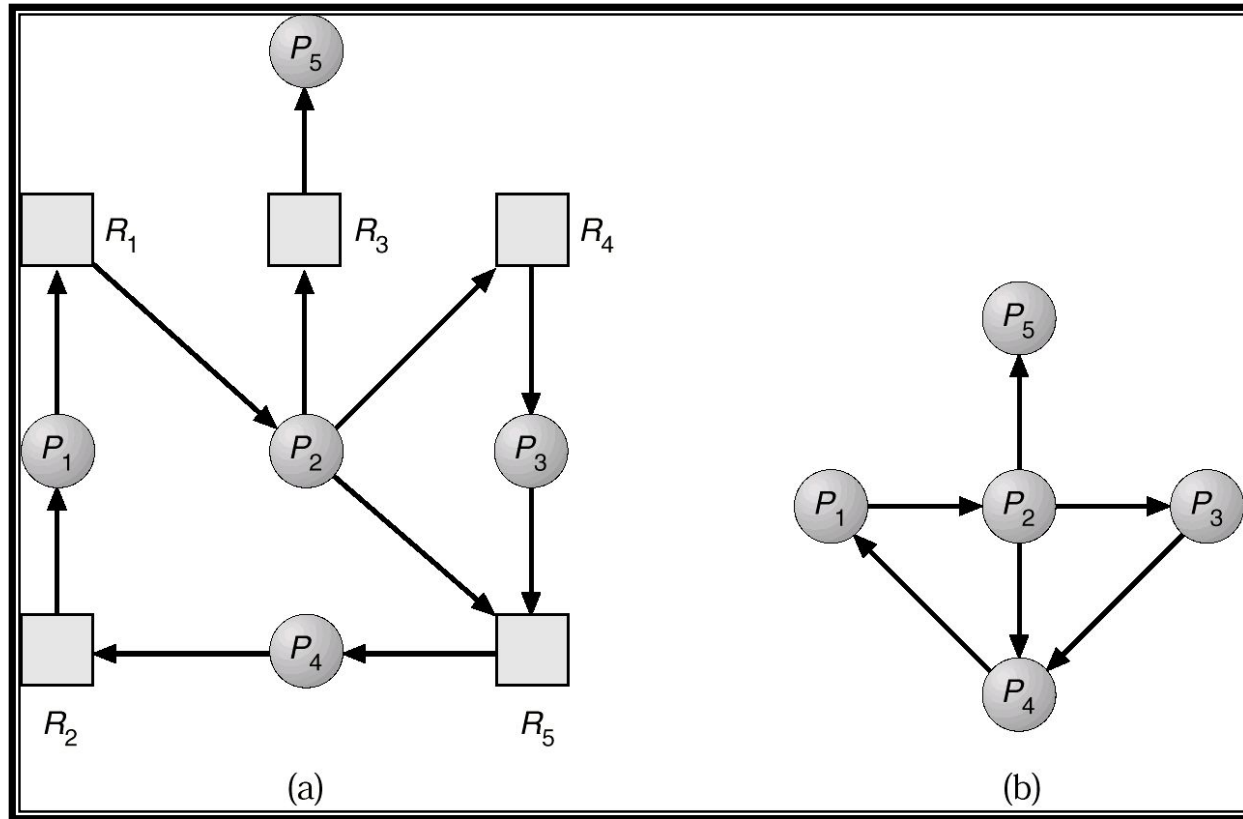
# Deadlock Detection

- 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  $C$ .

Request

$A \ B \ C$

$P_0$  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.