

**PDFZilla – Unregistered**

**PDFZilla - Unregistered**

**PDFZilla - Unregistered**

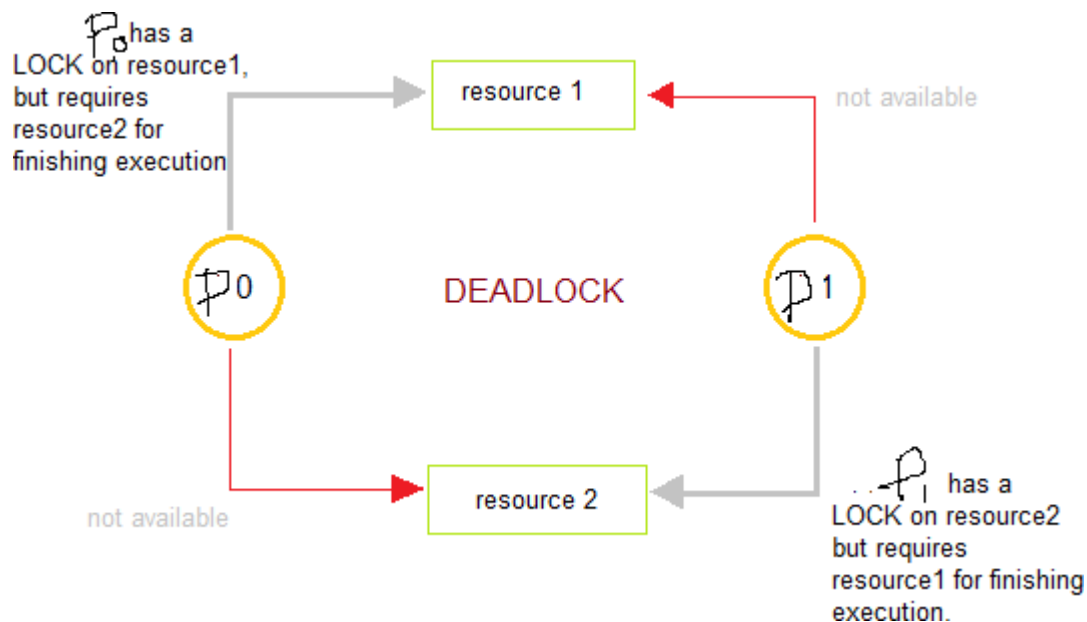
## Module – 3

### Deadlocks

- System model;
- Deadlock characterization;
- Methods for handling deadlocks;
- Deadlock prevention;
- Deadlock avoidance;
- Deadlock detection and recovery from deadlock.

In a multiprogramming system, numerous processes get competed for a finite number of resources. Any process requests resources and as the resources aren't available at that time, the process goes into a waiting state. At times, a waiting process is not at all able again to change its state as the resources it has requested are detained by other waiting processes. That condition is termed as deadlock.

Deadlocks are a set of blocked processes each holding a resource and waiting to acquire a resource held by another process.



### 4.1 System Model

A system consists of finite number of resources and is distributed among number of processes. A process must **request** a resource before using it and it must **release** the resource after using it. It can request any number of resources to carry out a designated task. The amount of resource requested may not exceed the total number of resources available.

- ✓ A process may utilize the resources in only the following **sequence**,
  1. **Request:** If the request is not granted immediately then the requesting process must wait it can acquire the resources.
  2. **Use:** The process can operate on the resource.
  3. **Release:** The process releases the resource after using it.

- ✓ To illustrate deadlock, consider a system with one printer and one tape drive. If a process  $P_i$  currently holds a printer and a process  $P_j$  holds the tape drive. If process  $P_i$  request a tape drive and process  $P_j$  request a printer then a deadlock occurs.
- ✓ Multithread programs are good candidates for deadlock because they compete for shared resources.

## 4.2 Deadlock Characterization

### • Necessary Conditions

A deadlock situation can occur if the following **4 conditions** occur simultaneously in a system.

- **Mutual Exclusion:** At least one resource must be held in a non-sharable mode; that is, only one process at a time can use the resource. If another process requests for the resource, the requesting process must be delayed until the resource has been released.
- **Hold and Wait:** A process must be holding at least one resource and waiting to acquire additional resources that are currently being held by the other process.
- **No Preemption:** Resources cannot be preempted i.e., only the process holding the resources must release it after the process has completed its task.
- **Circular Wait:** A set  $\{P_0, P_1, \dots, P_n\}$  of waiting process must exist such that  $P_0$  is waiting for a resource held by  $P_1$ ,  $P_1$  is waiting for a resource held by  $P_2$ ,  $P_{n-1}$  is waiting for resource held by process  $P_n$  and  $P_n$  is waiting for the resource held by  $P_0$ .

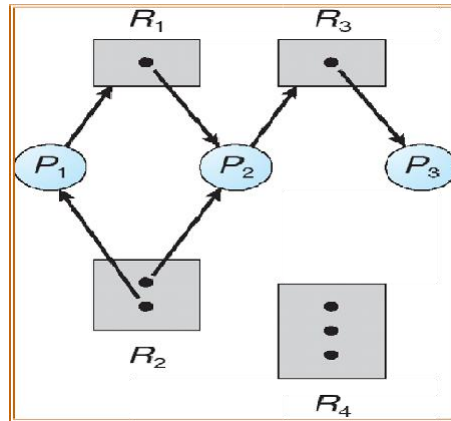
### • Resource Allocation Graph

- ✓ Deadlocks are described by using a directed graph called **system resource allocation graph**. The graph consists of set of **vertices (V)** and set of **edges (E)**.
- ✓ The set of vertices (V) can be described into two different types of nodes,
- ✓  $P = \{P_1, P_2, \dots, P_n\}$  a set consisting of all active processes and  $R = \{R_1, R_2, \dots, R_n\}$  a set consisting of all resource types in the system.
- ✓ A directed edge from process  $P_i$  to resource type  $R_j$  is denoted by  $P_i \rightarrow R_j$ , it indicates that  $P_i$  has requested an instance of resource type  $R_j$  and is waiting for that resource. This edge is called **Request edge**.
- ✓ A directed edge  $R_j \rightarrow P_i$  signifies that an instance of resource type  $R_j$  has been allocated to process  $P_i$ . This is called **Assignment edge**.
- ✓ Process  $P_i$  is represented as **circle** and each resource type  $R_j$  as a **rectangle**.
- ✓ Since resource type  $R_j$  may have more than one instance, we represent each such instance as a **dot** within the rectangle.
- ✓ Request edge points to only the rectangle  $R_j$ , whereas an assignment edge must also designate one of the dots in the rectangle.
- ✓ The below **figure 7.1** shows the Resource allocation graph which denotes,
- ✓ The sets P, R and E:
 
$$P = \{P_1, P_2, P_3\}$$

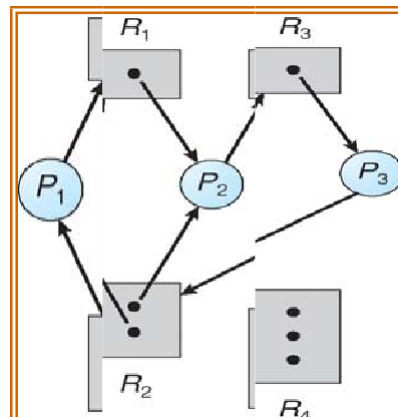
$$R = \{R_1, R_2, R_3, R_4\}$$

$$E = \{P_1 \rightarrow R_1, P_2 \rightarrow R_3, R_1 \rightarrow P_2, R_2 \rightarrow P_2, R_2 \rightarrow P_1, R_3 \rightarrow P_3\}$$
- ✓ Resource instances:
  - One instance of resource type  $R_1$
  - Two instances of resource type  $R_2$
  - One instance of resource type  $R_3$
  - Three instances of resource type  $R_4$
- ✓ Process states:

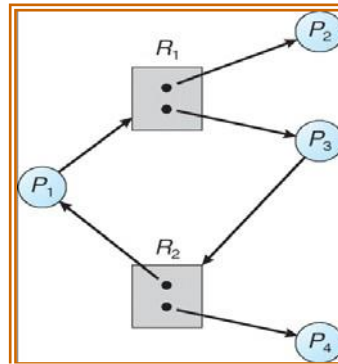
- Process  $P_1$  is holding an instance of resource type  $R_2$  and is waiting for an instance of resource type  $R_1$ .
- Process  $P_2$  is holding an instance of  $R_1$  and an instance of  $R_2$  and is waiting for an instance of  $R_3$ .
- Process  $P_3$  is holding an instance of  $R_3$ .



- ✓ If the graph contains no cycle, then no process in the system is deadlocked. If the graph contains a cycle then a deadlock may exist.
- ✓ If each resource type has exactly one instance then a cycle implies that a deadlock has occurred.
- ✓ If each resource type has several instances then a cycle does not necessarily imply that a deadlock has occurred.
- ✓ Consider the resource-allocation graph shown in **figure 7.1**.
- ✓ Suppose, process  $P_3$  requests an instance of resource type  $R_2$ . Since no resource instance is currently available, a request edge  $P_3 \rightarrow R_2$  is added to the graph which results in below **figure 7.2**.



- ✓ Now, two minimal cycles exist in the system:
  - $P_1 \rightarrow R_1 \rightarrow P_2 \rightarrow R_3 \rightarrow P_3 \rightarrow R_2 \rightarrow P_1$
  - $P_2 \rightarrow R_3 \rightarrow P_3 \rightarrow R_2 \rightarrow P_2$
- ✓ Processes  $P_1$ ,  $P_2$ , and  $P_3$  are deadlocked.
- ✓ Process  $P_2$  is waiting for the resource  $R_3$ , which is held by process  $P_3$ .
- ✓ Process  $P_3$  is waiting for either process  $P_1$  or process  $P_2$  to release resource  $R_2$ .
- ✓ In addition, process  $P_1$  is waiting for process  $P_2$  to release resource  $R_1$ .
- ✓ Consider the resource-allocation graph in below **figure 7.3**, which also has a cycle  $P_1 \rightarrow R_1 \rightarrow P_3 \rightarrow R_2 \rightarrow P_1$ , but there is no deadlock.



- ✓ P<sub>4</sub> may release its instance of resource type R<sub>2</sub>. That resource can then be allocated to P<sub>3</sub>, breaking the cycle.
- ✓ If a resource-allocation graph does not have a cycle, then the system is not in a deadlocked state. If there is a cycle, then the system may or may not be in a deadlocked state.

### 4.3 Methods for Handling Deadlocks

- ✓ **There are three ways to deal with deadlock problem:**
  - **Protocol can be implemented** to prevent or avoid deadlocks, ensuring that the system will never enter into the deadlock state.
  - Allow a system to enter into deadlock state, **detect it and recover** from it.
  - **Ignore** the problem and pretend that the deadlock never occur in the system. This is used by most OS including UNIX.
- ✓ To ensure that the deadlock never occurs, the system can use either deadlock avoidance or deadlock prevention.
- ✓ Deadlock prevention is a set of method for ensuring that at least one of the necessary conditions does not occur.
- ✓ Deadlock avoidance requires the OS is given advance information about which resource a process will request and use during its lifetime.
- ✓ If a system does not use either deadlock avoidance or deadlock prevention then a deadlock situation may occur. In this situation the system can provide an algorithm that examines the state of the system to determine whether a deadlock has occurred and an algorithm to recover from deadlock.
- ✓ Undetected deadlock will result in deterioration of the system performance.

### 4.4 Deadlock Prevention

- ✓ For a deadlock to occur each of the four necessary conditions must hold. If at least one of these conditions **does not hold** then we can prevent occurrence of deadlock.
  - **Mutual Exclusion:** This holds for non-sharable resources. For ex, A printer can be used by only one process at a time. Mutual exclusion is not possible in sharable resources and thus they cannot be involved in deadlock. Read-only files are good examples for sharable resources.
  - **Hold and Wait:** This condition can be eliminated by forcing a process to release all its resources held by it when it requests a resource. Two possible **solutions(protocols)** to achieve this are,
    - One protocol can be used is that each process is allocated with all of its resources before it starts execution.

- Another protocol that can be used is to allow a process to request a resource when the process has none.
- To illustrate the difference between these two protocols, we consider a process that copies data from a DVD drive to a file on disk, sorts the file, and then prints the results to a printer. If all resources must be requested at the beginning of the process, then the process must initially request the DVD drive, disk file, and printer. It will hold the printer for its entire execution, even though it needs the printer only at the end.
- The second method allows the process to request initially only the DVD drive and disk file. It copies from the DVD drive to the disk and then releases both the DVD drive and the disk file. The process must then again request the disk file and the printer. After copying the disk file to the printer, it releases these two resources and terminates.
- Both protocols have two main **disadvantages**.
  1. First, resource utilization is low, since resources may be allocated but unused for a long period.
  2. Second, starvation is possible. A process that needs several popular resources may have to wait indefinitely, because at least one of the resources that it needs is always allocated to some other process.
- **No Preemption:** To ensure that this condition never occurs the resources must be preempted. The following protocols can be used.
  - If a process is holding some resource and request another resource that cannot be immediately allocated to it, then all the resources currently held by the requesting process are preempted and added to the list of resources for which other processes may be waiting. The process will be restarted only when it regains the old resources and the new resources that it is requesting.
  - When a process request resources, we check whether they are available or not. If they are available we allocate them else we check that whether they are allocated to some other waiting process. If so we preempt the resources from the waiting process and allocate them to the requesting process. Otherwise, the requesting process must wait.
- **Circular Wait:** One way to ensure that this condition never holds is to impose total ordering of all resource types and each process requests resource in an increasing order. For ex, Let  $R = \{R_1, R_2, \dots, R_m\}$  be the set of resource types. We assign each resource type with a unique integer value. This allows us to compare two resources and determine whether one precedes the other in ordering. We can define a one to one function  $F: R \rightarrow N$  as follows,
 
$$F(\text{disk drive})=5, F(\text{printer})=12, F(\text{tape drive})=1$$

Deadlock can be prevented by using the following protocols.

- Each process can request the resource in **increasing order**. A process can request any number of instances of resource type say  $R_i$  and it can request instances of resource type  $R_j$  only  $F(R_j) > F(R_i)$ .
- Alternatively when a process requests an instance of resource type  $R_j$ , it has released any resource  $R_i$  such that  $F(R_i) \geq F(R_j)$ .

If these two protocols are used then the circular wait cannot hold.

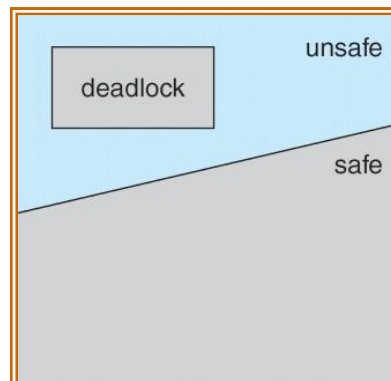
#### 4.5 Deadlock Avoidance

- ✓ Deadlock prevention algorithm may lead to low device utilization and reduces system throughput.
- ✓ Avoiding deadlocks requires additional information about how resources are to be requested. With the knowledge of the complete sequences of requests and releases we can decide for each requests whether the process should wait or not.

- ✓ For each requests it requires checking of the resources **currently available**, resources that are **currently allocated** to each processes, future **requests and release** of each process to decide whether the current requests can be satisfied or must wait to avoid future possible deadlock.
- ✓ A deadlock avoidance algorithm dynamically examines the resources allocation state to ensure that a circular wait condition never exists. The resource allocation state is defined by the number of available and allocated resources and the maximum demand of each process.

- **Safe State**

- ✓ A state is a **safe state** in which there exists at least one order in which all the process will run completely without resulting in a deadlock.
- ✓ A system is in safe state if there exist a **safe sequence**.
- ✓ A sequence of processes  $\langle P_1, P_2, \dots, P_n \rangle$  is a safe sequence for the current allocation state if, for each  $P_i$ , the resources requests that  $P_i$  can still make and it can be satisfied by the currently available resources plus the resources held by all  $P_j$ , with  $j < i$ .
- ✓ If the resources that  $P_i$  requests are not currently available then  $P_i$  can wait until all  $P_j$  have finished. When they have finished,  $P_i$  can obtain all of its needed resource to complete its designated task.
- ✓ A safe state is not a deadlocked state. But, a deadlocked state is an unsafe state.
- ✓ Not all unsafe states are deadlocked. An unsafe state may lead to a deadlock. It is shown in below **figure 7.4**.



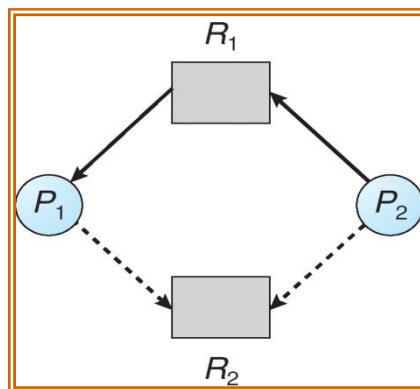
- ✓ **For Ex**, Consider a system with 12 magnetic tape drives and three processes  $P_0$ ,  $P_1$ , and  $P_2$ . Process  $P_0$  requires ten tape drives, process  $P_1$  may need as many as four tape drives, and process  $P_2$  may need up to nine tape drives. Suppose that, at time  $t_0$ , process  $P_0$  is holding five tape drives, process  $P_1$  is holding two tape drives, and process  $P_2$  is holding two tape drives. (there are 3 free tape drives.)

	Maximum Needs	Current Needs (allocated)
$P_0$	10	5
$P_1$	4	2
$P_2$	9	2

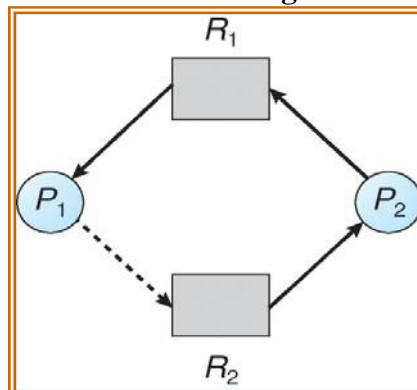
- ✓ At time  $t_0$ , the system is in a safe state. The sequence  $\langle P_1, P_0, P_2 \rangle$  satisfies the safety condition.
- ✓ A system can go from a safe state to an unsafe state. Suppose that, at time  $t_1$ , process  $P_2$  requests and is allocated one more tape drive. The system is no longer in a safe state. Only process  $P_1$  can be allocated all its tape drives.
- ✓ Whenever a process request a resource that is currently available, the system must decide whether resources can be allocated immediately or whether the process must wait. The request is granted only if the allocation leaves the system in safe state.

## Resource Allocation Graph Algorithm

- ✓ This algorithm is used only if we have **one instance** of a resource type. In addition to the request edge and the assignment edge a new edge called **claim edge** is used. A claim edge  $P_i \odot R_j$  indicates that process  $P_i$  may request  $R_j$  in future. The claim edge is represented by a **dotted line**.
- ✓ When a process  $P_i$  requests the resource  $R_j$ , the claim edge is converted to a request edge. When resource  $R_j$  is released by process  $P_i$ , the **assignment edge**  $R_j \odot P_i$  is replaced by the claim edge  $P_i \odot R_j$ .
- ✓ When a process  $P_i$  requests resource  $R_j$  the request is granted only if converting the request edge  $P_i \odot R_j$  to as assignment edge  $R_j \odot P_i$  do not result in a cycle.
- ✓ Cycle detection algorithm is used to detect the cycle. If there are no cycles then the allocation of the resource to process leave the system in safe state.
- ✓ To illustrate this algorithm, we consider the resource-allocation graph shown in below **figure 7.5**.



- ✓ Suppose that  $P_2$  requests  $R_2$  but we cannot allocate it to  $P_2$  even if  $R_2$  is currently free, because this will create a cycle in the graph as shown in below **figure 7.6**.



- ✓ A cycle indicates that the system is in an unsafe state. If  $P_1$  requests  $R_2$ , and  $P_2$  requests  $R_1$ , then a deadlock will occur.

### • Banker's Algorithm

- ✓ This algorithm is applicable to the system with **multiple instances** of each resource types, but this is less efficient than the resource allocation graph algorithm.
- ✓ When a new process enters the system it must declare the maximum number of resources that it may need. This number may not exceed the total number of resources in the system. The system must determine that whether the allocation of the resources will leave the system in a safe state or not. If it is so resources are allocated else it should wait until the process release enough resources.



- ✓ Several **data structures** are used to implement the banker's algorithm. Let 'n' be the number of processes in the system and 'm' be the number of resource types. The following data structures are needed.
  - **Available:** A vector of length m indicates the number of available resources. If  $\text{Available}[j]=k$ , then k instances of resource type  $R_j$  is available.
  - **Max:** An  $n \times m$  matrix defines the maximum demand of each process. If  $\text{Max}[i][j]=k$ , then  $P_i$  may request at most k instances of resource type  $R_j$ .
  - **Allocation:** An  $n \times m$  matrix defines the number of resources of each type currently allocated to each process. If  $\text{Allocation}[i][j]=k$ , then  $P_i$  is currently allocated k instances of resource type  $R_j$ .
  - **Need:** An  $n \times m$  matrix indicates the remaining resources need of each process. If  $\text{Need}[i][j]=k$ , then  $P_i$  may need k more instances of resource type  $R_j$  to complete its task. So  $\text{Need}[i][j] = \text{Max}[i][j] - \text{Allocation}[i][j]$ .

#### ▪ Safety Algorithm

- ✓ This algorithm is used to find out whether or not a system is in safe state or not. The algorithm can be described as follows,
  - Step 1.** Let Work and Finish be two vectors of length m and n respectively. Initialize  $\text{work} = \text{available}$  and  $\text{Finish}[i] = \text{false}$  for  $i=1,2,3,\dots,n$
  - Step 2.** Find i such that both
    - $\text{Finish}[i] == \text{false}$
    - $\text{Need}_i \leq \text{Work}$
 If no such i exists, then go to step 4
  - Step 3.**  $\text{Work} = \text{Work} + \text{Allocation}_i$   
 $\text{Finish}[i] = \text{true}$   
 Go to step 2
  - Step 4.** If  $\text{Finish}[i] == \text{true}$  for all i, then the system is in safe state.
- ✓ This algorithm may require an order of  $m \times n^2$  operation to decide whether a state is safe.

#### ▪ Resource Request Algorithm

- ✓ Let  $\text{Request}_i$  be the request vector of process  $P_i$ . If  $\text{Request}_i[j] = k$ , then process  $P_i$  wants k instances of the resource type  $R_j$ . When a request for resources is made by process  $P_i$  the following actions are taken.
  - If  $\text{Request}_i \leq \text{Need}_i$ , go to step 2. Otherwise raise an error condition, since the process has exceeded its maximum claim.
  - If  $\text{Request}_i \leq \text{Available}$ , go to step 3. Otherwise  $P_i$  must wait, since the resources are not available.
  - The system pretend to have allocated the requested resources to process  $P_i$ , then modify the state as follows.
    - $\text{Available} = \text{Available} - \text{Request}_i$
    - $\text{Allocation}_i = \text{Allocation}_i + \text{Request}_i$
    - $\text{Need}_i = \text{Need}_i - \text{Request}_i$
- ✓ If the resulting resource allocation state is safe, the transaction is complete and  $P_i$  is allocated its resources. If the new state is unsafe, then  $P_i$  must wait for  $\text{Request}_i$  and old resource allocation state is restored.

### ▪ An Illustrative Example

- ✓ To illustrate the use of the banker's algorithm, consider a system with five Processes  $P_0$  through  $P_4$  and three resource types A, B, and C. Resource type A has ten instances, resource type B has five instances, and resource type C has seven instances. Suppose that, at time  $T_0$ , the following snapshot of the system has been taken.

	Allocation			Max			Available		
	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			

- ✓ The content of the matrix **Need** is defined to be **Max - Allocation** and is as follows:

	Need		
	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

$P_0 \rightarrow 7\ 4\ 3 \leq 3\ 3\ 2$  is false,

$P_1 \rightarrow 1\ 2\ 2 \leq 3\ 3\ 2$  is true, so  $work = work + allocation$

$$work = 3\ 3\ 2 + 2\ 0\ 0 = 5\ 3\ 2$$

$P_2 \rightarrow 6\ 0\ 0 \leq 5\ 3\ 2$  is false,

$P_3 \rightarrow 0\ 1\ 1 \leq 5\ 3\ 2$  is true, so  $work = 5\ 3\ 2 + 2\ 1\ 1 = 7\ 4\ 3$

$P_4 \rightarrow 4\ 3\ 1 \leq 7\ 4\ 3$  is true, so  $work = 7\ 4\ 3 + 0\ 0\ 2 = 7\ 4\ 5$

$P_2 \rightarrow 6\ 0\ 0 \leq 7\ 4\ 5$  is true, so  $work = 7\ 4\ 5 + 3\ 0\ 2 = 10\ 4\ 7$

$P_0 \rightarrow 7\ 4\ 3 \leq 10\ 4\ 7$  is true, so  $work = 10\ 4\ 7 + 0\ 1\ 0 = 10\ 5\ 7$

- ✓ We claim that the system is currently in a safe state.
- ✓ The sequence  $\langle P_1, P_3, P_4, P_2, P_0 \rangle$  satisfies the safety criteria.
- ✓ Suppose now the process  $P_1$  requests one additional instance of resource type A and two instances of resource type C, so  $Request_1 = (1, 0, 2)$ .
- ✓ To decide whether this request can be immediately granted, we first check that **Request<sub>1</sub> ≤ Need<sub>1</sub>**, that is, **(1,0,2) ≤ (1,2,2)**, which is true then, **Request<sub>1</sub> ≤ Available<sub>1</sub>**, that is, **(1,0,2) ≤ (3,3,2)**, which is true. Then we arrive at the following new state:

	Allocation			Need			Available		
	A	B	C	A	B	C	A	B	C
$P_0$	0	1	0	7	4	3	2	3	0
$P_1$	3	0	2	0	2	0			
$P_2$	3	0	2	6	0	0			
$P_3$	2	1	1	0	1	1			

P<sub>4</sub>      0 0 2      4 3 1

- ✓ Now we must determine whether this new system state is safe. We execute safety algorithm and find that the sequence  $\langle P_1, P_3, P_4, P_0, P_2 \rangle$  satisfies the safety requirement. Hence the request can be immediately granted.
- ✓ From Resource Request Algorithm, we must also see that when the system is in this state, a request for (3,3,0) by P<sub>4</sub> cannot be granted, since the resources are not available,  
i.e.,  $(3,3,0) \leq (4,3,1) \dots \text{true}$   
 $(3,3,0) \leq (2,3,0) \dots \text{false}$
- ✓ Similarly, a request for (0,2,0) by P<sub>0</sub> cannot be granted, even though the resources are available, because the resulting state is unsafe.  
i.e.,  $(0,2,0) \leq (7,4,3) \dots \text{true}$   
 $(0,2,0) \leq (2,3,0) \dots \text{true}$
- ✓ Now the snapshot changes as follows,

	Allocation	Need	Available
	A B C	A B C	A B C
P <sub>0</sub>	0 3 0	7 2 3	2 1 0
P <sub>1</sub>	3 0 2	0 2 0	
P <sub>2</sub>	3 0 2	6 0 0	
P <sub>3</sub>	2 1 1	0 1 1	
P <sub>4</sub>	0 0 2	4 3 1	

- ✓ Then we are supposed to apply safety algorithm to this snapshot, but no safe sequence is generated, hence the request for (0,2,0) by P<sub>0</sub> cannot be granted.

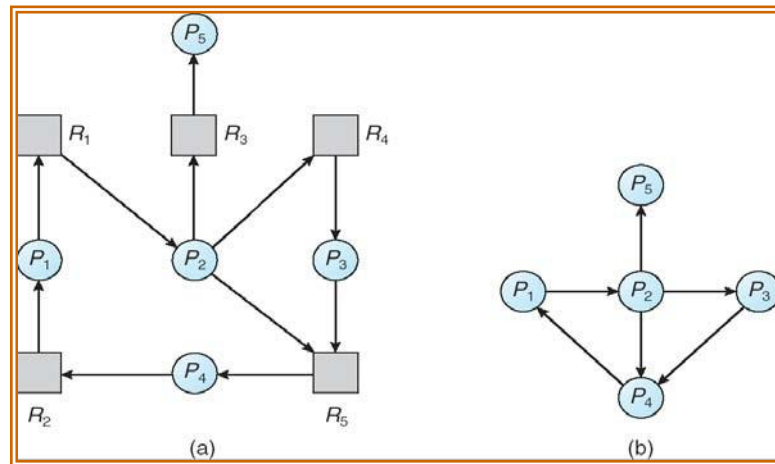
#### 4.6 Deadlock Detection

If a system does not employ either deadlock prevention or a deadlock avoidance algorithm then a deadlock situation may occur. In this environment the system must provide,

- An algorithm that examines the state of the system to determine whether a deadlock has occurred.
- An algorithm to recover from the deadlock.

- **Single Instances of each Resource Type**

- ✓ If all the resources have only a single instance then we can define deadlock detection algorithm that uses a **variant of resource allocation graph** as shown in below **figure 7.7(a)** called a **wait-for graph** as shown in below **figure 7.7(b)**. This graph is obtained by removing the resource nodes and collapsing appropriate edges.



- ✓ An edge from  $P_i$  to  $P_j$  in wait for graph implies that  $P_i$  is waiting for  $P_j$  to release a resource that  $P_i$  needs.
- ✓ An edge from  $P_i$  to  $P_j$  exists in wait for graph if and only if the corresponding resource allocation graph contains the edges  $P_i \odot R_q$  and  $R_q \odot P_j$ .
- ✓ Deadlock exists within the system if and only if there is a cycle. To detect deadlock the system needs an algorithm that searches for cycle in a graph.

#### • Several Instances of Resource Type

- ✓ The wait-for graph scheme is not applicable to a resource-allocation system with multiple instances of each resource type.
- ✓ The deadlock detection algorithm includes following time-varying data structures.
  - **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  $\text{Request}[i][j] = k$  then  $P_i$  is requesting  $k$  more instances of resources type  $R_j$ .
- ✓ The deadlock detection algorithm can be defined as follows,

**Step 1.** Let Work and Finish be vectors of length  $m$  and  $n$  respectively. Initialize  $\text{Work} = \text{Available}$ . For  $i = 0, 1, 2, \dots, n-1$ , if  $\text{allocation}_i \neq 0$  then  $\text{Finish}[i] = \text{false}$ , else  $\text{Finish}[i] = \text{true}$ .

**Step 2.** Find an index  $i$  such that both

$\text{Finish}[i] = \text{false}$

$\text{Request}_i \leq \text{Work}$

If no such  $i$  exists, go to step 4.

**Step 3.**  $\text{Work} = \text{Work} + \text{Allocation}_i$

$\text{Finish}[i] = \text{true}$

Go to step 2.

**Step 4.** If  $\text{Finish}[i] == \text{false}$ , for some  $i$  where  $0 \leq i < n$ , then a system is in a deadlock state.

- ✓ To illustrate this algorithm, we consider a system with five processes  $P_0$  through  $P_4$  and three resource types A, B, and C. Resource type A has seven instances, resource type B has two instances, and resource type C has six instances. Suppose that, at time  $T_0$ , we have the following resource-allocation state:

Allocation	Request	Available
------------	---------	-----------

	A B C	A B C	A B C
P <sub>0</sub>	0 1 0	0 0 0	0 0 0
P <sub>1</sub>	2 0 0	2 0 2	
P <sub>2</sub>	3 0 3	0 0 0	
P <sub>3</sub>	2 1 1	1 0 0	
P <sub>4</sub>	0 0 2	0 0 2	

- ✓ From the above algorithm, the sequence <P<sub>0</sub>, P<sub>2</sub>, P<sub>3</sub>, P<sub>1</sub>, P<sub>4</sub>> will result in Finish[i]== true for all i.
- ✓ If P<sub>2</sub> requests an additional instance of type C, the Request matrix is modified as follows,

	Request A B C
P <sub>0</sub>	0 0 0
P <sub>1</sub>	2 0 1
P <sub>2</sub>	0 0 1
P <sub>3</sub>	1 0 0
P <sub>4</sub>	0 0 2

- ✓ The system is now deadlocked. Even though we can reclaim resources held by process P<sub>0</sub>, but number of available resources is not sufficient to fulfill the requests of other processes. Thus, deadlock exists, consisting of processes P<sub>1</sub>, P<sub>2</sub>, P<sub>3</sub>, and P<sub>4</sub>.

- **Detection algorithm usage**

- ✓ This algorithm helps to find,
  - How **often** a deadlock is likely to occur?
  - How **many** processes will be affected by deadlock when it happens?
- ✓ We can invoke the deadlock detection algorithm when a request for allocation cannot be granted immediately.
- ✓ If detection algorithm is invoked for every resource request, this will cause a computational time overhead. So the algorithm must be invoked less frequently.

### 4.7 Recovery from Deadlock

There are **two** options for breaking a deadlock,

- One is to abort one or more processes to break the circular wait.
- The other is to preempt some resources from one or more of the deadlocked processes.

- **Process Termination**

- ✓ To eliminate deadlocks by aborting a process, **one of two methods** can be used. In both methods, the system reclaims all resources allocated to the terminated processes.
  - **Abort all deadlocked processes.** This method breaks the deadlock cycle, but at great expense.
  - **Abort one process at a time until the deadlock cycle is eliminated.** This method causes overhead, since after each process is aborted, a deadlock-detection algorithm must be invoked to determine whether any processes are still deadlocked.
- ✓ Aborting a process is not a easy task. If the process was in the middle of updating a file, terminating it will leave that file in an incorrect state.
- ✓ If the **partial termination method** is used, then we must determine which deadlocked process (or processes) should be terminated.

- ✓ We should abort those processes whose termination will incur the minimum cost. The following **factors** are considered to **select** the process.

1. What the priority of the process is?
2. How long the process has computed and how much longer the process will compute before completing its designated task?
3. How many and what types of resources the process has used?
4. How many more resources the process needs in order to complete?
5. How many processes will need to be terminated?
6. Whether the process is interactive or batch?

- **Resource Preemption**

- ✓ To eliminate deadlocks using resource preemption, we successively preempt some resources from processes and give these resources to other processes until the deadlock cycle is broken.
- ✓ The **three issues** need to be addressed are,
  - **Selecting a victim.** Which resources and which processes are to be preempted? We must determine the order of preemption to minimize cost.
  - **Rollback.** We must roll back the preempted process to some safe state and restart it from that state.
  - **Starvation.** We must ensure that a process can be picked as a victim only a small number of times. The most common solution is to include the number of rollbacks in the cost factor.

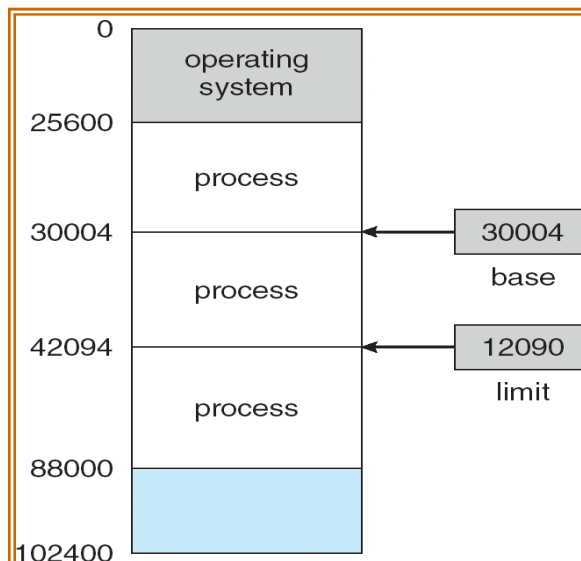
## MEMORY MANAGEMENT STRATEGIES

### ➤ Background

- ✓ Memory management is concerned with managing the primary memory.
- ✓ Memory consists of array of bytes or words each with its own address.
- ✓ We can ignore how a program generates a memory address. We are interested only in the sequence of memory addresses generated by the running program.

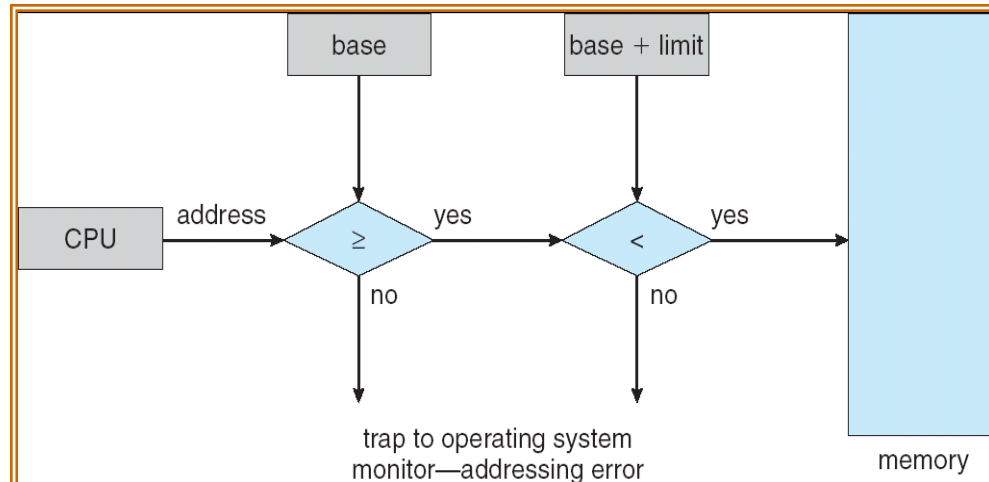
### • Basic Hardware

- ✓ **Main memory and the registers** in the processor are the only storage that the CPU can access directly. Hence the program and data must be brought from disk into main memory for CPU to access.
- ✓ Registers can be accessed in **one CPU** clock cycle. But main memory access can take **many CPU** clock cycles. Hence processor needs to stall, since it does not have the data required to complete the instruction that is executing.
- ✓ To overcome above situation a fast memory called **cache** is placed between main memory and CPU registers.
- ✓ We must ensure correct operation to **protect the operating system** from access by user processes and also to protect user processes from one another. This protection must be provided by the hardware. It can be implemented in several ways and **one such possible implementation** is,
  - We first need to make sure that each process has a separate memory space.
  - To do this, we need the ability to determine the range of legal addresses that the process may access and to ensure that the process can access only these legal addresses.
  - We can provide this protection by using two registers, **a base and a limit**, as illustrated in below **figure**.



- The **base register** holds the smallest legal **physical memory address**; the **limit register** specifies the size of **the range**. **For example**, if the base register holds 300040 and the limit register is 120900, then the program can legally access all addresses from 300040 through 420940.

- ✓ **Protection of memory space** is accomplished by having the CPU hardware compare every address generated in user mode with the registers.
- ✓ Any attempt by a program executing in user mode to access operating-system memory or other users' memory results in a **trap to the operating system**, which treats the attempt as a **fatal error** as shown in below **figure**.
- ✓ This prevents a user program from (accidentally or deliberately) modifying the code or data structures of either the operating system or other users.
- ✓ The **base and limit registers** uses a special privileged instructions which can be executed only in **kernel mode**, and since only the operating system executes in kernel mode, **only the operating system can load the base and limit registers**.

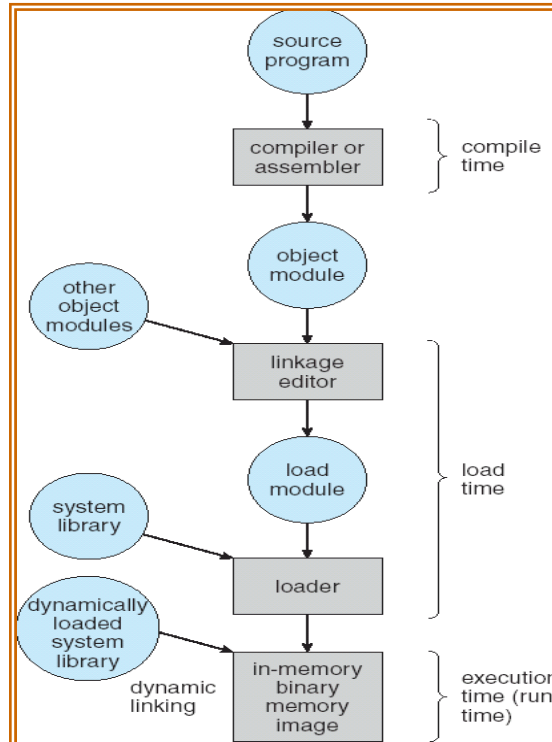


### • Address Binding

- ✓ Programs are stored on the secondary storage disks as binary executable files.
- ✓ When the programs are to be executed they are brought in to the main memory and placed within a process.
- ✓ The collection of processes on the disk waiting to enter the main memory forms the **input queue**.
- ✓ One of the processes which are to be executed is fetched from the queue and is loaded into main memory.
- ✓ During the execution it fetches instruction and data from main memory. After the process terminates it returns back the memory space.
- ✓ During execution the process will go through several steps as shown in below **figure** and in each step the address is represented in different ways.
- ✓ In source program the address is symbolic. The compiler **binds** the symbolic address to re-locatable address. The loader will in turn bind this re-locatable address to absolute address.
- ✓ Binding of instructions and data to memory addresses can be done at **any step** along the way:
  - **Compile time:** If we know at compile time where the process resides in memory, then **absolute code** can be generated. **For example**, if we know that a user process will reside starting at location *R*, then the generated compiler code will start at that location and extend up from there. If, at some later time, the starting location changes, then it will be necessary to recompile this code.

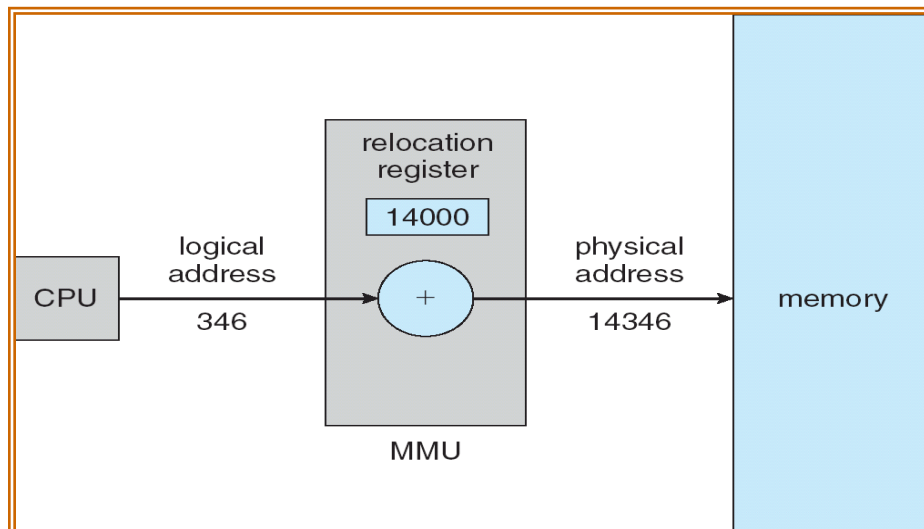


- **Load time:** If it is not known at compile time where the process will reside in memory, then the compiler must generate **relocatable code**. In this case, final binding is delayed **until load time**.
- **Execution time:** If the process is moved during its execution from one memory segment to another then the binding is delayed until run time. Special hardware is used for this. Most of the general purpose operating system uses this method.



### • Logical versus physical address

- ✓ The address generated by the CPU is called **logical address or virtual address**.
- ✓ The address seen by the memory unit i.e., the one loaded in to the memory register is called the **physical address**.
- ✓ Compile time and load time address binding methods generate **same logical and physical address**. The execution time addressing binding generate **different logical and physical address**.
- ✓ Set of logical address space generated by the programs is the **logical address space**. Set of physical address corresponding to these logical addresses is the **physical address space**.
- ✓ The **mapping** of virtual address to physical address during run time is done by the hardware device called **Memory Management Unit (MMU)**.
- ✓ The **base register** is now called **re-location register**.
- ✓ Value in the re-location register is added to every address generated by the user process at the time it is sent to memory as shown in below **figure 8.4**.
- ✓ **For example**, if the base is at **14000**, then an attempt by the user to address **location 0** is dynamically relocated to location 14000; an access to location **346** is mapped to location **14346**. The user program never sees the real physical addresses.



- **Dynamic Loading**

- ✓ For a process to be executed it should be loaded in to the physical memory. The size of the process is limited to the size of the physical memory. Dynamic loading is used to obtain better memory utilization.
- ✓ In dynamic loading the routine or procedure will not be loaded **until it is called**.
- ✓ Whenever a routine is called, the calling routine first checks whether the called routine is already loaded or not. If it is not loaded it calls the loader to load the desired program in to the memory and updates the programs address table to indicate the change and control is passed to newly **invoked or called** routine.
- ✓ The **advantages** are ,
  - Gives better memory utilization.
  - Unused routine is never loaded.
  - Do not need special operating system support.
  - Useful when large amount of codes are needed to handle infrequently occurring cases, such as error routines.

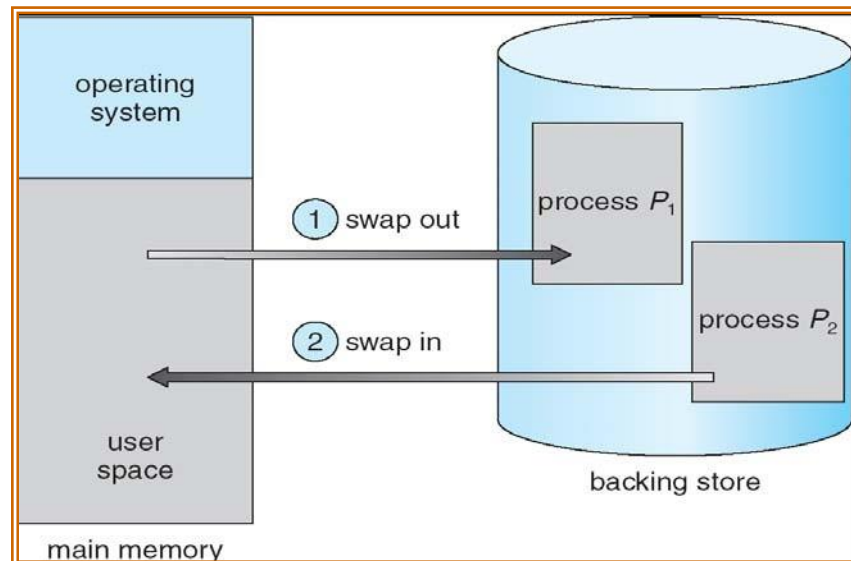
- **Dynamic linking and Shared libraries**

- ✓ Some operating system supports only the **static linking**.
- ✓ In dynamic linking only the main program is loaded in to the memory. If the main program requests a procedure, the procedure is loaded and the link is established at the time of references. This linking is postponed until the execution time.
- ✓ With dynamic linking a “**stub**” is used in the image of each library referenced routine. A “**stub**” is a **piece of code** which is used to **indicate how to locate the appropriate memory resident library routine or how to load library if the routine is not already present**.
- ✓ When “**stub**” is executed it checks whether the routine is present in memory or not. If not it loads the routine in to the memory.
- ✓ This feature can be used to update libraries i.e., library is replaced by a new version and all the programs can make use of this library.
- ✓ More than one version of the library can be loaded in memory at a time and each program uses its version of the library. Only the program that is compiled with the new version is affected by

the changes incorporated in it. Other programs linked before new version was installed will continue using older library. This system is called “**shared libraries**”.

## ➤ Swapping

- ✓ A process can be **swapped** temporarily out of the memory to a **backing store** and then brought back in to the memory for continuing the execution. This process is called swapping. **Ex.** In a multi-programming environment with a round robin CPU scheduling whenever the time quantum expires then the process that has just finished is swapped out and a new process swaps in to the memory for execution as shown in below **figure**.



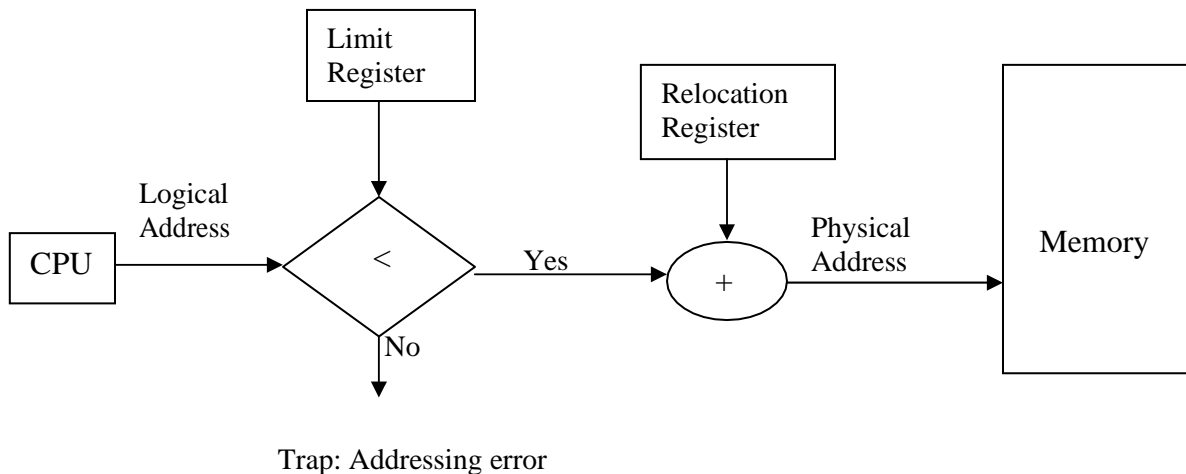
- ✓ A variant of this swapping policy is priority based scheduling. When a low priority is executing and if a high priority process arrives then a low priority will be swapped out and high priority is allowed for execution. This process is also called as **Roll out and Roll in**.
- ✓ Normally the process which is swapped out will be swapped back to the same memory space that is occupied previously and this depends upon address binding.
- ✓ The system maintains a **ready queue** consisting of all the processes whose memory images are on the backing store or in memory and are ready to run.
- ✓ The **context-switch time** in a swapping system is **high**. **For ex**, assume that the user process is 10 MB in size and the backing store is a standard hard disk with a transfer rate of 40MB per second. The actual transfer of the 40MB process to or from main memory takes,
 
$$\frac{10\text{MB}(10000\text{KB})}{40\text{MB}(40000\text{KB}) \text{ per second}} = \frac{1}{4} \text{ second} = 250 \text{ milliseconds}$$
- ✓ Assuming an **average latency** of **8 milliseconds**, the swap time is **258 milliseconds**. Since we must both swap out and swap in, the total swap time is about **516 milliseconds**.
- ✓ Swapping is constrained by other factors,
  - To swap a process, it should be completely idle.
  - If a process is waiting for an I/O operation, then the process cannot be swapped.

## ➤ Contiguous Memory Allocation

- ✓ The main memory must accommodate both the operating system and the various user processes. One common **method** to allocate main memory in the most efficient way is **contiguous memory allocation**.
- ✓ The memory is divided into two partitions, **one for the resident of operating system and one for the user processes**.

- **Memory mapping and protection**

- ✓ Relocation registers are used to protect user processes from each other, and to protect from changing OS code and data.
- ✓ The relocation register contains the value of the smallest physical address and the limit register contains the range of logical addresses.
- ✓ With relocation and limit registers, each logical address must be less than the limit register.
- ✓ The MMU maps the logical address dynamically by adding the value in the relocation register. This mapped address is sent to main memory as shown in below **figure**.
- ✓ The relocation-register scheme provides an effective way to allow the operating system's size to change dynamically.



- **Memory Allocation**

- ✓ One of the simplest methods for memory allocation is to divide memory into several **fixed partition**. Each partition contains exactly one process. The degree of multi-programming depends on the number of partitions.
- ✓ In **multiple partition method**, when a partition is free, process is selected from the input queue and is loaded into the free partition of memory. When process terminates, the memory partition becomes available for another process.
- ✓ The OS keeps a table indicating which part of the memory is free and is occupied.
- ✓ Initially, all memory is available for user processes and is considered one large block of available memory called a **hole**.
- ✓ When a process requests, the OS searches for a large hole for this process. If the hole is too large, it is **split** into two. One part is allocated to the requesting process and the other is returned to the set of holes.
- ✓ The set of holes are searched to determine which hole is best to allocate.
- ✓ **Dynamic storage allocation problem** is one which concerns about how to satisfy a request of size  $n$  from a list of free holes. There are three **strategies/solutions** to select a free hole,

- **First fit:** Allocates first hole that is big enough. This algorithm scans memory from the beginning and selects the first available block that is large enough to hold the process.
- **Best fit:** It chooses the hole i.e., closest in size to the request. It allocates the smallest hole i.e., big enough to hold the process.
- **Worst fit:** It allocates the largest hole to the process request. It searches for the largest hole in the entire list.
- ✓ First fit and best fit are the most popular algorithms for dynamic memory allocation. **All these algorithms suffer from fragmentation.**

- **Fragmentation**

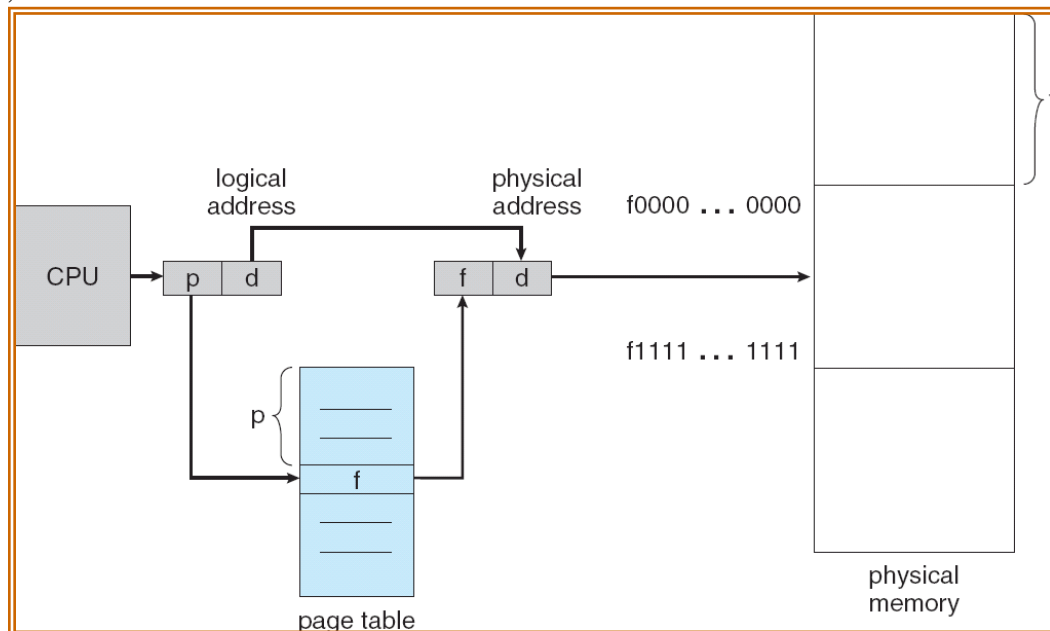
- ✓ **External Fragmentation** exists when there is enough memory space exists to satisfy the request, but it is not contiguous. Storage is fragmented into a large number of small holes.
- ✓ External Fragmentation may be either minor or a major problem.
- ✓ Statistical analysis of first fit reveals that, even with some optimization, given  $N$  allocated blocks, another  $0.5 N$  blocks will be lost to fragmentation. That is, one-third of memory may be unusable. This property is known as the **50-percent rule**.
- ✓ **Internal fragmentation:** Memory fragmentation can be internal as well as external. Consider a multiple-partition allocation scheme with a hole of 18,464 bytes. Suppose that the next process requests 18,462 bytes. If we allocate exactly the requested block, we are left with a hole of 2 bytes. The difference between these two numbers is **internal fragmentation** that is internal to a partition.
- ✓ The overhead to keep track of this hole will be substantially larger than the hole itself.
- ✓ The general approach to avoid this problem is to break the physical memory into fixed-sized blocks and allocate memory in units based on block size. With this approach, the memory allocated to a process may be slightly larger than the requested memory.
- ✓ One solution to over-come external fragmentation is **compaction**. The goal is to move all the free memory together to form a large block. Compaction is possible only if the re-location is dynamic and done at execution time.
- ✓ **Another solution** to the external fragmentation problem is to permit the logical address space of a process to be non-contiguous.

➤ **Paging**

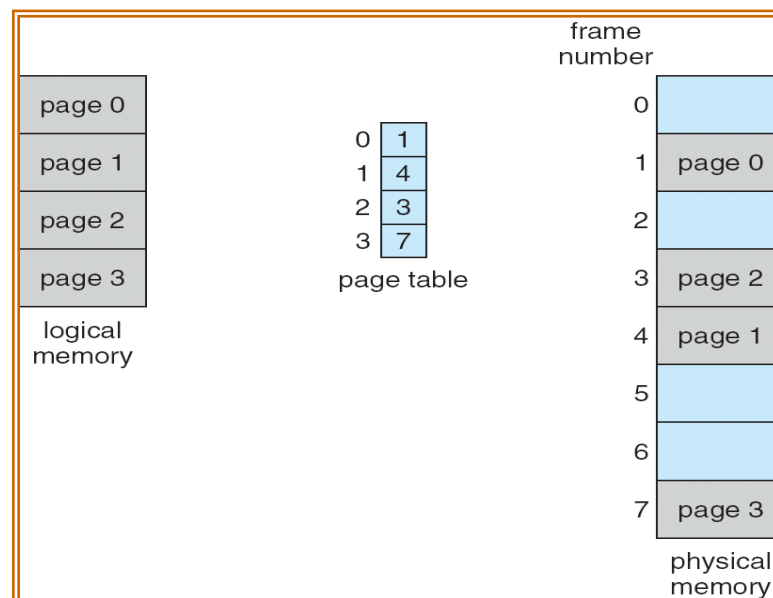
- ✓ Paging is a **memory management scheme** that permits the physical address space of a process to be **non-contiguous**. Support for paging is handled by **hardware**.
- ✓ Paging avoids the considerable problem of fitting the varying sized memory chunks on to the backing store.

- **Basic Method**

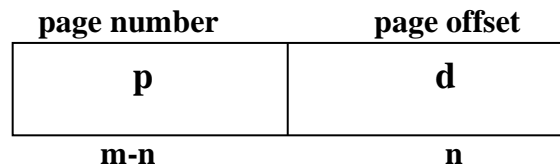
- ✓ Physical memory is broken in to fixed sized blocks called **frames (f)** and Logical memory(secondary memory) is broken in to blocks of same size called **pages (p)**.
- ✓ When a process is to be executed its pages are loaded in to available frames from the backing store. The backing store is also divided in to **fixed-sized blocks of same size as memory frames**.
- ✓ The below **figure** shows paging hardware.



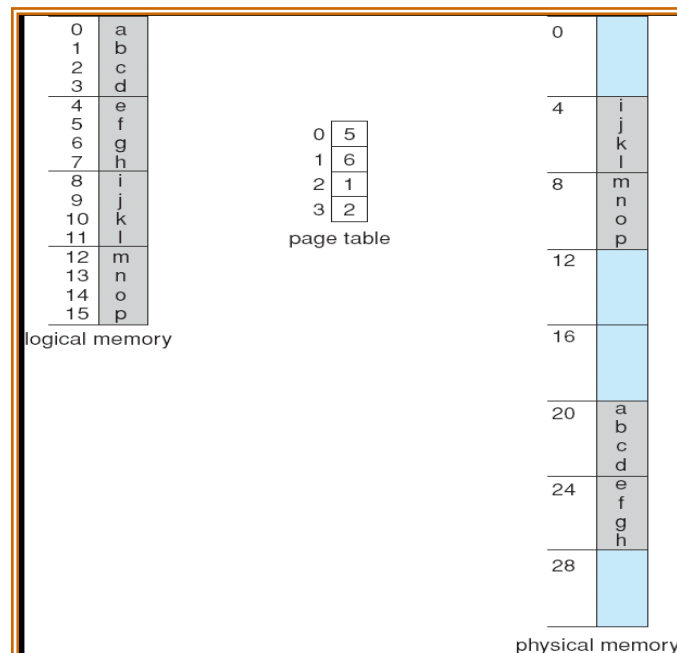
- ✓ Logical address generated by the CPU is divided in to **two parts: page number (p) and page offset (d)**.
- ✓ The page number (p) is used as **index** to the page table. The page table contains base address of each page in physical memory. This base address is combined with the page offset to define the physical memory i.e., sent to the memory unit. The paging model memory is shown in below **figure**.



- ✓ The page size is defined by the hardware. The size is the power of 2, varying between **512 bytes and 16Mb per page**.
- ✓ If the size of logical address space is  $2^m$  address unit and page size is  $2^n$ , then high order **m-n** designates the **page number** and **n** low order bits represents **page offset**. Thus logic address is as follows.

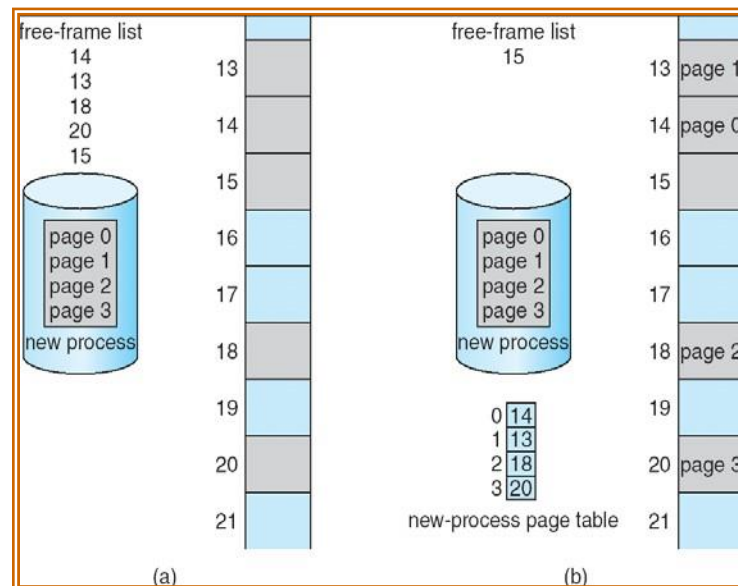


- ✓ **Ex:** To show how to map logical memory in to physical memory, consider a page size of 4 bytes and physical memory of 32 bytes (8 pages) as shown in below **figure 8.9**.
- Logical address 0** is page 0 and offset 0 and Page 0 is in frame 5. The **logical address 0** maps to physical address  $[(5*4) + 0]=20$ .
  - Logical address 3** is page 0 and offset 3 and Page 0 is in frame 5. The **logical address 3** maps to **physical address**  $[(5*4) + 3]= 23$ .
  - Logical address 4** is page 1 and offset 0 and page 1 is mapped to frame 6. So logical address 4 maps to **physical address**  $[(6*4) + 0]=24$ .
  - Logical address 13** is page 3 and offset 1 and page 3 is mapped to frame 2. So logical address 13 maps to **physical address**  $[(2*4) + 1]=9$ .



- ✓ In paging scheme, we have **no external fragmentation**. Any free frame can be allocated to a process that needs it. But we may have some **internal fragmentation**.
- ✓ If the memory requirements of a process do not happen to coincide with page boundaries, the last frame allocated may not be completely full.
- ✓ **For example**, if page size is 2,048 bytes, a process of 72,766 bytes will need 35 pages plus 1,086 bytes. It will be allocated 36 frames, resulting in **internal fragmentation** of  $2,048 - 1,086 = 962$  bytes.
- ✓ When a process arrives in the system to be executed, its size expressed in pages is examined. Each page of the process needs one frame. Thus, if the process requires **n** pages, at least **n** frames must be available in memory. If **n** frames are available, they are allocated to this arriving process.
- ✓ The first page of the process is loaded in to one of the allocated frames, and the frame number is put in the page table for this process. The next page is loaded into another frame and its frame

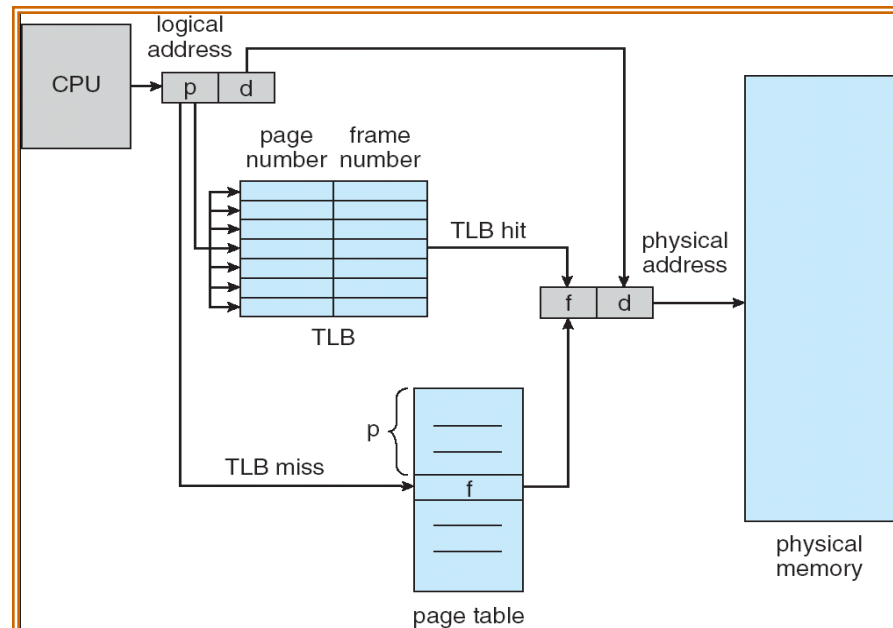
number is put into the page table and so on, as shown in below **figure (a) before allocation, (b) after allocation.**



## • Hardware Support

- ✓ The **hardware implementation** of the page table can be done in several ways. The simplest method is that the page table is implemented as a set of **dedicated registers**.
- ✓ The use of registers for the page table is satisfactory if the page table is reasonably small (for example, 256 entries). But most computers, allow the page table to be very large (for example, 1 million entries) and for these machines, the use of fast registers to implement the page table is not feasible.
- ✓ So the page table is kept in the main memory and a **page table base register (PTBR)** points to the page table and **page table length register (PTLR)** indicates size of page table. Here two memory accesses are needed to access a byte and thus memory access is slowed by a factor of 2.
- ✓ The only solution is to use a special, fast lookup hardware **cache** called **Translation look aside buffer (TLB)**. TLB is associative, with high speed memory. Each entry in TLB contains **two** parts, **a key and a value**. When an associative register is presented with an item, it is compared with all the key values, if found the corresponding value field is returned. Searching is fast but hardware is expensive.
- ✓ TLB is used with the page table **as follows**:
  - TLB contains only few page table entries.
  - When a logical address is generated by the CPU, its page number is presented to TLB. If the page number is found its frame number is immediately available and is used to access the actual memory. If the page number is not in the TLB (**TLB miss**) the memory reference to the page table must be made.
  - When the frame number is obtained we can use it to access the memory as shown in below **figure**. The page number and frame number are added to the TLB, so that they will be found quickly on the next reference.
  - If the TLB is full of entries the OS must select anyone for **replacement**.
  - Some TLBs allow entries to be **wired down** meaning that they cannot be removed from the TLB.





- ✓ Some TLBs store **Address Space Identifiers (ASIDs)** in each TLB entry. An ASID uniquely identifies each process and is used to provide address-space protection for that process.
- ✓ The percentage of time that a page number is found in the TLB is called **hit ratio**.
- ✓ **For example**, an **80-percent hit ratio** means that we find the desired page number in the TLB 80 percent of the time. If it takes **20 nanoseconds** to search the TLB and **100 nanoseconds to access memory**, then a mapped-memory access takes **120 nanoseconds when the page number is in the TLB**. If we fail to find the page number in the TLB (20 nanoseconds), then we must first access memory for the page table and frame number (100 nanoseconds) and then access the desired byte in memory (100 nanoseconds), for a total of **220 nanoseconds**. Thus the effective access time is,

$$\text{Effective Access Time (EAT)} = 0.80 \times 120 + 0.20 \times 220 \\ = 140 \text{ nanoseconds.}$$

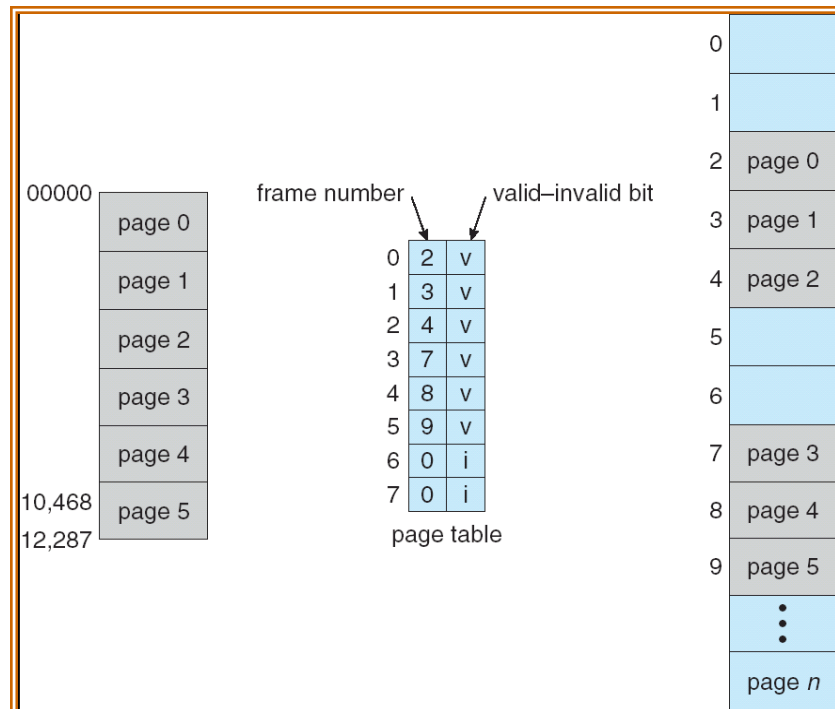
In this example, we suffer a **40-percent slowdown** in memory-access time (from 100 to 140 nanoseconds).

- ✓ For a **98-percent hit ratio** we have
- $$\text{Effective Access Time (EAT)} = 0.98 \times 120 + 0.02 \times 220 \\ = 122 \text{ nanoseconds.}$$
- ✓ This increased hit rate produces only a **22 percent slowdown** in access time.

## • Protection

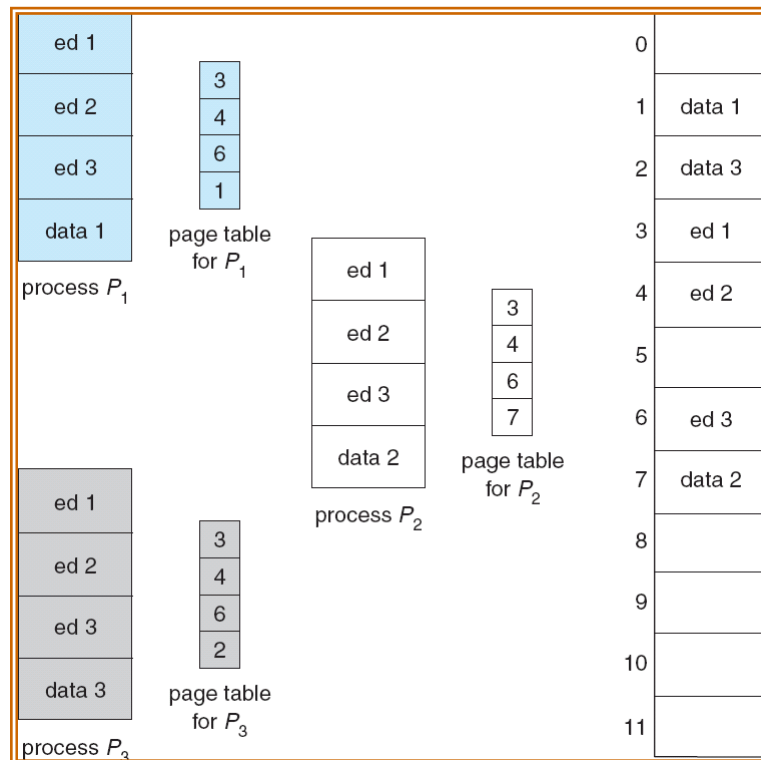
- ✓ Memory protection in paged environment is done by **protection bits** that are associated with each frame. These bits are kept in page table.
- ✓ One bit can define a page to be read-write or read-only.
- ✓ One more bit is attached to each entry in the page table, a **valid-invalid** bit.
- ✓ A valid bit indicates that associated page is in the process's logical address space and thus it is a legal or valid page.
- ✓ If the bit is invalid, it indicates the page is not in the process's logical address space and is illegal. Illegal addresses are trapped by using the valid-invalid bit.
- ✓ The OS sets this bit for each page to allow or disallow accesses to that page.

- ✓ **For example**, in a system with a 14-bit address space (0 to 16383), we have a program that should use only addresses 0 to 10468. Given a page size of 2 KB, we have the situation shown in below **figure**. Addresses in pages 0, 1, 2, 3, 4, and 5 are mapped normally through the page table. Any attempt to generate an address in pages 6 or 7, we find that the valid-invalid bit is set to **invalid**, and the computer will trap to the operating system (**invalid page reference**).



### • Shared Pages

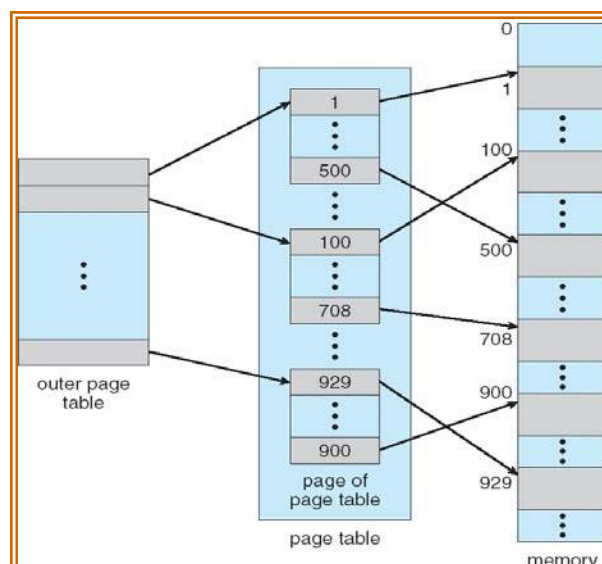
- ✓ An advantage of paging is the possibility of **sharing common code**. This consideration is particularly important in a time-sharing environment.
- ✓ Consider a system that supports 40 users, each of whom executes a text editor. If the text editor consists of 150 KB of code and 50 KB of data space, we need 8,000 KB to support 40 users. If the code is **reentrant code (or pure code)** it can be shared as shown in below **figure**. Here there are three-page editor-each page 50 KB in size and are being shared among three processes. Each process has its own data page.
- ✓ Reentrant code is **non-self-modifying code (Read only)**. It never changes during execution. Thus, two or more processes can execute the same code at the same time.
- ✓ Each process has its own copy of registers and data storage to hold the data for the process's execution.
- ✓ Thus, to support 40 users, we need only one copy of the editor (150 KB), plus 40 copies of the 50 KB of data space per user. The total space required is now 2150 KB instead of 8,000 KB. (**i.e.,  $150+40*50=2150$  KB**).



## ➤ Structure of the Page Table

### • Hierarchical paging

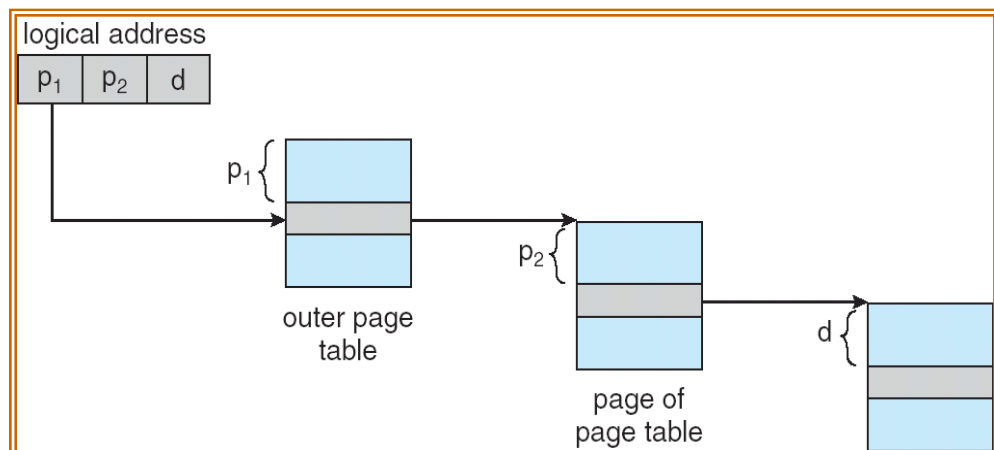
- ✓ Recent computer system support a large logical address space from  $2^{32}$  to  $2^{64}$  and thus page table becomes large. So it is very difficult to allocate contiguous main memory for page table. One **simple solution** to this problem is to divide page table in to smaller pieces.
- ✓ One way is to use **two-level paging algorithm** in which the page table itself is also paged as shown in below **figure**.



- ✓ **Ex.** In a 32-bit machine with page size of 4kb, a logical address is divided into a page number consisting of 20 bits and a page offset of 12 bit (page size  $4\text{kb} = 2^{12}$ ). The page table is further divided since the page table is paged, the page number is further divided into 10 bit page number and a 10 bit offset. So the logical address is,

Page number		page offset
P <sub>1</sub>	P <sub>2</sub>	d
10	10	12

- ✓ P<sub>1</sub> is an index into the outer page table and P<sub>2</sub> is the displacement within the page of the outer page table. The **address-translation method** for this architecture is shown in below **figure 8.15**. Because address translation works from the outer page table inward, this scheme is also known as a **forward-mapped page table**.



- ✓ For a system with a 64-bit logical address space, a two-level paging scheme is no longer appropriate. Suppose the page size in such a system is 4 KB the page table consists of up to  $2^{52}$  entries. If we use a two-level paging scheme, then the inner page tables can be one page long, or contain  $2^{10}$  4-byte entries. The addresses look like this,

Outer page	inner page	offset
P <sub>1</sub>	P <sub>2</sub>	d
42	10	12

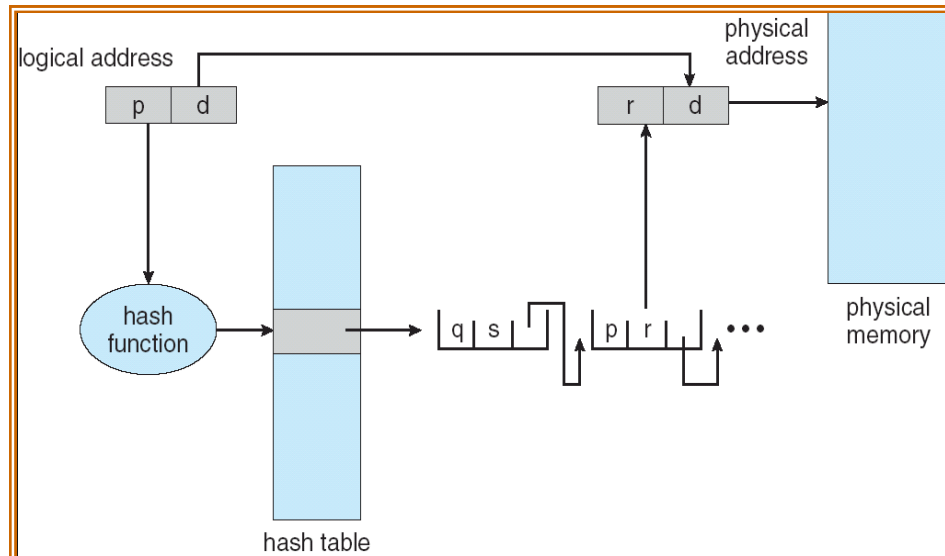
- ✓ The outer page table consists of  $2^{42}$  entries, or  $2^{44}$  bytes. The one way to avoid such a large table is to divide the outer page table into smaller pieces.
- ✓ We can avoid such a large table using **three-level paging scheme**.

2 <sup>nd</sup> outer page	outer page	inner page	offset
P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>	d
32	10	10	12

- ✓ The outer page table is still  $2^{34}$  bytes in size. The next step would be a **four-level paging scheme**.

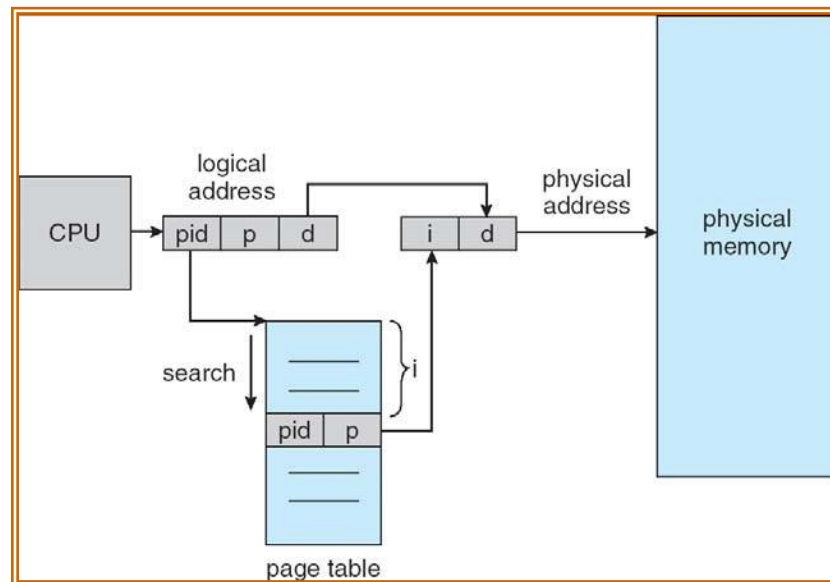
- **Hashed page table**

- ✓ Hashed page table handles the address space larger than 32 bit. The virtual page number is used as **hash value**. Linked list is used in the hash table which contains a list of elements that hash to the same location.
- ✓ Each element in the hash table contains the following three fields,
  - **Virtual page number**
  - **Mapped page frame value**
  - **Pointer to the next element in the linked list**
- ✓ The algorithm works as follows,
  - Virtual page number is taken from virtual address and is hashed in to the hash table.
  - Virtual page number is compared with **field 1** in the first element in the linked list.
  - If there is a match, the corresponding page frame (**field 2**) is used to form the desired physical address. If there is no match, subsequent entries in the linked list are searched for a matching virtual page number. This scheme is shown in below **figure**.
  - **Clustered pages** are similar to hash table but one difference is that each entity in the hash table refer to several pages.



- **Inverted Page Tables**

- ✓ Page tables may consume large amount of physical memory just to keep track of how other physical memory is being used.
- ✓ To solve this problem, we can use an inverted page table has one entry for each real page (or frame) of memory. Each entry consists of the virtual address of the page stored in that real memory location with information about the process that owns the page.
- ✓ Thus, only one page table is in the system, and it has only one entry for each page of physical memory.
- ✓ The below **figure** shows the operation of an inverted page table.

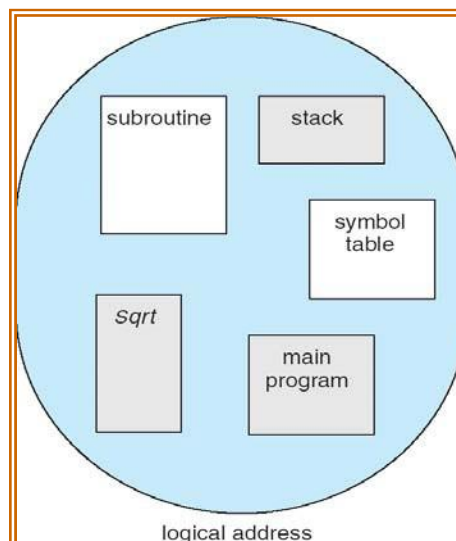


- ✓ The inverted page table entry is a pair **<process-id, page number>**. Where **process-id assumes the role of the address-space identifier**. When a memory reference is made, the part of virtual address consisting of **<process-id, page number>** is presented to memory sub-system.
- ✓ The inverted page table is searched for a match. If a match is found at **entry i**, then the physical address **<i, offset>** is generated. If no match is found then an illegal address access has been attempted.
- ✓ This scheme **decreases the amount of memory** needed to store each page table, but increases the amount of time needed to search the table when a page reference occurs.

## ➤ Segmentation

### • Basic method

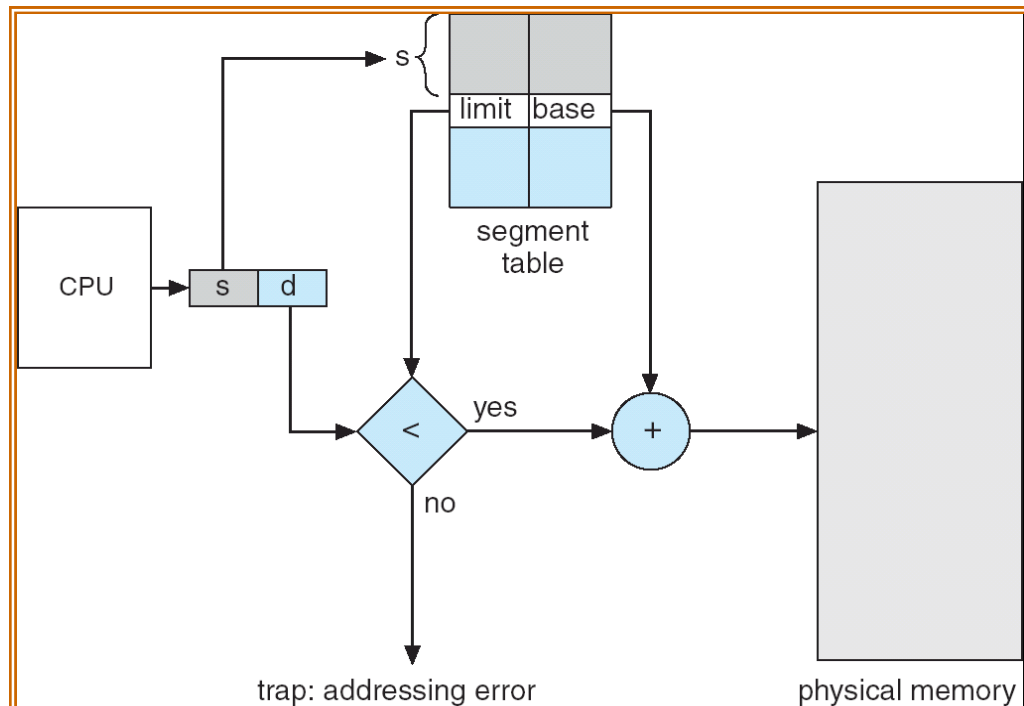
- ✓ Users prefer to view memory as a collection of **variable-sized segments, with no ordering among segments** as shown in below **figure**.



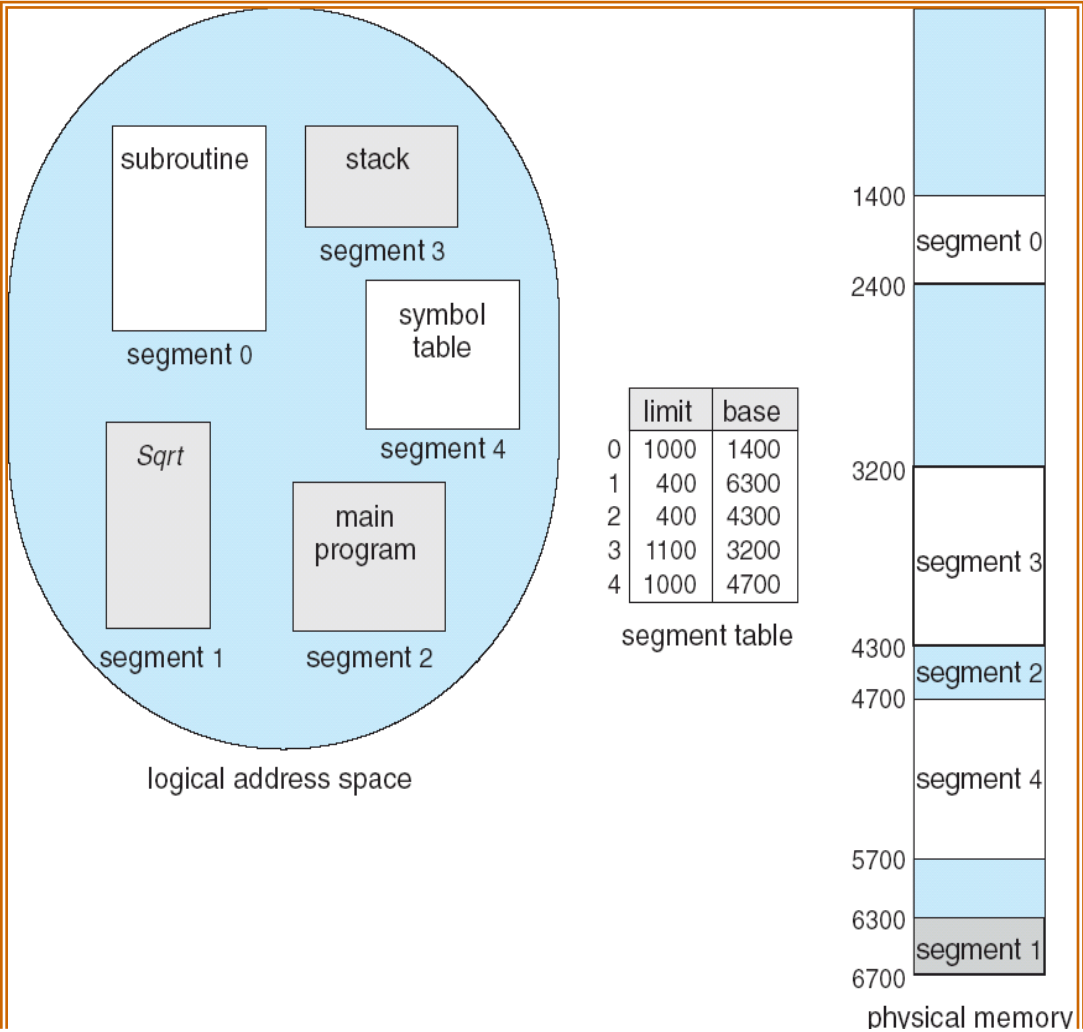
- ✓ **Segmentation** is a memory-management scheme that supports the user view of memory.
- ✓ A **logical address** is a collection of segments. Each segment has a **name and length**. The address specifies both the **segment name and the offset** within the segments.
- ✓ The segments are numbered and are referred by a segment number. So the logical address consists of **<segment number, offset>**.

- **Hardware**

- ✓ **Segment table** maps **2-Dimensional user defined address** in to **1-Dimensional physical address**.
- ✓ Each entry in the segment table has a **segment base and segment limit**.
- ✓ The segment **base contains the starting physical address where the segment resides and limit specifies the length of the segment**.
- ✓ The use of segment table is shown in the below **figure**.



- ✓ Logical address consists of two parts, **segment number s and an offset d**.
- ✓ The segment number is used as an **index** to segment table. The offset must be in between **0 and limit**, if not an error is reported to OS.
- ✓ If legal the offset is added to the base to generate the actual physical address.
- ✓ The segment table is an **array of base-limit register pairs**.
- ✓ **For example**, consider the below **figure**. We have **five segments** numbered from 0 through 4. Segment 2 is 400 bytes long and begins at location 4300. Thus, a reference to byte 53 of segment 2 is mapped onto location  $4300 + 53 = 4353$ . A reference byte 852 of segment 3, is mapped to 3200 (the base of segment 3) + 852 = 4052. A reference to byte 1222 of segment 0 would result in a trap to the operating system, as this segment is only 1,000 bytes long.





## **QUESTION BANK**

### **Process Synchronization**

1. What are semaphores? Explain two primitive semaphore operations. What are its advantages?
2. Explain any one synchronization problem for testing newly proposed sync scheme.
3. Explain three requirements that a solution to critical –section problem must satisfy.
4. State dining philosopher's problem and give a solution using semaphores. Write structure of philosopher i.
5. What do you mean by binary semaphore and counting semaphore? With C struct, explain implementation of wait() and signal().
6. Describe term monitor. Explain solution to dining philosopher's problem using monitor.
7. What are semaphores? Explain its usage and implementation.
8. What are monitors? Solve Dining philosophers problem by using monitors
9. Define spin lock.
10. Explain synchronization problem of bounded buffer and reader writers and give a solution code by using semaphores.
11. Explain Peterson's solution to solve critical section problem.
12. Give the need of priority inheritance protocol.
13. Explain synchronization hardware and mutex lock implementation with sample code.
14. What are monitors? Explain.
15. Give a solution for producer and consumer problem using monitors

### **File Implemetation**

1. What is a file? Describe different access methods on files.
2. What is file mounting? Explain.
3. Draw neat diagram and explain fixed file allocation. Is FAT linked allocation?
4. Explain following: file types, file operations, file attributes.
5. Explain methods to implement directories .
6. What is free space list? With example, explain any two methods to implement free space list.
7. What are major methods to allocate disk space? Explain each with examples.
8. Explain different file access methods.
9. Explain various directory structures.
10. Explain different disk space allocation methods with example.
11. Explain the various storage mechanisms available to store files, with neat diagram.

## Deadlocks:

1. Explain with an example how resource allocation graph is used to define deadlock?
2. What are two options for breaking deadlock?
3. What is wait-for graph? How is it useful for detection of deadlock?

4.

	Allocation	Request	Available
	1. ABC	ABC	ABC
ii. P <sub>0</sub>	0 1 0	0 0 0	0 0 0
iii. P <sub>1</sub>	2 0 0	2 0 2	
iv. P <sub>2</sub>	3 0 3	0 0 0	
v. P <sub>3</sub>	2 1 1	1 0 0	
vi. P <sub>4</sub>	0 0 2	0 0 2	

Show the system is not deadlocked by one safe sequence. At t<sub>2</sub>, p<sub>2</sub> makes one additional request for type C, show that system is deadlocked if request is granted.

5. Define hardware instructions test() , set() and swap(). Give algorithms to implement mutual exclusion with these instructions.
6. Describe necessary conditions for a deadlock situation to arise and how to handle.
7. List any 4 examples of deadlock that are not related to computer systems.
8. Explain the safety algorithm used in banker's algorithm, with suitable data structures.
9. List any 4 examples of deadlock that are not related to computer systems.
10. Explain the safety algorithm used in banker's algorithm, with suitable data structures.
11. Explain banker's algorithm for deadlock avoidance.
12. Explain process termination and resource preemption to recover from deadlock.
13. Consider given chart and answer i) what is content of matrix need? ii) is system safe? If yes give safe sequence. iii) if request comes from P<sub>1</sub>, arrives for ( 0,4,2,0), can it be granted ?

	Allocation	Max	Available
	ABCD	ABCD	ABCD
a. P <sub>0</sub>	0 0 1 2	0 0 1 2	1 5 2 0
b. P <sub>1</sub>	1 0 0 0	1 7 5 0	
c. P <sub>2</sub>	1 3 5 4	2 3 5 6	
d. P <sub>3</sub>	0 6 3 2	0 6 5 2	
e. P <sub>4</sub>	0 0 1 4	0 6 5 6	

The content of the matrix **Need** is defined to be **Max - Allocation** and is as follows:

### i. Need

ii. **A B C D**

b. $P_0$	0 0 0 0
c. $P_1$	0 7 5 0
d. $P_2$	1 0 0 2
e. $P_3$	0 0 2 0
f. $P_4$	0 6 4 2

ii.  $P_0 \rightarrow 0000 \leq 1520$  is true, so  $work = work + allocation$

$$a. work = 1520 + 0012 = 1532$$

iii.  $P_1 \rightarrow 0750 \leq 1532$  is false,

iv.  $P_2 \rightarrow 1002 \leq 1532$  is true, so  $work = 1532 + 1354 = 2886$

v.  $P_3 \rightarrow 0020 \leq 2886$  is true, so  $work = 2886 + 0632 = 3518$

vi.  $P_4 \rightarrow 0642 \leq 3518$  is true, so  $work = 3518 + 0014 = 3532$

vii.  $P_1 \rightarrow 0750 \leq 3532$  is true, so  $work = 3532 + 1000 = 4532$

We claim that the system is currently in a safe state. The sequence  $\langle P_0, P_2, P_3, P_4, P_1 \rangle$  satisfies the safety criteria. Suppose now the process  $P_1$  requests for **(0,4,2,0)**. To decide whether this request can be immediately granted, we first check that

f. **Request<sub>1</sub>**  $\leq$  **Need<sub>1</sub>**, that is, **(0,4,2,0)**  $\leq$  **(0,7,5,0)** which is true then,

g. **Request<sub>1</sub>**  $\leq$  **Available<sub>1</sub>**, that is, **(0,4,2,0)**  $\leq$  **(1,5,2,0)**, which is true. Then we arrive at the following new state:

1. **Allocation**    **Max**    **Available**

2. **ABCD**        **ABCD**    **ABCD**

ii. $P_0$	0 0 1 2	0 0 1 2	1 1 0 0
iii. $P_1$	1 4 2 0	1 7 5 0	
iv. $P_2$	1 3 5 4	2 3 5 6	
v. $P_3$	0 6 3 2	0 6 5 2	
vi. $P_4$	0 0 1 4	0 6 5 6	

i. **Need**

ii. **A B C D**

b. $P_0$	0 0 0 0
c. $P_1$	0 3 3 0
d. $P_2$	1 0 0 2
e. $P_3$	0 0 2 0
f. $P_4$	0 6 4 2

Now we must determine whether this new system state is safe. We execute safety algorithm and find that the sequence  $\langle P_1, P_3, P_4, P_0, P_2 \rangle$  satisfies the safety requirement. Hence the request can be immediately granted.

14. Consider a system consisting of  $m$  resources of the same type being shared by  $n$  processes. A process can request or release only one resource at a time. Show that the system is deadlock free if the following two conditions hold:
- The maximum need of each process is between one resource and  $m$  resources.
  - The sum of all maximum needs is less than  $m + n$ .

**Solution:**

The given conditions can be written as

$$\text{Max}_i \geq 1 \text{ for all } i$$

$n$

$$\sum_{i=1}^n \text{Max}_i < m + n$$

The need value can be calculated as,  $\text{Need}_i = \text{Max}_i - \text{Allocation}_i$

If there exists a deadlock then,

$n$

$$\sum_{i=1}^n \text{Allocation}_i = m$$

Therefore,  $\sum \text{Need}_i + \sum \text{Allocation}_i = \sum \text{Max}_i < m + n$

We get,  $\sum \text{Need}_i + m < m + n$

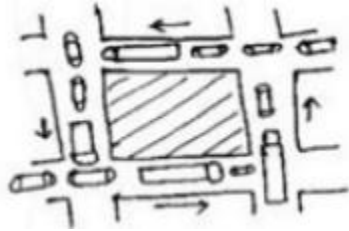
Hence,  $\sum \text{Need}_i < n$

This implies there exist  $P_i$  such that  $\text{Need}_i = 0$ . Since  $\text{Max}_i \geq 1$ ,  $P_i$  has atleast one resource to release. Hence no deadlocks.

**For ex,** consider  $m=5$ ,  $n=3$  and below snapshot then we can trace out with above steps and prove that no deadlock occurs

	Max	Allocation	Available	Need
$P_0$	2	1	0	1
$P_1$	3	2		1
$P_2$	2	2		0

15. Define deadlock. Consider the traffic deadlock depicted in the figure given below, Analyze the four necessary conditions for deadlock indeed hold in the example below.



16. What is wait for graph? Explain Deadlock detection algorithm for several instances of resource type.

### **Memory Management:**

1. What do you mean by fragmentation? Explain difference between internal and external fragmentation.
2. Differentiate between internal and external fragmentation. How are they overcome?
3. What is paging and swapping?
4. What is address binding? Explain with necessary steps, binding instructions and data to memory addresses.
5. Mention the problem with simple paging scheme. How TLB is used to solve this problem? Explain with supporting h/w dig with example.
6. Draw and explain the multistep processing of a user program.
7. In the paging scheme with TLB it takes 20 ns to search the TLB and 100 ns to access memory. Find the effective access time and percentage slowdown in memory access time if
  - i) Hit Ratio is 80%
  - ii) Hit Ratio is 98%
8. What are drawbacks of contiguous memory allocation? Given five memory partitions of 100KB, 500 KB, 200 KB, 300 KB and 600 KB . How would each of first fir, best fir and worst fir algorithms work to place processes of 212KB, 417 KB, 112KB and 426KB(in order)? Which is more efficient?

**Solution:**

#### **First Fit**

100K	
	212K
500K	112K
200K	
300K	

#### **Best Fit**

100K	
500K	417K
200K	112K
300K	212K
600K	426K

#### **Worst Fit**

100K	
500K	417K
200K	
300K	212K
600K	112K

600K | 417K |  
426 must  
wait

426 must  
wait

The Best Fit is the efficient algorithm.

9. What is locality of reference? Differentiate between paging and segmentation.
10. Explain the differences between.
  - i) Logical and Physical address space.
  - ii) Internal and External fragmentation.
11. What is paging? Give advantages and disadvantages.
12. What is segmentation? Give advantages & disadvantages.
13. What are the different methods of implementing page table?
14. How can you ensure protection & sharing in paging?
15. How to satisfy a request of size n from a list of free holes? Explain.
16. Explain the following:
  - a) Privileged instruction
  - b) Transient code
  - c) 50-percent rule
  - d) Roll in, roll out
  - e) Compaction
17. How can hardware address protection be ensured with base & limit registers?
18. What is dynamic loading? Give advantages.
19. Explain the following:
  - a) Frame table
  - b) Hit ratio
  - c) Re-entrant code
  - d) Legal page
  - e) TLB
20. Give a brief idea on dynamic linking & shared libraries.

## **MODULE -2**

### **MULTI-THREADED PROGRAMMING**

#### **Contents:**

- 2.1 Overview
- 2.2 Multithreading models
- 2.3 Thread Libraries
- 2.4 Threading issues
- 2.5 Process Scheduling
- 2.6 Basic concepts
- 2.7 Scheduling Criteria
- 2.8 Scheduling Algorithms
- 2.9 Multiple-processor scheduling
- 2.10 Thread scheduling

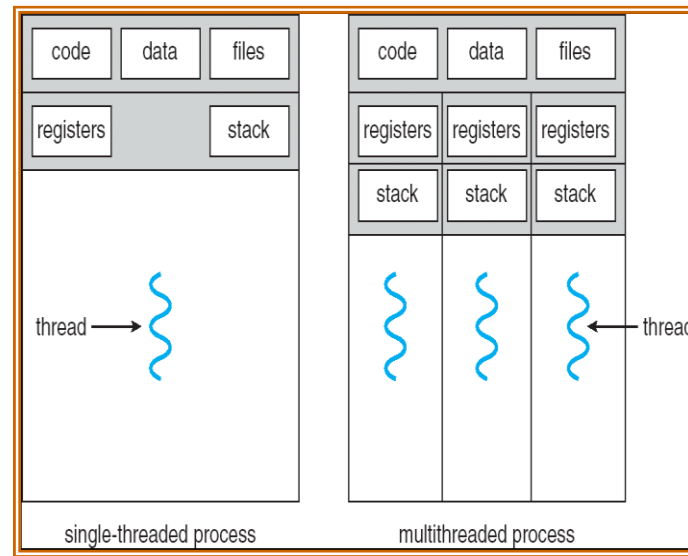
#### **Process Synchronization:**

- 2.11 Synchronization
- 2.12 The critical section problem
- 2.13 Peterson's solution
- 2.14 Synchronization hardware
- 2.15 Semaphores
- 2.16 Classical problems of synchronization
- 2.17 Monitors.

## 2.1 Overview

### ➤ Threads

- A thread is a basic unit of CPU utilization.
- It comprises a **thread ID**, a **program counter**, a **register set**, and a **stack**. It shares its **code section**, **data section**, and other operating-system resources, such as **open files** and **signals** with **other threads** belonging to the same process.
- A traditional (or **heavyweight**) process has a **single thread** of control. If a process has **multiple threads** of control, it can perform more than one task at a time.
- The below **figure** illustrates the difference between a **traditional single threaded** process and a **multithreaded** process.



### • Motivation

- Many software packages that run on modern desktop PCs are multithreaded.
- An application is implemented as a separate process with several threads of control. Eg: A Web browser might have one thread to display images or text while another thread retrieves data from the network.
- As **process creation** takes **more time** than **thread creation** it is more efficient to use process that contains multiple threads. So, that the amount of time that a client have to wait for its request to be serviced from the web server will be less.
- Threads also play an important role in **remote procedure call**.

### • Benefits

The benefits of multithreaded programming

1. **Responsiveness:** Multithreading allows program to continue running even if part of it is blocked or is performing a lengthy operation, thereby increasing responsiveness to the user.



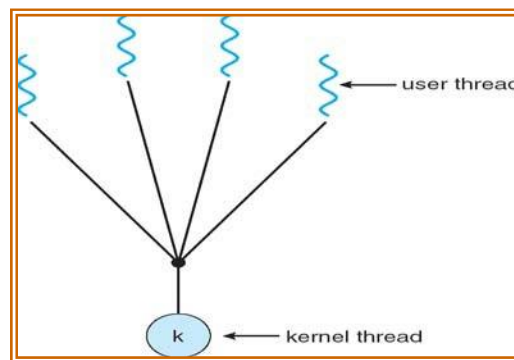
2. **Resource Sharing:** Threads share the memory and resources of the process to which they belong. The benefit of sharing code and data is that it allows an application to have several different threads of activity within the same address space.
3. **Economy:** Because of resource sharing context switching and thread creation are fast when working with threads.
4. **Utilization of multiprocessor architectures:** Threads can run in parallel on different processors. Multithreading on multi-CPU machine increases concurrency.

### ➤ Multithreading Models

- Support for threads may be provided either at user level for **user threads** or by the kernel for **kernel threads**.
- User threads are created and managed by Thread libraries without the kernel support.
- Kernel threads are directly managed by OS.
- There **must be a relationship** between user threads and kernel threads. **Three** common ways of establishing this relationship are :
  1. **Many-to-One**
  2. **One-to-One**
  3. **Many-to-Many**

- **Many-to-One**

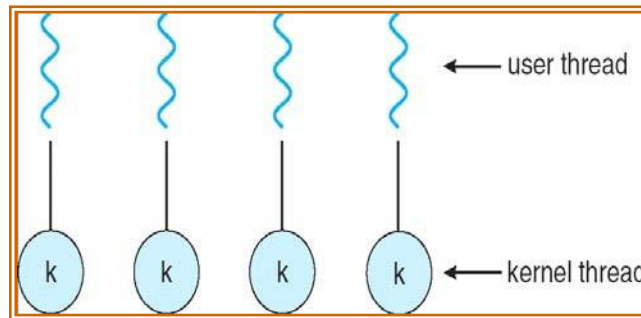
- Many user-level threads are mapped to single kernel thread as shown in below **figure**.
- This model is efficient as the thread management is done by the thread library in user space, but the entire process will block if a thread makes a blocking system call.
- As only one thread can access the kernel thread at a time, multiple threads are unable to run in parallel on multiprocessors.
- **Examples:** Solaris Green Threads, GNU Portable Threads



- **One-to-One**

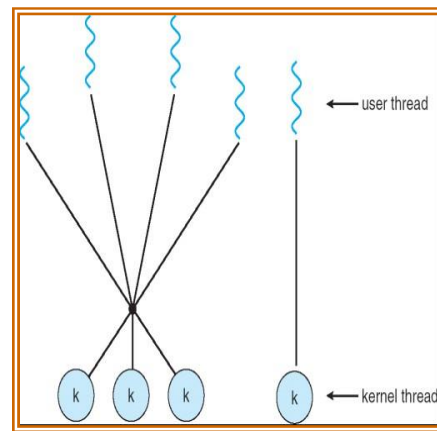
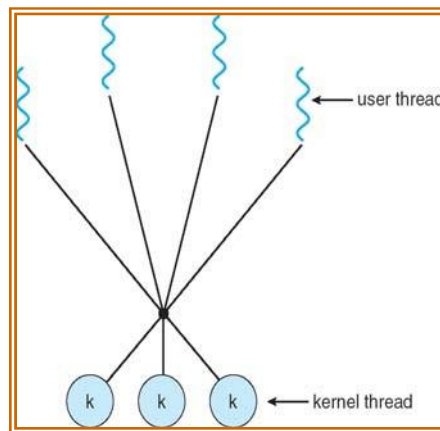
- Each user-level thread maps to kernel thread as shown in **figure 4.3**
- It provides more concurrency than Many-to-One model by allowing thread to run when a thread makes a blocking system call.
- It allows multiple threads to run in parallel on multiprocessors.

- The only **drawback** is, creating a user thread requires creating the corresponding kernel thread and it burdens performance of an application.
- **Examples:** Windows NT/XP/2000, Linux



- **Many-to-Many Model**

- **One-to-One** model **restricts** creating more user threads and **Many-to-One** model allows creating more user threads but kernel can schedule only one thread at a time. These drawbacks can be **overcome** by Many-to-Many model as shown in below **figure 4.4. (Left side)**
- Many-to-Many model allows many user level threads to be mapped to many kernel threads.
- It allows the operating system to create a sufficient number of kernel threads.
- When thread performs a blocking system call, the kernel can schedule another thread for execution.
- It allows user-level thread to be bound to a kernel thread and this is referred as **two- level model** as shown in below **figure 4.5. (Right side)**
- **Examples:** IRIX, HP-UX, Solaris OS.



➤ **Thread Libraries**

- A thread library provides the programmer an **API** for creating and managing threads.
- There are **two primary ways** of implementing a thread library.
- The first approach is to provide a **library entirely in user space** with **no kernel support**. All code and data structures for the library exist in user space.

- The second approach is to implement a **kernel-level library supported directly by the operating system**.
- Three primary **thread libraries** are
  - **POSIX Pthreads**: extension of posix standard, they may be provided as either a user or kernel library.
  - **Win32 threads**: is a kernel level library available on windows systems.
  - **Java threads**: API allows creation and management directly in Java programs. However, on windows java threads are implemented using win32 and on UNIX and Linux using Pthreads

- **Pthreads**

- ✓ Pthreads, the threads extension of the POSIX standard, may be provided as either a user or kernel-level library.
- ✓ Pthreads refers to the POSIX standard (IEEE 1003.1c) defining an API for thread creation and synchronization.
- ✓ This is a specification for thread behavior, not an implementation.
- ✓ Operating system designers may implement the specification in any way they wish. Numerous systems implement the Pthreads specification, including Solaris, Linux, Mac OS X, and Tru64 UNIX.
- ✓ **Shareware** implementations are available in the public domain for the various Windows operating systems as well.
- ✓

### **Multithreaded C program using the Pthreads API**

```
#include <pthread.h>
#include <stdio.h>
int sum;                                /* this data is shared by the thread(s) */
void *runner(void *param);             /* the thread */
int main(int argc, char *argv[])
{
    pthread_t tid;                      /* the thread identifier */
    pthread_attr_t attr;                /* set of thread attributes */

    if (argc != 2)
    {
        fprintf(stderr, "usage: a.out <integer value>\n");
        return -1;
    }
    if (atoi(argv[1]) < 0)
    {
        fprintf(stderr, "%d must be >= 0\n", atoi(argv[1]));
        return -1;
    }
    pthread_attr_init(&attr);            /* get the default attributes */
    pthread_create(&tid, &attr, runner, argv[1]); /* create the thread */
    pthread_join(tid, NULL);             /* wait for the thread to exit */
}
```

```
printf("sum = %d\n",sum);
}
```

```
void *runner(void *param)          /* The thread will begin control in this function */
{
int i, upper= atoi(param);
sum = 0;
for (i = 1; i <= upper; i++)
    sum += i;
pthread_exit(0) ;
}
```

- **Win32 Threads**

- The Win32 thread library is a kernel-level library available on Windows systems.
- The technique for creating threads using the Win32 thread library is similar to the Pthreads technique in several ways. We must include the windows.h header file when using the Win32 API.
- Threads are created in the Win32 API using the CreateThread() function and a set of attributes for the thread is passed to this function.
- These attributes include security information, the size of the stack, and a flag that can be set to indicate if the thread is to start in a suspended state.
- The parent thread waits for the child thread using the WaitForSingleObject() function, which causes the creating thread to block until the summation thread has exited.

### **Multithreaded C program using the Win32 API**

```
#include <Windows.h>
#include <stdio.h>
DWORD Sum;                                /* data is shared by the thread(s) */
                                           /* the thread runs in this separate function */
DWORD WINAPI Summation(LPVOID Param)
{
    DWORD Upper = *(DWORD*)Param;
    for (DWORD i = 0; i <= Upper; i++)
        Sum += i;
    return 0;
}

int main(int argc, char *argv[])
{
    DWORD ThreadId;
    HANDLE ThreadHandle;
    int Param;

    if (argc != 2)                        /* perform some basic error checking */
    {
        fprintf(stderr, " An integer parameter is required\n");
    }
}
```

```
    return -1;
}
Param = atoi(argv[1]);
if (Param < 0)
{
    fprintf(stderr, "An integer >= 0 is required\n");
    return -1;
}

/*create the thread*/
ThreadHandle = CreateThread(NULL, 0, Summation, &Param, 0, &ThreadId);

// NULL: default security attributes
// 0: default stack size
// Summation: thread function
// &Param: parameter to thread function
// 0: default creation flags
//ThreadId: returns the thread identifier

if (ThreadHandle != NULL)
{
    WaitForSingleObject(ThreadHandle, INFINITE); // now wait for the thread to finish
    CloseHandle(ThreadHandle); // close the thread handle
    printf("sum = %d\n", Sum);
}
}
```

- **Java Threads**

- The Java thread API allows thread creation and management directly in Java programs.
- Threads are the fundamental model of program execution in a Java program, and the Java language and its API provide a rich set of features for the creation and management of threads.
- All Java programs comprise at least a single thread of control and even a simple Java program consisting of only a **main()** method runs as a single thread in the JVM.
- There are **two techniques** for creating threads in a Java program. One approach is to create a new class that is derived from the **Thread class** and to override its **run()** method. An alternative and more commonly used technique is to define a class that implements the **Runnable interface**. The **Runnable interface** is defined as follows:

```
public interface Runnable
{
    public abstract void run () ;
}
```

- When a class implements Runnable, it must define a **run()** method. The code implementing the run() method runs as a separate thread.

- Creating a Thread object does not specifically create the new thread but it is the **start()** method that actually creates the new thread. Calling the **start()** method for the new object does two things:
  - It allocates memory and initializes a new thread in the JVM.
  - It calls the **run()** method, making the thread eligible to be run by the JVM.
- As Java is a **pure object-oriented** language, it has **no notion of global data**. If two or more threads have to share data means then the sharing occurs by passing reference to the shared object to the appropriate threads.
- This shared object is referenced through the appropriate **getSum()** and **setSum()** methods.
- As the Integer class is **immutable**, that is, once its value is set it cannot change, a new **sum** class is designed.
- The parent threads in Java uses **join()** method to wait for the child threads to finish before proceeding.

**Java program for the summation of a non-negative integer.**

```
class Sum
{
    private int sum;
    public int getSum()
    {
        return sum;
    }

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

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

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

    public void run()
    {
        int sum = 0;
        for (int i = 0; i <= upper; i++)
```

```
        sum += i;
        sumValue.setSum(sum);
    }
}

public class Driver
{
    public static void main(String[] args)
    {
        if (args.length > 0)
        {
            if (Integer.parseInt(args[0]) < 0)
                System.err.println(args[0] + "must be>= 0.");
            else
            {
                Sum sumObject = new Sum();    //create the object to be shared
                int upper= Integer.parseInt(args[0]);
                Thread thrd =new Thread(new Summation(upper, sumObject));
                thrd.start();
                try
                {
                    thrd. join () ;
                    System.out.println("The sum of "+upper+" is "+sumObject.getSum());
                }
                catch (InterruptedException ie) { }
            }
        }
        else
            System.err.println("Usage: Summation <integer value>");
    }
}
```

## ➤ Threading Issues

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

- ✓ The semantics of the **fork()** and **exec()** system calls change in a multithreaded program. Some UNIX systems have chosen to have **two versions** of **fork()**, **one** that duplicates all threads and **another** that duplicates only the thread that invoked the fork() system call.
- ✓ If a thread invokes the **exec()** system call, the program specified in the parameter to **exec()** will replace the entire process including all threads.
- ✓ Which of the two versions of **fork()** to use depends on the application. If **exec()** is called immediately after forking, then duplicating all threads is unnecessary, as the program specified in the parameters to exec() will replace the process. In this instance, duplicating only the calling thread is appropriate.

- **Thread Cancellation**

- ✓ Thread cancellation is the task of terminating a thread before it has completed. **For example**, if multiple threads are concurrently searching through a database and one thread returns the result, the remaining threads might be cancelled.
- ✓ A thread that is to be cancelled is often referred to as the **target thread**. Cancellation of a target thread may occur in **two different** scenarios:
  - **Asynchronous cancellation:** One thread immediately terminates the target thread.
  - **Deferred cancellation:** The target thread periodically checks whether it should terminate, allowing it an opportunity to terminate itself in an orderly fashion.
- ✓ The **difficulty** with asynchronous cancellation occurs in situations where resources have been allocated to a cancelled thread or where a thread is cancelled while in the midst of updating data it is sharing with other threads. This becomes especially troublesome with asynchronous cancellation. Often, the operating system will reclaim system resources from a cancelled thread but will not reclaim all resources. Therefore, cancelling a thread asynchronously may not free a necessary system-wide resource.
- ✓ With deferred cancellation, one thread indicates that a target thread is to be cancelled, but cancellation occurs only after the target thread has checked a flag to determine if it should be cancelled or not. This allows a thread to check whether it should be cancelled at a point when it can be cancelled safely. Pthreads refers to such points as **cancellation points**.

- **Signal Handling**

- ✓ A signal is used in UNIX systems to notify a process that a particular event has occurred.
- ✓ A signal may be received either synchronously or asynchronously, depending on the **source of** and **the reason** for the event being signaled.
- ✓ All signals, whether synchronous or asynchronous, follow the same pattern:
  - A signal is generated by the occurrence of a particular event.
  - A generated signal is delivered to a process.
  - Once delivered, the signal must be handled.
- ✓ **Examples** of synchronous signals include illegal memory access and division by 0. If a running program performs either of these actions, a signal is generated. Synchronous signals are delivered to the same process that performed the operation that caused the signal.
- ✓ When a signal is generated by an event external to a running process, that process receives the signal **asynchronously**. **Examples** of such signals include terminating a process with specific keystrokes (such as <control><C>) and having a **timer expires**. An asynchronous signal is sent to another process.
- ✓ Every signal may be handled by one of **two possible handlers**,
  - A default signal handler
  - A user-defined signal handler
- ✓ Every signal has a **default signal handler** that is run by the kernel when handling that signal.
- ✓ This default action can be overridden by a **user-defined signal handler** that is called to handle the signal.



- ✓ Signals may be handled in different ways. Some signals (such as changing the size of a window) may simply be ignored; others (such as an illegal memory access) may be handled by terminating the program.
- ✓ Delivering signals is more complicated in multithreaded programs. The following **options exist** to deliver a signal:
  - Deliver the signal to the thread to which the signal applies.
  - Deliver the signal to every thread in the process.
  - Deliver the signal to certain threads in the process.
  - Assign a specific thread to receive all signals for the process.

- **Thread Pools**

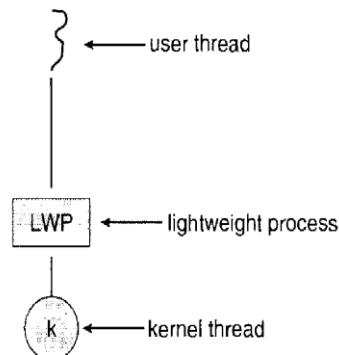
- ✓ The idea behind a thread pool is to **create a number of threads** at process startup and place them into a pool, where they sit and wait for work.
- ✓ When a server receives a request, it awakens a thread from this pool and passes the request to it to service.
- ✓ Once the thread completes its service, it returns to the pool and waits for more work.
- ✓ If the pool contains no available thread, the server waits until one becomes free.
- ✓ The **benefits** of Thread pools are,
  - Servicing a request with an existing thread is usually faster than waiting to create a thread.
  - A thread pool limits the number of threads that exist at any one point. This is particularly important on systems that cannot support a large number of concurrent threads.
- ✓ The **number of threads** in the pool can be set based on **factors** such as the number of CPUs in the system, the amount of physical memory, and the expected number of concurrent client requests.

- **Thread-Specific Data**

- ✓ Threads belonging to a process share the data of the process. This sharing of data provides one of the benefits of multithreaded programming. But, in some circumstances, each thread might need its own copy of certain data. Such data is called as **thread-specific data**. **For example**, in a transaction-processing system, we might service each transaction in a separate thread. Furthermore, each transaction may be assigned a unique identifier.
- ✓ Most thread libraries including Win32 and Pthreads provide support for thread-specific data.

- **Scheduler Activations**

- ✓ Many systems implementing either the many-to-many or two-level model place an **intermediate data structure** between the user and kernel threads. This data structure is known as a lightweight process, or LWP as shown in the following figure 4.6



- ✓ An application may require any number of LWPs to run efficiently.
- ✓ In a CPU-bound application running on a single processor only one thread can run at once, so one LWP is sufficient. An application that is I/O- intensive may require multiple LWPs to execute.
- ✓ One scheme for communication between the user-thread library and the kernel is known as

**scheduler activation.** It works as follows: The kernel provides an application with a set of **virtual processors (LWPs)**, and the application can schedule user threads onto an available virtual processor. The kernel must inform an application about certain events. This procedure is known as an **upcall**.

- ✓ Upcalls are handled by the thread library with an upcall handler, and upcall handlers must run on a virtual processor.
- ✓ One event that triggers an upcall occurs when an application thread is about to block. In this

situation, the kernel makes an upcall to the application informing it that a thread is about to block and identifying the specific thread. The kernel then allocates a new virtual processor to the

application.

- ✓ The application runs an upcall handler on this new virtual processor, which saves the state of the blocking thread and gives up the virtual processor on which the blocking thread is running. The upcall handler then schedules another thread that is eligible to run on the new virtual processor.
- ✓ When the event that the blocking thread was waiting for occurs, the kernel makes another upcall to the thread library informing it that the previously blocked thread is now eligible to run.
- ✓ The upcall handler for this event also requires a virtual processor, and the kernel may allocate a new virtual processor or preempt one of the user threads and run the upcall handler on its virtual processor.
- ✓ After marking the unblocked thread as eligible to run, the application schedules an eligible thread to run on an available virtual processor.

### ➤ Differences

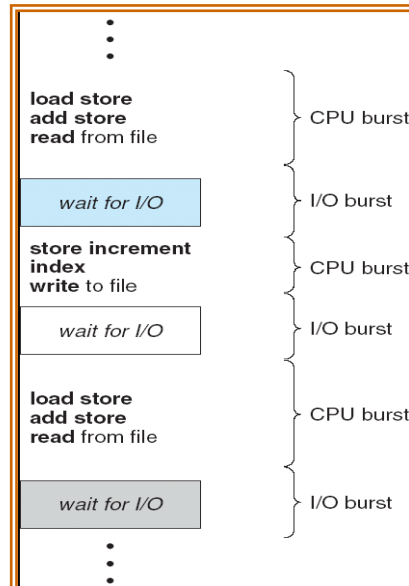
- ✓ The differences between process and thread are,

	Process	Thread
1.	It is called heavyweight process.	It is called lightweight process.
2.	Process switching needs interface with OS.	Thread switching does not need interface with OS.
3.	Multiple processes use more resources than multiple threads.	Multiple threaded processes use fewer resources than multiple processes.
4.	In multiple process implementations each process executes same code but has its own memory and file resources.	All threads can share same set of open files.
5.	If one server process is blocked no other server process can execute until the first process unblocked.	While one server thread is blocked and waiting, second thread in the same task could run.
6.	In multiple processes each process operates independently of others.	One thread can read, write or even completely wipeout another threads stack.

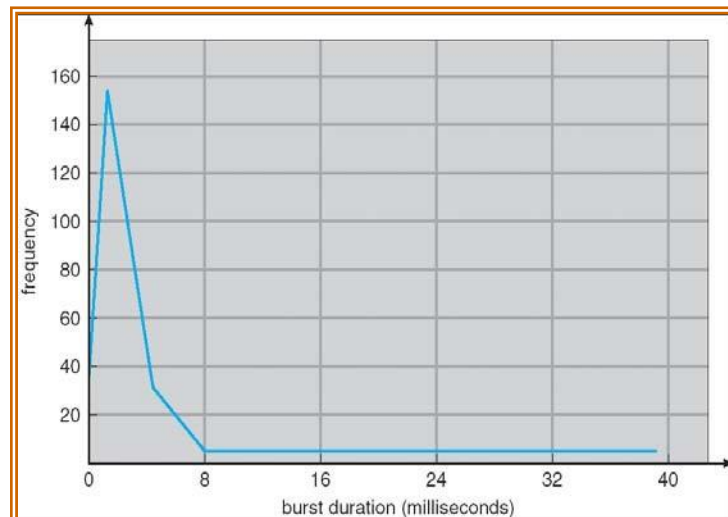
## PROCESS SCHEDULING

### ➤ Basic Concepts

- ✓ In a single-processor system, only one process can run at a time and others must wait until the CPU is free and can be rescheduled.
- ✓ The objective of **multiprogramming** is to have some process running at all times, to **maximize CPU utilization**.
- ✓ With multiprogramming, several processes are kept in memory at one time. When one process has to wait, the operating system takes the CPU away from that process and gives the CPU to another process.
- ✓ The CPU is one of the primary computer resources. Thus, its scheduling is central to
- ✓ operating-system design.
- **CPU-I/O Burst Cycle**
  - ✓ Process execution consists of a **cycle** of CPU execution and I/O wait
  - ✓ Process execution starts with **CPU burst** and this is followed by **I/O burst** as shown in below **figure 5.1**.
  - ✓ The final CPU burst ends with a system request to terminate execution.



- ✓ The duration of CPU bursts vary from process to process and from computer to computer.
- ✓ The **frequency curve** is as shown below **figure 5.2**.



## • CPU Scheduler

- ✓ Whenever the CPU becomes idle, the operating system must select one of the processes in the ready queue to be executed. This selection is carried out by the **short-term scheduler (or CPU scheduler)**.
- ✓ The scheduler selects a process from the processes in memory that are ready to execute and allocates the CPU to that process.

- ✓ A ready queue can be implemented as a FIFO queue, a priority queue, a tree, or simply an unordered linked list. But all the processes in the ready queue are lined up waiting for a chance to run on the CPU. The **records** in the queues are **Process Control Blocks (PCBs)** of the processes.
- **Preemptive scheduling**
  - ✓ **CPU-scheduling decisions** may take place under the following **four circumstances**.
    1. When a process switches from the running state to the waiting state (for example, as the result of an I/O request or an invocation of wait for the termination of one of the child processes)
    2. When a process switches from the running state to the ready state (for example, when an interrupt occurs)
    3. When a process switches from the waiting state to the ready state (for example, at completion of I/O)
    4. When a process terminates
  - ✓ When scheduling takes place only under circumstances **1 and 4**, we say that the scheduling scheme is **nonpreemptive** or **cooperative**; otherwise, it is **preemptive**.
  - ✓ Under nonpreemptive scheduling, once the CPU has been allocated to a process, the process keeps the CPU until it releases the CPU either by terminating or by switching to the waiting state. A scheduling algorithm is preemptive if, once a process has been given the CPU and it can be taken away.
- **Dispatcher**
  - ✓ Another component involved in the CPU-scheduling function is the dispatcher.
  - ✓ The dispatcher is the module that gives control of the CPU to the process selected by the short-term scheduler.
  - ✓ This function involves the following:
    - Switching context
    - Switching to user mode
    - Jumping to the proper location in the user program to restart that program
  - ✓ The dispatcher should be as fast as possible, since it is invoked during every process switch.
  - ✓ The time it takes for the dispatcher to stop one process and start another running is known as the **dispatch latency**.

## ➤ Scheduling Criteria

- ✓ Many criteria have been suggested for comparing CPU scheduling algorithms. The **criteria** include the following:
  - **CPU utilization:** CPU must be kept as busy as possible. CPU utilization can range from 0 to 100 percent. In a real system, it should range from 40 percent (for a lightly loaded system) to 90 percent (for a heavily used system).

- **Throughput:** If the CPU is busy executing processes, then work is being done. One **measure of work** is the number of processes that are completed per time unit, called throughput. For long processes, this rate may be one process per hour; for short transactions, it may be 10 processes per second.
- **Turnaround time:** The interval from the time of submission of a process to the time of completion is the turnaround time. Turnaround time is the sum of the periods spent waiting to get into memory, waiting in the ready queue, executing on the CPU, and doing I/O.
- **Waiting time:** The CPU scheduling algorithm does not affect the amount of time during which a process executes or does I/O; it affects only the amount of time that a process spends waiting in the ready queue. Waiting time is the sum of the periods spent waiting in the ready queue.
- **Response time:** The measure of the time from the submission of a request until the first response is produced. This measure, called response time, is the time it takes to start responding, not the time it takes to output the response.

### ➤ Scheduling Algorithms

CPU Scheduling deals with the problem of **deciding** which of the processes in the ready queue is to be allocated the CPU. Following are some **scheduling algorithms**,

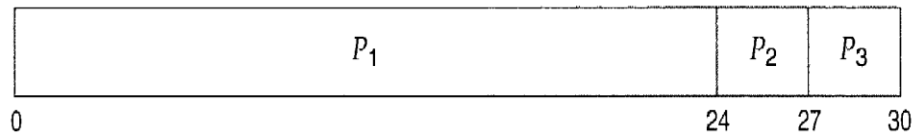
- FCFS Scheduling.
- Round Robin Scheduling.
- SJF Scheduling.
- Priority Scheduling.
- Multilevel Queue Scheduling.
- Multilevel Feedback Queue Scheduling.

#### • First-Come-First-Served (FCFS) Scheduling

- ✓ The simplest CPU-scheduling algorithm is the first-come, first-served (FCFS) scheduling algorithm.
- ✓ With this scheme, the process that requests the CPU first is allocated the CPU first.
- ✓ The implementation of the FCFS policy is easily managed with a FIFO queue.
- ✓ When a process enters the ready queue, its PCB is linked onto the tail of the queue.
- ✓ When the CPU is free, it is allocated to the process at the head of the queue. The running process is then removed from the queue.
- ✓ The average waiting time under the FCFS policy is often quite long. Consider the following set of processes that arrive at time 0, with the length of the CPU burst given in milliseconds:

Process	Burst Time
$P_1$	24
$P_2$	3
$P_3$	3

- ✓ If the processes arrive in the order  $P_1, P_2, P_3$ , and are served in FCFS order, we get the result shown in the following Gantt chart:



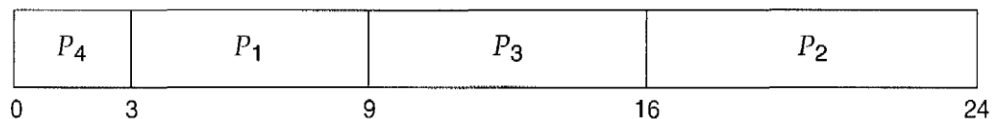
- ✓ The waiting time is 0 milliseconds for process P<sub>1</sub>, 24 milliseconds for process P<sub>2</sub>, and 27 milliseconds for process P<sub>3</sub>. Thus, the average waiting time is  $(0 + 24 + 27)/3 = 17$  milliseconds.
- ✓ The FCFS scheduling algorithm is **nonpreemptive**. Once the CPU has been allocated to a process, that process keeps the CPU until it releases the CPU, either by terminating or by requesting I/O. The FCFS algorithm is thus particularly troublesome for time-sharing systems, where it is important that each user get a share of the CPU at regular intervals.

- **Shortest-Job-First Scheduling**

- ✓ This algorithm associates with each process the length of the process's next CPU burst.
- ✓ When the CPU is available, it is assigned to the process that has the smallest next CPU burst.
- ✓ If the next CPU bursts of two processes are the same, FCFS scheduling is used to break the tie.
- ✓ As an example of SJF scheduling, consider the following set of processes, with the length of the CPU burst given in milliseconds:

Process	Burst Time
P <sub>1</sub>	6
P <sub>2</sub>	8
P <sub>3</sub>	7
P <sub>4</sub>	3

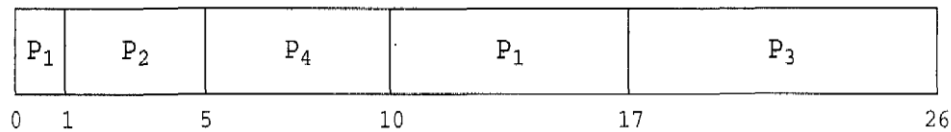
- ✓ Using SJF scheduling, we would schedule these processes according to the following Gantt chart:



- ✓ The waiting time is 3 milliseconds for process P<sub>1</sub>, 16 milliseconds for process P<sub>2</sub>, 9 milliseconds for process P<sub>3</sub>, and 0 milliseconds for process P<sub>4</sub>. Thus, the average waiting time is  $(3 + 16 + 9 + 0)/4 = 7$  milliseconds.
- ✓ The SJF scheduling algorithm is **optimal**, it gives the minimum average waiting time for a given set of processes. Moving a short process before a long one, decreases the waiting time of the short process more than it increases the waiting time of the long process. Consequently, the average waiting time decreases.
- ✓ The SJF algorithm can be either **preemptive or nonpreemptive**. The choice arises when a new process arrives at the ready queue while a previous process is still executing. The next CPU burst of the newly arrived process may be shorter than what is left of the currently executing process. A preemptive SJF algorithm will preempt the currently executing process, whereas a nonpreemptive SJF algorithm will allow the currently running process to finish its CPU burst. Preemptive SJF scheduling is sometimes called **shortest-remaining-time-first scheduling**.
- ✓ As an example, consider the following four processes, with the length of the CPU burst given in milliseconds:

Process	Arrival Time	Burst Time
$P_1$	0	8
$P_2$	1	4
$P_3$	2	9
$P_4$	3	5

- ✓ If the processes arrive at the ready queue at the times shown and need the indicated burst times, then the resulting preemptive SJF schedule is as depicted in the following Gantt chart:



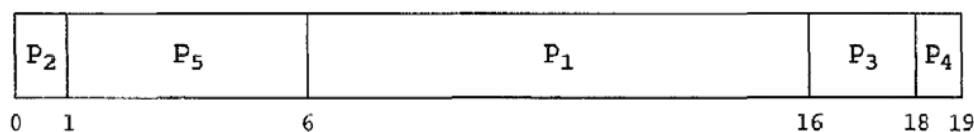
- ✓ Process  $P_1$  is started at time 0, since it is the only process in the queue. Process  $P_2$  arrives at time 1.
- ✓ The remaining time for process  $P_1$  (7 milliseconds) is larger than the time required by process  $P_2$  (4 milliseconds), so process  $P_1$  is preempted, and process  $P_2$  is scheduled. The average waiting time for this example is  $((10 - 1) + (1 - 1) + (17 - 2) + (5 - 3))/4 = 26/4 = 6.5$  milliseconds. Nonpreemptive SJF scheduling would result in an average waiting time of 7.75 milliseconds.

## • Priority Scheduling

- ✓ The SJF algorithm is a special case of the general priority scheduling algorithm.
- ✓ A priority is associated with each process, and the CPU is allocated to the process with the highest priority.
- ✓ Equal-priority processes are scheduled in FCFS order.
- ✓ An SJF algorithm is simply a priority algorithm where the priority ( $p$ ) is the inverse of the (predicted) next CPU burst. The larger the CPU burst, the lower the priority, and vice versa.
- ✓ As an example, consider the following set of processes, assumed to have arrived at time 0, in the order  $P_1, P_2, \dots, P_5$ , with the length of the CPU burst given in milliseconds:

Process	Burst Time	Priority
$P_1$	10	3
$P_2$	1	1
$P_3$	2	4
$P_4$	1	5
$P_5$	5	2

- ✓ Using priority scheduling, we would schedule these processes according to the following Gantt chart:



- ✓ The average waiting time is 8.2 milliseconds.



- ✓ Priority scheduling can be either **preemptive or nonpreemptive**. When a process arrives at the ready queue, its priority is compared with the priority of the currently running process. A preemptive priority scheduling algorithm will preempt the CPU if the priority of the newly arrived process is higher than the priority of the currently running process. A nonpreemptive priority scheduling algorithm will simply put the new process at the head of the ready queue.
- ✓ A **major problem** with priority scheduling algorithms is **indefinite blocking, or starvation**. A process that is ready to run but waiting for the CPU can be considered blocked. A priority scheduling algorithm can leave some low-priority processes waiting indefinitely. In a heavily loaded computer system, a steady stream of higher-priority processes can prevent a low-priority process from ever getting the CPU.
- ✓ A **solution** to the problem of indefinite blockage of low-priority processes is **aging**. Aging is a technique of gradually increasing the priority of processes that wait in the system for a long time.

- **Round-Robin Scheduling**

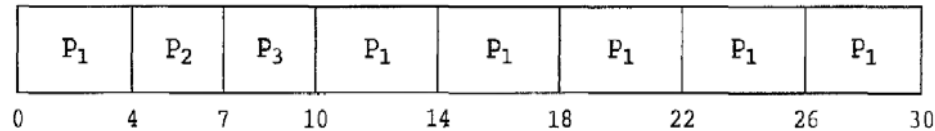
- ✓ The round-robin (RR) scheduling algorithm is designed especially for timesharing systems.
- ✓ It is similar to FCFS scheduling, but preemption is added to switch between processes.
- ✓ A small unit of time, called a **time quantum or time slice**, is defined. A time quantum is generally from 10 to 100 milliseconds.
- ✓ The ready queue is treated as a circular queue. The CPU scheduler goes around the ready queue, allocating the CPU to each process for a time interval of up to 1 time quantum.
- ✓ To implement RR scheduling, we keep the ready queue as a FIFO queue of processes. New processes are added to the tail of the ready queue.
- ✓ The CPU scheduler picks the first process from the ready queue, sets a timer to interrupt after 1 time quantum, and dispatches the process. One of two things will then happen.
  - The process may have a CPU burst of less than 1 time quantum. In this case, the process itself will release the CPU voluntarily. The scheduler will then proceed to the next process in the ready queue.
  - Otherwise, if the CPU burst of the currently running process is longer than 1 time quantum, the timer will go off and will cause an interrupt to the operating system. A context switch will be executed, and the process will be put at the tail of the ready queue. The CPU scheduler will then select the next process in the ready queue.
- ✓ The average waiting time under the RR policy is often long. Consider the following set of processes that arrive at time 0, with the length of the CPU burst given in milliseconds:

Process	Burst Time
$P_1$	24
$P_2$	3
$P_3$	3

- ✓ If we use a time quantum of 4 milliseconds, then process  $P_1$  gets the first 4 milliseconds. Since it requires another 20 milliseconds, it is preempted after the first time quantum, and the CPU is given

to the next process in the queue ie. process P<sub>2</sub>. Since process P<sub>2</sub> does not need 4 milliseconds, it quits before its time quantum expires.

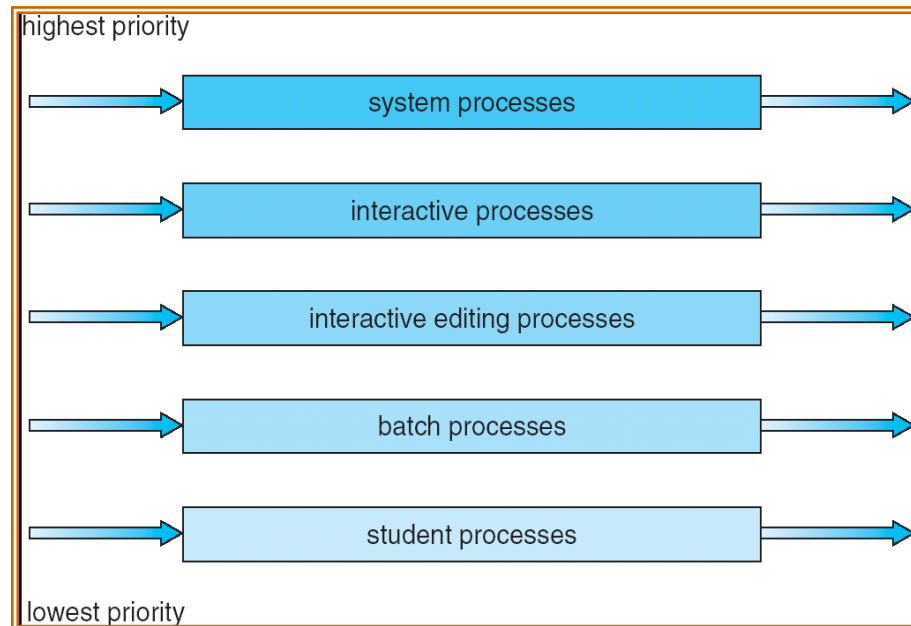
- ✓ The CPU is then given to the next process P<sub>3</sub>. Once each process has received 1 time quantum, the CPU is returned to process P<sub>1</sub> for an additional time quantum. The resulting RR schedule is



- ✓ The average waiting time is  $17/3 = 5.66$  milliseconds.
- ✓ The RR scheduling algorithm is thus preemptive.
- ✓ If there are n processes in the ready queue and the time quantum is q, then each process gets 1/n of the CPU time in chunks of at most q time units. Each process must wait no longer than  $(n-1) * q$  time units until its next time quantum. **For example**, with five processes and a time quantum of 20 milliseconds, each process will get up to 20 milliseconds every 100 milliseconds.

### • Multilevel Queue Scheduling

- ✓ Another class of scheduling algorithms has been created for situations in which processes are easily classified into different groups. For example, a common division is made between **foreground (interactive) processes and background (batch) processes**.
- ✓ These two types of processes have different **response-time** requirements and may have **different scheduling** needs.
- ✓ Foreground processes have priority over background processes.
- ✓ A multilevel queue scheduling algorithm partitions the ready queue into several separate queues as shown in **figure 5.3**.
- ✓ The processes are permanently assigned to one queue based on some property of the process, such as memory size, process priority, or process type.
- ✓ Each queue has its own scheduling algorithm. For example, separate queues might be used for foreground and background processes.
- ✓ The foreground queue might be scheduled by an **RR algorithm**, while the background queue is scheduled by an **FCFS algorithm**.

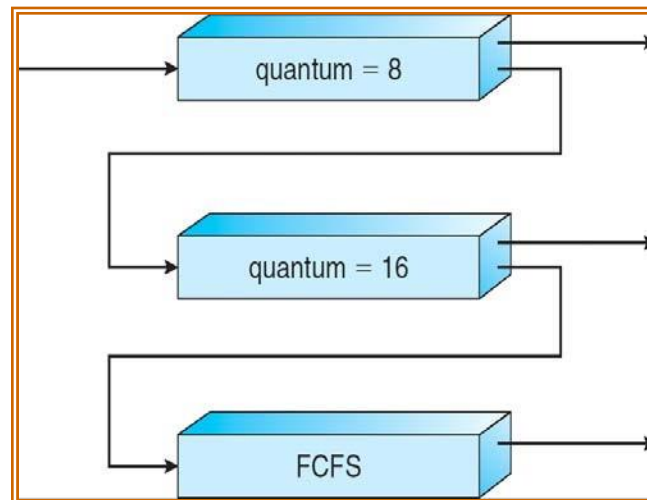


- ✓ There must be scheduling among the queues, which is commonly implemented as **fixed-priority preemptive** scheduling. For example, the foreground queue may have absolute priority over the background queue.
- ✓ An example of a multilevel queue scheduling algorithm with **five queues**, listed below in order of priority:
  1. System processes
  2. Interactive processes
  3. Interactive editing processes
  4. Batch processes
  5. Student processes
- ✓ Each queue has absolute priority over lower-priority queues. No process in the batch queue could run unless the queues for system processes, interactive processes, and interactive editing processes were all empty.
- ✓ If an interactive editing process entered the ready queue while a batch process was running, the batch process would be preempted.
- ✓ Another possibility is to **time-slice** among the queues. So, each queue gets a certain portion of the CPU time, which it can then schedule among its various processes. For example, the foreground queue can be given **80 percent** of the CPU time for RR scheduling among its processes, whereas the background queue receives **20 percent** of the CPU to give to its processes on an FCFS basis.

- **Multilevel Feedback-Queue Scheduling**

- ✓ When the multilevel queue scheduling algorithm is used, processes are permanently assigned to a queue when they enter the system.
- ✓ The multilevel feedback-queue scheduling algorithm allows a process to move between queues. The idea is to separate processes according to the characteristics of their CPU bursts.

- ✓ If a process uses too much CPU time, it will be moved to a lower-priority queue. This scheme leaves I/O-bound and interactive processes in the higher-priority queues.
- ✓ A process that waits too long in a lower-priority queue may be moved to a higher-priority queue. This form of aging prevents starvation.
- ✓ For example, consider a multilevel feedback-queue scheduler with three queues, numbered from 0 to 2 as shown in below **figure 5.4**.
- ✓ The scheduler first executes all processes in queue 0. Only when queue 0 is empty it will execute processes in queue 1.
- ✓ Similarly, processes in queue 2 will only be executed if queues 0 and 1 are empty. A process that arrives for queue 1 will preempt a process in queue 2.
- ✓ A process in queue 1 will in turn be preempted by a process arriving for queue 0.



- ✓ A process entering the ready queue is put in queue 0. A process in queue 0 is given a time quantum of 8 milliseconds. If it does not finish within this time, it is moved to the tail of queue 1. If queue 0 is empty, the process at the head of queue 1 is given a quantum of 16 milliseconds. If it does not complete, it is preempted and is put into queue 2. Processes in queue 2 are scheduled on an FCFS basis but they run only when queue 0 and 1 are empty.
- ✓ This scheduling algorithm gives highest priority to any process with a CPU burst of 8 milliseconds or less. Such a process will quickly get the CPU, finish its CPU burst, and go off to its next I/O burst. Processes that need more than 8 but less than 24 milliseconds are also served quickly, although with lower priority than shorter processes. Long processes automatically sink to queue 2 and are served in FCFS order with any CPU cycles left over from queues 0 and 1.
- ✓ A multilevel feedback-queue scheduler is defined by the following **parameters**:
  - The number of queues.
  - The scheduling algorithm for each queue.
  - The method used to determine when to upgrade a process to a higher-priority queue.
  - The method used to determine when to demote a process to a lower-priority queue.
  - The method used to determine which queue a process will enter when that process needs service.

## ➤ Multiple-Processor Scheduling

- **Approaches to Multiple-Processor Scheduling**

- ✓ One approach to CPU scheduling in a multiprocessor system is where all scheduling decisions, I/O processing, and other system activities are handled by a single processor i.e., **the master server**. The other processors execute only user code. This **asymmetric multiprocessing** is simple because only one processor accesses the system data structures, reducing the need for data sharing.
- ✓ A second approach uses **symmetric multiprocessing (SMP)**, where each processor is self-scheduling. All processes may be in a common ready queue, or each processor may have its own private queue of ready processes.

- Some of the **issues** related to SMP are,

### **a. Processor Affinity**

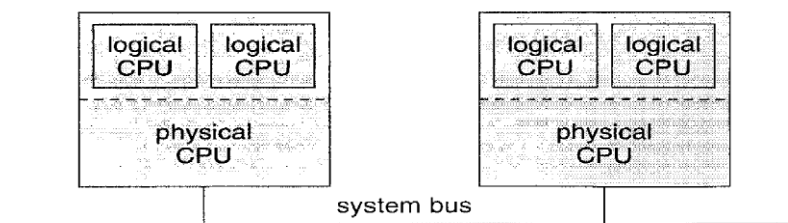
- ✓ The data most recently accessed by the process is populated in the **cache** for the processor and successive memory accesses by the process are often satisfied in cache memory.
- ✓ If the process **migrates** to another processor, the contents of cache memory must be invalidated for the processor being **migrated from**, and the cache for the processor being **migrated to** must be re-populated. Because of the **high cost** of invalidating and re-populating caches, most SMP systems try to avoid migration of processes from one processor to another and instead tries to keep a process running on the same processor. This is known as **processor affinity**, i.e., a process has an affinity for the processor on which it is currently running.
- ✓ Processor affinity takes **several forms**. When an operating system has a policy of attempting to keep a process running on the same processor but not guaranteeing that it will do so, a situation is known as **soft affinity**. Here, it is possible for a process to migrate between processors.
- ✓ Some systems such as Linux provide system calls that support **hard affinity**, thereby allowing a process to specify that it must not migrate to other processors.

### **b. Load Balancing**

- ✓ On SMP systems, it is important to keep the workload balanced among all processors to utilize the benefits of having more than one processor. Otherwise, one or more processors may sit idle while other processors have high workloads along with lists of processes awaiting the CPU.
- ✓ Load balancing attempts to keep the workload evenly distributed across all processors in an SMP system.
- ✓ There are two general approaches to load balancing: **push migration** and **pull migration**.
- ✓ With push migration, a specific task periodically checks the load on each processor and if it finds an imbalance it evenly distributes the load by moving (or pushing) processes from overloaded to idle or less-busy processors.
- ✓ Pull migration occurs when an idle processor pulls a waiting task from a busy processor.

### c. Symmetric Multithreading

- ✓ SMP systems allow several threads to run concurrently by providing multiple physical processors.
- ✓ An alternative strategy is to provide **multiple logical processors** rather than physical processors. Such a strategy is known as symmetric multithreading (or SMT).
- ✓ The idea behind SMT is to create multiple logical processors on the same physical processor, presenting a view of several logical processors to the operating system, even on a system with only a single physical processor.
- ✓ Each logical processor has its own architecture state, which includes general-purpose and machine-state registers and is responsible for its own interrupt handling, meaning that interrupts are delivered to and handled by logical processors rather than physical ones. Otherwise, each logical processor shares the resources of its physical processor, such as cache memory and buses.
- ✓ The following **figure 5.5** illustrates a typical **SMT architecture** with two physical processors, each housing two logical processors. From the operating system's perspective, four processors are available for work on this system.



### ➤ Thread Scheduling

- ✓ On operating systems that support user-level and kernel-level threads, the kernel-level threads are being scheduled by the operating system.
- ✓ User-level threads are managed by a thread library, and the kernel is unaware of them. To run on a CPU, user-level threads must be mapped to an associated kernel-level thread, although this mapping may be indirect and may use a lightweight process (LWP).
- ✓ One **distinction** between user-level and kernel-level threads lies in how they are **scheduled**. On systems implementing the many-to-one and many-to-many models, the thread library schedules user-level threads to run on an available LWP, a scheme known as **process-contention scope (PCS)**, since competition for the CPU takes place among threads belonging to the same process.
- ✓ To decide which kernel thread to schedule onto a CPU, the kernel uses **system-contention scope (SCS)**. Competition for the CPU with SCS scheduling takes place among all threads in the system.
- ✓ PCS is done according to priority. The scheduler selects the runnable thread with the highest priority to run. User-level thread priorities are set by the programmer. PCS will preempt the currently running thread in favour of a higher-priority thread.

- **Pthread Scheduling**

- ✓ Pthreads identifies the following contention scope values:
  - PTHREAD\_SCOPE\_PROCESS schedules threads using PCS scheduling.
  - PTHREAD\_SCOPE\_SYSTEM schedules threads using SCS scheduling.
- ✓ On systems implementing the many-to-many model, the PTHREAD\_SCOPE\_PROCESS policy schedules user-level threads onto available LWPs.
- ✓ The number of LWPs is maintained by the thread library using scheduler activations. The PTHREAD\_SCOPE\_SYSTEM scheduling policy will create and bind an LWP for each user-level thread on many-to-many systems, effectively mapping threads using the one-to-one policy.
- ✓ The Pthread IPC provides the following **two functions** for getting and setting the contention scope policy;
  - pthread\_attr\_setscope (pthread\_attr\_t \*attr, int scope)
  - pthread\_attr\_getscope (pthread\_attr\_t \*attr, int \*scope)
- ✓ The first parameter for both functions contains a pointer to the attribute set for the thread.
- ✓ The second parameter for the pthread\_attr\_setscope () function is passed either the PTHREAD\_SCOPE\_SYSTEM or PTHREAD\_SCOPE\_PROCESS value, indicating how the contention scope is to be set. In the case of pthread\_attr\_getscope(), this second parameter contains a pointer to an int value that is set to the current value of the contention scope. If an error occurs, each of these functions returns non-zero values.

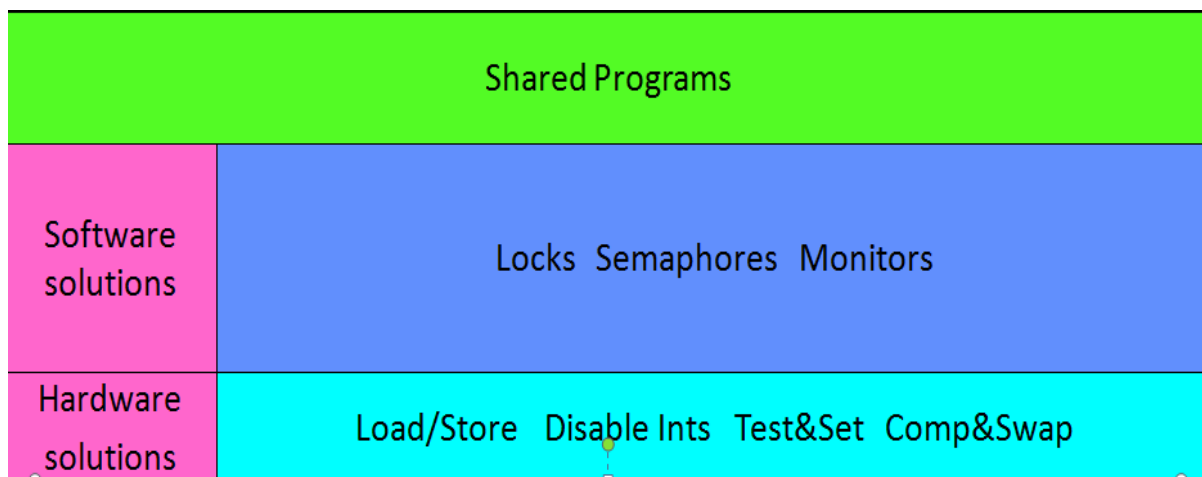
## **Module 2 (Continued)**

### **Process Synchronization**

- The critical section problem;
- Peterson's solution;
- Synchronization hardware;
- Semaphores;
- Classical problems of synchronization;
- Monitors.

#### **2.3.1 Introduction:**

A co-operating process is one that can affect or be affected by other process executing in the system. It can use shared memory or message passing for communication. If shared memory is implemented, the access to the share data must be taken care and will be discussed in detail in this chapter.



#### **2.3.2 Race Condition:**

A situation where several processes access and manipulate the same data concurrently and the outcome of the execution depends on the particular order in which the access takes place, is called a **race condition**. To guard against race condition, we require that the processes must be synchronized in some way.

Consider an example:

```

a=10
{
    Read a
    a++
    write a
}
```

If P1 and P2 are different process and sharing the above code section, then P1 & P2 executes one after other P1 prints **a** as **11** and P2 prints **a** as **12**. Assume if P1 and P2 are executing concurrently, then outcome of execution depends on order in which the access take place, which is a race condition and leads to data inconsistency.



### 2.3.3 Critical-section problem:

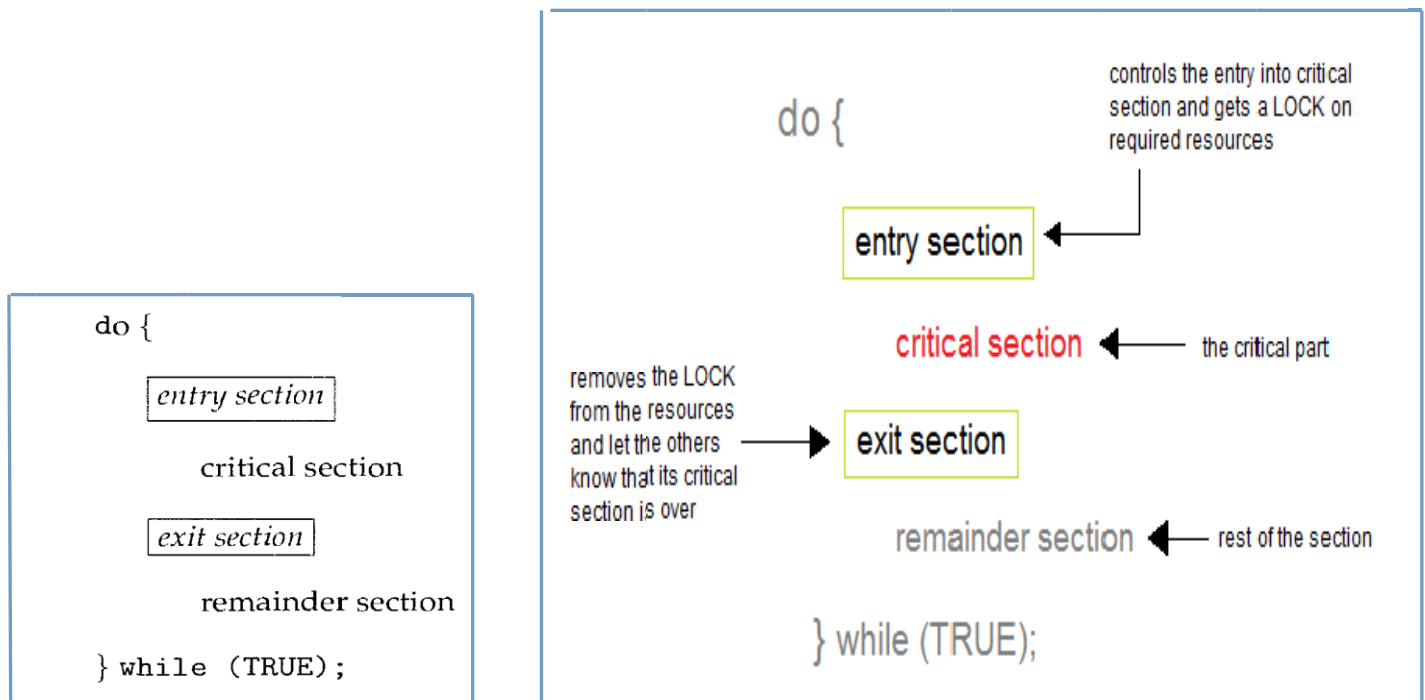
Critical section is a code /part of a program which contain access to shared variables and has to be executed as an atomic action.

**Example:** Consider the air line ticket reservation system, if only one sit is available and if 3 persons are booking a ticket at a time, only one should be allowed to block a sit, else it leads to data inconsistency. So the code which allows for updation of the seat variable(shared data) is the critical section.

Consider a system consisting of  $n$  processes  $\{P_0, P_1 \dots P_{n-1}\}$ . Each process has a segment of code, called a **critical section**, in which the process may be changing the common variables, updating a table, writing a file, and so on. The important feature of the system is that, when one process is executing in its critical section, no other process is to be allowed to execute in its critical section. That is, no two processes are executing in their critical sections at the same time.

The **critical-section problem** is to design a protocol that the processes can use to cooperate. Each process must request permission to enter its **critical section**. The section of code implementing this request is the **entry section**. The critical section may be followed by an **exit section**. The remaining code is the **remainder section**.

The **general structure** of a typical process  $P_i$ , is shown in the following figure



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

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 some processes wish to enter their critical sections, then only those processes that are not executing in their remainder sections can participate in the decision on which will enter its critical section next, and this selection cannot be postponed indefinitely.

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

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

1. Pre-emptive kernel
2. Non preemptive kernel.

A **pre-emptive kernel** allows a process to be pre-empted while it is running in kernel mode. It is **difficult to design** for SMP architectures, since in these environments it is possible for two kernel-mode processes to run simultaneously on different processors. It is more **suitable** for **real-time programming**, as it will allow a real-time process to pre-empt a process currently running in the kernel. Also, pre-emptive kernel may be **more responsive**, since there is less risk that a kernel-mode process will run for an arbitrarily long period before giving up the processor to waiting processes.

A **Nonpreemptive kernel** does not allow a process running in kernel mode to be pre-empted; a kernel-mode process will run until it exits kernel mode, blocks, or voluntarily yields control of the CPU. It is essentially free from **race conditions** on kernel data structures, as only one process is active in the kernel at a time.

#### 2.3.4 Peterson's solution

- A classic **software-based solution** to the critical-section problem is known as **Peterson's solution**. Peterson's solution is restricted to **two processes** that alternate execution between their critical sections and remainder sections.
- The processes are numbered  $P_0$  and  $P_1$  or  $P_i$  and  $P_j$  where  $j=1-i$ .  
Peterson's solution requires **two data items** to be shared between the two processes,
  - `int turn;`
  - `boolean flag[2];`
- The variable **turn** indicates whose turn it is to enter its critical section. That is, if  $turn = i$ , then process  $P_i$  is allowed to execute in its critical section. The **flag array** is used to indicate if a process is ready to enter its critical section. **For example**, if `flag[i]` is true, this value indicates that  $P_i$  is ready to enter its critical section.
- To enter the critical section, process  $P_i$  first sets **flag[i]** to be true and then sets **turn** to the value  $j$ , thereby asserting that if the other process wishes to enter the critical section, it can do so.
- If both processes try to enter at the same time, **turn** will be set to both  $i$  and  $j$  at roughly the same time. Only one of these assignments will last; the other will occur but will be overwritten immediately.
- The eventual value of **turn** decides which of the two processes is allowed to enter its critical section first.
- The following algorithm describes the structure of  $P_i$  in Peterson's solution.

```
do {
```

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

```
    critical section
```

```
    flag[i] = FALSE;
```

```
    remainder section
```

```
} while (TRUE);
```

- To prove that this solution is correct, we need to show that,
  1. Mutual exclusion is preserved.
  2. The progress requirement is satisfied.
  3. The bounded-waiting requirement is met.
- To prove **property 1**, we note that each  $P_i$  enters its critical section only if **either**  $\text{flag}[j] = \text{false}$  or  $\text{turn} = i$ . For  $P_0$  to enter, turn must be equal to 0 and for  $P_1$  to enter, turn must be equal to 1 because  $\text{flag}[0] = \text{flag}[1] = \text{true}$ . Since the value of turn can be either 0 or 1 but cannot be both, hence  $P_0$  and  $P_1$  cannot enter into critical section simultaneously.
- To prove **properties 2 and 3**, we note that a process  $P_i$  can be prevented from entering the critical section only if it is stuck in the while loop with the condition  $\text{flag}[j] = \text{true}$  and  $\text{turn} = j$ . If  $P_j$  is not ready to enter the critical section, then  $\text{flag}[j] = \text{false}$ , and  $P_i$  can enter its critical section. If  $\text{turn} = j$ , then  $P_j$  will enter the critical section. However, once  $P_j$  exits its critical section, it will reset  $\text{flag}[j]$  to false, allowing  $P_i$  to enter its critical section. If  $P_j$  resets  $\text{flag}[j]$  to true, it must also set turn to  $i$ . Thus,  $P_i$  will enter the critical section (progress) after at most one entry by  $P_i$  (bounded waiting).

### 2.3.5 Synchronization hardware:

- Any solution to the critical-section problem requires a simple tool called a **lock**. **Race conditions** are prevented by protecting the critical regions by locks. That is, a process must acquire a lock before entering a critical section; it releases the lock when it exits the critical section.
  - Lock before entering critical section and before accessing shared data
  - Unlock when leaving, after accessing shared data
  - Wait if locked
    - Important idea: all synchronization involves waiting
    - Should sleep if waiting for a long time

This is illustrated below,

```
do {
    acquire lock

    critical section

    release lock

    remainder section
} while (TRUE);
```

- Hardware features can make any programming task easier and improve **system efficiency**.
- The critical-section problem can be solved simply in a uniprocessor environment if we could

prevent interrupts from occurring while a shared variable was being modified. This solution is **not feasible** in a multiprocessor environment. Disabling interrupts on a multiprocessor can be **time consuming**, as the message is passed to all the processors. This message passing delays entry into each critical section, and **system efficiency decreases**.

- Many modern computer systems therefore provide special hardware instructions that allow us either to **test** and **modify** the content of a word or to **swap** the contents of two words **atomically** that is, as one uninterruptible unit.
- We can use these special instructions to solve the critical-section problem in a relatively simple manner. The TestAndSet() instruction can be defined as below,

```
boolean TestAndSet(boolean *target)    // This code returns false for first time i.e rv
{                                       //value and sets lock as true for next
    boolean rv = *target;              //execution
    *target = TRUE;
    return rv;
}
```

- The important characteristic is that this instruction is executed atomically. Thus, if two TestAndSet() instructions are executed simultaneously (each on a different CPU), they will be executed sequentially in some arbitrary order.
- The structure of process  $P_i$  is shown below,

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

    // critical section
```

```

        lock = FALSE;           //sets lock as False so allows next to
                                //take chance
    // remainder section

}while (TRUE);

```

- The Swap() instruction, in contrast to the TestAndSet() instruction, operates on the contents of two words as shown below,

```

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

```

- It is executed atomically. If the machine supports the Swap() instruction, then mutual exclusion can be provided as follows.
- A global Boolean variable **lock** is declared and is initialized to false. In addition, each process has a local Boolean variable **key**. The structure of process  $P_i$  is shown below,

```

do {
    key = TRUE;
    while (key == TRUE)
        Swap(&lock, &key);
};

// critical section

lock = FALSE;

// remainder section
}while (TRUE);

```

- Although these algorithms satisfy the mutual-exclusion requirement, they do not satisfy the bounded-waiting requirement.
- Another algorithm is given below using the TestAndSet() instruction that satisfies all the critical-section requirements. The common **data structures** are,

```

boolean waiting[n];
boolean lock;

```

- These data structures are initialized to **false**.

```

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

```

```

// critical section

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

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

// remainder section
} while (TRUE);

```

### 2.3.6 Semaphores:

The **hardware-based solutions** to the critical-section problem are **complicated** for application programmers to use. To **overcome this difficulty**, we can use a **synchronization tool** called a **semaphore**.

Dijkstra's worked on semaphores established over 30 years ago the foundation of modern techniques for accomplishing synchronization

- A semaphore, S, is a integer variable that is changed or tested only by one of the two following indivisible operations

**P(S):** while(S<=0)

*Dono-operation;*

S--;

**V(S):** S++;

- In Dijkstra's original paper, the P operation was an abbreviation for the Dutch word Proberen, meaning "to test" and the V operation was an abbreviation for the word Verhogen, meaning "to increment"

Now, P () and V() is normally called wait() and signal()

A **semaphore S** is an integer variable it is accessed only through **two standard atomic operations: wait() and signal()**. The wait() operation was termed as **P**; signal() was called **V**. The definition of **wait()** is as follows,

```

wait(S)
{
    (while S<= 0) ; // no-operation
    S--;
}

```

- ✓ The definition of **signal()** is as follows,

```

signal(S)
{
    S++ ;
}

```

- ✓ The wait() and signal() operations must be executed indivisibly. That is, when one process modifies the semaphore value, no other process can simultaneously modify that same semaphore value.

### Usage:

Operating systems often distinguish between **counting** and **binary** semaphores. The value of a **counting semaphore** can range over an unrestricted domain. The value of a **binary semaphore** can range only between 0 and 1.

**Binary semaphores** are known as **mutex locks**, as they are locks that provide mutual exclusion. We can use binary semaphores to deal with the critical-section problem for multiple processes. The  $n$  processes share a semaphore, **mutex**, initialized to 1. Each process  $P_i$  is organized as shown,

```
do {
    waiting(mutex);

    // critical section

    signal (mutex) ,

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

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

We can also use semaphores to solve various synchronization problems. For example, consider two concurrently running processes:  $P_1$  with a statement  $S_1$  and  $P_2$  with a statement  $S_2$ . Suppose we require that  $P_2$  be executed only after  $S_1$  has completed. We can implement this by letting  $P_1$  and  $P_2$  to share a common semaphore **synch**, initialized to 0, and by inserting the statements

```
S1;
signal(synch);
```

in process  $P_1$ , and the statements

```
wait(synch);
S2;
```

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

### **Disadvantage:**

- The main **disadvantage** of the semaphore definition given here is that it requires **busy waiting**. While a process is in its critical section, any other process that tries to enter its critical section must loop continuously in the entry code.

- This type of semaphore is also called a **spin lock** because the process "spins" while waiting for the lock.

### **Implementation :**

To overcome the need for busy waiting, we can **modify** the definition of the **wait()** and **signal()** semaphore operations. When a process executes the wait() operation and finds that the semaphore value is not positive, then rather than engaging in busy waiting, the process can **block** itself. The block operation places a process into a **waiting queue** associated with the semaphore, and the state of the process is switched to the **waiting state**.

A process that is blocked, waiting on a semaphore S, should be **restarted** when some other process executes a signal() operation. The process is restarted by a **wakeup()** operation, which changes the process from the waiting state to the ready state. The process is then placed in the ready queue. To implement semaphores under this definition, we define a semaphore as a 'C' struct:

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

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

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

The signal () semaphore operation can now be defined as,

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

The block() operation suspends the process that invokes it. The wakeup(P) operation resumes the execution of a blocked process P. These two operations are provided by the operating system as basic system calls. This implementation, semaphore values may be negative, if a semaphore value is negative, its magnitude is the number of processes waiting on that semaphore.

The list of waiting processes can be easily implemented by a link field in each process control block (PCB). Each semaphore contains an integer value and a pointer to a list of PCBs.



**Deadlocks and Starvation**

- The implementation of a semaphore with a waiting queue may result in a **deadlock** situation where two or more processes are waiting indefinitely for an event that can be caused only by one of the waiting processes. When such a state is reached, these processes are said to be **deadlocked**.
  - ✓ To illustrate this, we consider a system consisting of two processes,  $P_0$  and  $P_1$ , each accessing two semaphores,  $S$  and  $Q$ , set to the value 1:

$P_0$	$P_1$
wait(S);	wait(Q);
wait(Q);	wait(S);
.	.
.	.
.	.
signal(S);	signal(Q);
signal(Q);	signal(S);

- Suppose that  $P_0$  executes wait(S) and then  $P_1$  executes wait(Q). When  $P_0$  executes wait(Q), it must wait until  $P_1$  executes signal(Q). Similarly, when  $P_1$  executes wait(S), it must wait until  $P_0$  executes signal(S). Since these signal() operations cannot be executed,  $P_0$  and  $P_1$  are deadlocked.
- Another problem related to deadlocks is **indefinite blocking, or starvation**, a situation in which processes wait indefinitely within the semaphore. Indefinite blocking may occur if we remove processes from the list associated with a semaphore in LIFO (last-in, first-out) order.

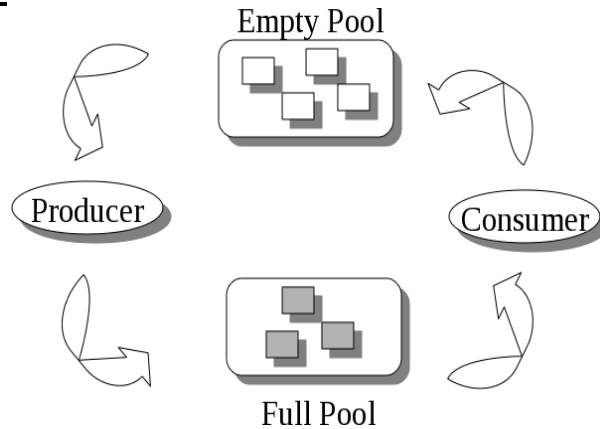
**Priority Inversion:**

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

As an example, assume we have three processes— $L$ ,  $M$ , and  $H$ —whose priorities follow the order  $L < M < H$ . Assume that process  $H$  requires resource  $R$ , which is currently being accessed by process  $L$ . Ordinarily, process  $H$  would wait for  $L$  to finish using resource  $R$ . However, now suppose that process  $M$  becomes runnable, thereby preempting process  $L$ . Indirectly, a process with a lower priority—process  $M$ —has affected how long process  $H$  must wait for  $L$  to relinquish resource  $R$ . This problem is known as **priority inversion**. It occurs only in systems with more than two priorities, so one solution is to have only two priorities. That is insufficient for most general-purpose operating systems, however. Typically these systems solve the problem by implementing a **priority-inheritance protocol**. According to this protocol, all processes that are accessing resources needed by a higher-priority process inherit the higher priority until they are finished with the resources in question. When they are finished, their priorities revert to their original values. In the example above, a priority-inheritance protocol would allow process  $L$  to temporarily inherit the priority of process  $H$ , thereby preventing process  $M$  from preempting its execution. When process  $L$  had finished using resource  $R$ , it would relinquish its inherited priority from  $H$  and assume its original priority. Because resource  $R$  would now be available, process  $H$ —not  $M$ —would run next.

**2.3.7 Classic problems of Synchronization**

- Bounded-buffer problem
- The Readers-Writers Problem
- The Dining-Philosophers Problem

**The Bounded-Buffer Problem:**

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

The code for the producer process is shown:

```
do{
    // produce an item in nextp
    ...

    wait(empty); // checks for empty buffer, if no then waits
    wait (mutex); // if empty buffer is available, then will fill it
    ...
    // add nextp to buffer
    ...
    signal(mutex); //releases the buffer
    signal (full); // increments full array
}while (TRUE);
```

The code for the consumer process is shown:

```
do
{
    wait (full); // if full is 0, it waits
    wait(mutex); //if full is non -zero , then reads
    ...
    //remove an item from buffer to nextc
    ...
```

```

        signal(mutex);
        signal(empty);
        //consume the item in nextc
    ...

}while (TRUE);

```

### **The Readers-Writers Problem:**

A database is to be shared among several concurrent processes. Some of these processes may want only to read the database, whereas others may want to update (that is, to read and write) the database.

- ✓ **Readers** - processes that only read the database.
- ✓ **Writers** - processes performing both read and write (update).
- ✓ **Problem:** If two readers access the shared data simultaneously, no problem will result. But if a writer and some other thread (either a reader or a writer) access the database simultaneously, problem arises.

This synchronization problem is referred to as the readers-writers problem.

The readers-writers problem has several variations. The **first** readers-writers problem, which requires that no reader must be kept waiting unless a writer has already obtained permission to use the shared object. The **second** readers-writers problem requires that, once a writer is ready, that writer performs its write as soon as possible.

A solution to either problem may result in starvation. In the first case, writers may starve; in the second case, readers may starve. In the solution to the first readers-writers problem, the reader processes share the following data structures:

```

semaphore mutex, wrt;
int readcount;

```

- ✓ The semaphores **mutex** and **wrt** are initialized to 1 and **readcount** is initialized to 0.
- ✓ The semaphore **wrt** is common to both reader and writer processes.
- ✓ The **mutex** semaphore is used to ensure **mutual exclusion** when the variable **readcount** is updated.
- ✓ The **readcount** variable keeps track of how many processes are currently reading the object.
- ✓ The semaphore **wrt** functions as a mutual-exclusion semaphore for the writers.
- ✓ The code for a writer process is,

```

do {
    wait(wrt);
    ....
    // writing is performed
    signal(wrt);
} while (TRUE);

```

- ✓ The code for a reader process is shown,

```

do {
    wait (mutex);
    readcount++;
    if (readcount == 1)    // First reader must

```

```

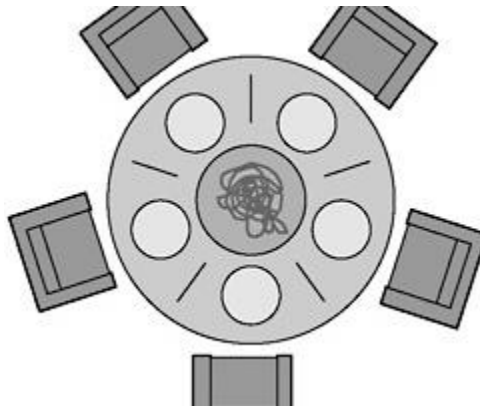
        wait (wrt);        // block the writers
    signal(mutex);
    .....
    // reading is performed
    .....
    wait(mutex);
    readcount--;
    if (readcount== 0)      //Last reader must
        signal(wrt);      //unblock the writers
    signal(mutex);
} while (TRUE);

```

- ✓ If a writer is in the critical section and n readers are waiting, then one reader is queued on **wrt**, and n-1 readers are queued on **mutex**.
- ✓ Also, when a writer executes **signal(wrt)**, we may resume the execution of either the waiting readers or a single waiting writer. The selection is made by the scheduler.

### **The Dining-Philosophers Problem:**

- ✓ Consider **five** philosophers who spend their lives thinking and eating.
- ✓ The philosophers share a circular table surrounded by five chairs, each belonging to one philosopher.
- ✓ In the centre of the table is a bowl of noodles, and the table is laid with five single chopsticks (below figure).
- ✓ When a philosopher is thinking, he does not interact with her colleagues.
- ✓ When a philosopher gets hungry, he tries to pick up the two chopsticks that are closest to him (the chopsticks that are between his left and right neighbors).
- ✓ When a hungry philosopher has both her chopsticks at the same time, she eats without releasing her chopsticks. When she is finished eating, she puts down both of her chopsticks and starts thinking again.



**The situation of the dining philosophers**

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

**semaphore chopstick[5];**  
 where all the elements of chopstick are initialized to 1.

- ✓ The structure of philosopher i is shown below,

```

do {
    wait(chopstick[i]);
    k[(i+1) % 5];
    .....
    // eat
    .....

    signal(chopstick[i]);
    signal(chopstick[(i+1) % 5]);
    .....
    // think
    .....
} while (TRUE);

```

- ✓ The disadvantage is it can create a deadlock. Suppose that all five philosophers become hungry simultaneously and each grabs left chopstick. All the elements of chopstick will now be equal to 0. When each philosopher tries to grab right chopstick, he will be delayed forever.
- ✓ Several possible remedies to the deadlock problem are available.
  - Allow at most four philosophers to be sitting simultaneously at the table.
  - Allow a philosopher to pick up chopsticks only if both chopsticks are available (to do this, he must pick them up in a critical section).
  - Use an asymmetric solution; that is, an odd philosopher picks up first left chopstick and then right chopstick, whereas an even philosopher picks up right chopstick and then left chopstick.

### 2.3.8 Monitors:

- ✓ All processes share a semaphore variable mutex, which is initialized to 1. Each process must execute wait(mutex) before entering the critical section and signal(mutex) afterward. If this sequence is not observed, two processes may be in their critical sections simultaneously.
- ✓ This may result in various difficulties.
  - Suppose that a process interchanges the order in which the wait() and signal() operations on the semaphore mutex are executed, resulting in the following execution:

```

signal(mutex);
.....
critical section
.....
wait(mutex);

```

- ✓ In this situation, several processes may be executing in their critical sections simultaneously, violating the mutual-exclusion requirement. This error may be discovered only if several processes are simultaneously active in their critical sections.
  - Suppose that a process replaces signal (mutex) with wait (mutex). That is, it executes
 

```
wait(mutex);
.....
```

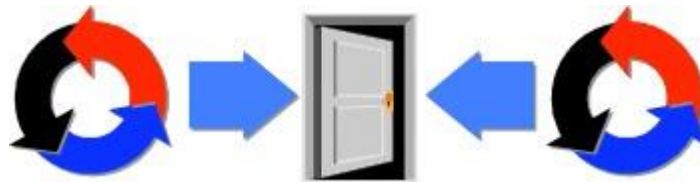
**critical section**

.....

**wait(mutex);**

In this case, a deadlock will occur.

- Suppose that a process omits the wait (mutex), or the signal (mutex), or both. In this case, either mutual exclusion is violated or a deadlock will occur.
- ✓ These **examples** illustrate that various types of errors can be generated easily when programmers use semaphores incorrectly to solve the critical-section problem. To deal with such errors, one fundamental high-level synchronization construct-the **monitor** is used.

**Usage:****Monitors**

- Monitors are a high-level synchronization primitive:

```

monitor monitor-name
{
    shared variable declarations

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

    {
        initialization code
    }
}

```

**Key points:**

- Shared variables are only accessible through the monitor's procedures.
- The language performs locking so that only one process can be running any monitor procedure at a time.

- ✓ Monitors encapsulate data structures that are not externally accessible ♦ Mutual exclusive access to data structure enforced by compiler or language run-time.
- ✓ An abstract data type encapsulates data with the set of functions to operate on that data are independent of any specific implementation of ADT.
- ✓ A monitor type is an ADT that includes a set of programmer-defined operations that are provided with mutual exclusion within a monitor. It declares the variables whose values define the state of an instance of that type as shown below:

```

monitor monitor name
{
    // shared variable declarations

```

```

procedure P1 ( . . . ) {
    .....
}
procedure P2 ( . . . ) {
    .....
}
.
.
procedure Pn ( . . . ) {
    .....
}
initialization code ( . . . ) {
    .....
}

```

- ✓ The functions defined within a monitor can access only those variables declared locally within the monitor and its formal parameters and the local variables of a monitor can be accessed by only the local functions.
- ✓ **Consider an example:**

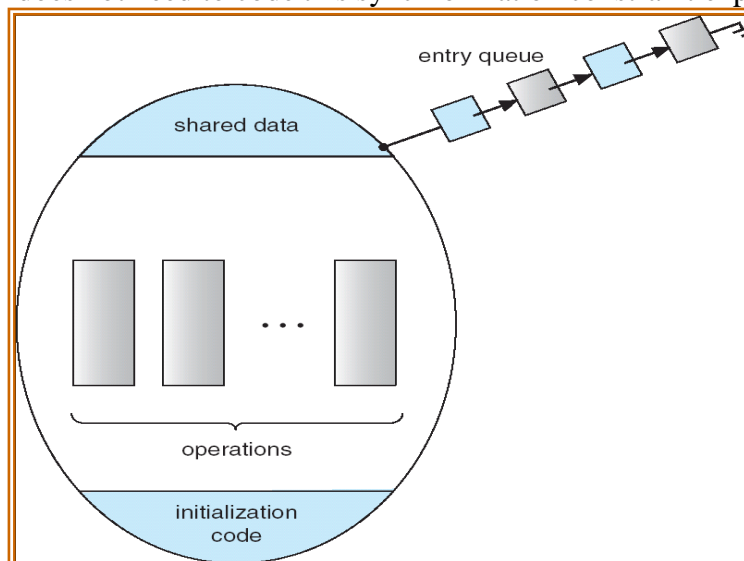
```

Monitor bank
{
    Condition balance;

    Withdraw(amount)
    {
        balance=balance-amount;
    }
}

```

- ✓ Balance variable can be updated only by calling withdraw method and monitors allows only one process to enter and ensures synchronization.
- ✓ The monitor construct ensures that only one process at a time is active within the monitor. The programmer does not need to code this synchronization constraint explicitly.

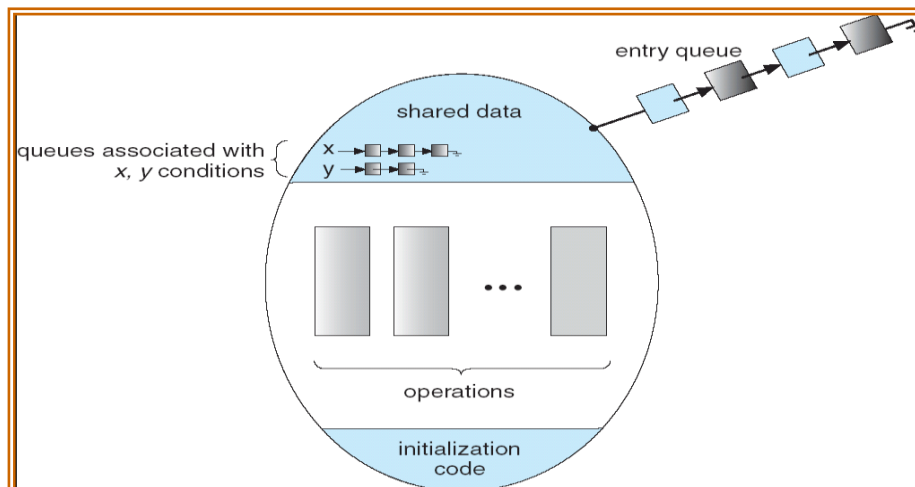


**Schematic view of a monitor**

- ✓ The monitor construct, as defined so far is not sufficiently powerful for modeling some synchronization schemes. For this purpose, we need to define additional synchronization mechanisms. These mechanisms are provided by the **condition** construct. A programmer can define one or more variables of the type **condition**,

**condition x, y;**

- ✓ For example in producer and consumer problem, Monitors allows only one process to enter ,but it will not check whether it is producer or consumer. If consumer enters and if there are empty buffers then , either it can read nor it allows other process (not even producer) to enter the monitor. So condition variables are used.
- ✓ The only operations that can be invoked on a condition variable are 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; that is, the state of x is the same as if the operation has never been executed. It is shown in below figure.



### Monitor with condition variables

- ✓ Condition Variables need wait and wakeup as in semaphores . Monitor uses Condition Variables Conceptually associated with some conditions Operations on condition variables:
  - wait(): suspends the calling thread and releases the monitor lock. When it resumes, reacquire the lock. Called when condition is not true
  - signal(): resumes one thread waiting in wait() if any. Called when condition becomes true and wants to wake up one waiting thread
- ✓ Suppose when the x. signal () operation is invoked by a process P, there exists a suspended process Q associated with condition x. If the suspended process Q is allowed to resume its execution, the signaling process P must wait. Otherwise, both P and Q would be active simultaneously within the monitor. However both processes can continue with their execution. Two possibilities exist,
  - **Signal and wait.** P either waits until Q leaves the monitor or waits for another condition.
  - **Signal and continue.** Q either waits until P leaves the monitor or waits for another condition.



**Dining-Philosophers Solution Using Monitors:**

- ✓ It is a deadlock-free solution to the dining-philosophers problem.
- ✓ This solution imposes the restriction that a philosopher may pick up her chopsticks only if both of them are available.
- ✓ We use the following data structure to distinguish among three states in which we may find a philosopher.

```
enum {thinking, hungry, eating}state [5];
```

- ✓ Philosopher i can set the variable **state[i] = eating** only if her two neighbors are not eating.
- ✓ We also need to declare
 

```
condition self [5];
```

 where philosopher i can delay herself when she is hungry but is unable to obtain the chopsticks she needs.
- ✓ The distribution of the chopsticks is controlled by the monitor **dp**, whose definition is shown below.

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

    void pickup(int i)
    {
        state[i] =HUNGRY;           //set philosopher as hungry
        test(i);                    //whether he gets both the chopsticks

        if (state [i] != EATING)    // after eating set philosopher to
                                   //wait until he is hungry & chopsticks
                                   //are available

            self [i] . wait() ;
    }

    void putdown(int i)
    {
        state[i] =THINKING;         //set philosopher to thinking
        test((i + 4) % 5);          //Test whether the neighbor philosophers
                                   //are hungry and waiting for chopsticks
        test((i + 1) % 5);
    }
}
```

//test function first checks whether both the chopsticks are available for the philosopher. If yes allows him //to eat by setting self[i].signal[self is a condition variable which allows to continue eating].

```
void test(int i)
{
```

```

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

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

```

- ✓ Each philosopher, before starting to eat, must invoke the operation pickup(). After the successful completion of the operation, the philosopher may eat. Following this, the philosopher invokes the putdown() operation. Thus, philosopher i must invoke the operations pickup() and putdown() in the following sequence,

```

dp.pickup(i);
...
eat
...
dp.putdown(i);

```

### Implementing a Monitor Using Semaphores:

- ✓ For each monitor, a semaphore **mutex** (initialized to 1) is provided.
- ✓ A process must execute **wait(mutex)** before entering the monitor and must execute **signal(mutex)** after leaving the monitor.
- ✓ Since a signaling process must wait until the resumed process either leaves or waits, an additional semaphore, **next** (initialized to 0) is introduced, on which the signaling processes may suspend themselves.
- ✓ An integer variable **next\_count** is also provided to count the number of processes suspended on **next**. Thus, each external procedure **F** is replaced by ,

```

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

```

Mutual exclusion within a monitor is ensured.

- ✓ Condition variables are implemented as follows. For each condition x, semaphore **x\_sem** and an integer variable **x\_count**, both initialized to 0 are introduced.
- ✓ The operation x.wait() can now be implemented as

```

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

```

- ✓ The operation x.signal() can be implemented as

```

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

```

### Resuming Processes within a Monitor

- ✓ If several processes are suspended on condition x, and an x.signal() operation is executed by some process, then the problem is how to determine which of the suspended processes should be resumed next.
- ✓ One simple solution is to use an FCFS ordering, so that the process that has been waiting for the longest time is resumed first.
- ✓ In many circumstances, such simple scheduling scheme is not adequate. For this purpose, the **conditional-wait** construct can be used; it has the form

x.wait(c);

where c is an integer expression that is evaluated when the wait() operation is executed. The value of c, which is called a **priority number** is then stored with the name of the process that is suspended.

- ✓ When x.signal() is executed, the process with the smallest priority number is resumed next.
- ✓ To illustrate this, consider the **ResourceAllocator** monitor shown below, which controls the **allocation of a single resource** among competing processes.
- ✓

```

monitor ResourceAllocator
{
    boolean busy;
    condition x;

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

```

```

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

initialization_code()
{
    busy = FALSE;
}
}

```

- ✓ Each process, when requesting an allocation of this resource, it specifies the maximum time it plans to use the resource.
- ✓ The monitor allocates the resource to the process that has the shortest time-allocation request. A process that needs to access the resource in question must observe the following sequence:

```

R.acquire(t);
.....
    access the resource;
.....
R. release() ;

```

where R is an instance of type ResourceAllocator.

- ✓ The monitor concept cannot guarantee that the preceding access sequence will be observed and the following problems can occur,
  - A process might access a resource without first gaining access permission to the resource.
  - A process might never release a resource once it has been granted access to the resource.
  - A process might attempt to release a resource that it never requested.
  - A process might request the same resource twice (without first releasing the resource).
- ✓ The same difficulties are encountered with the use of semaphores.

### **Solution to the Producer-Consumer problem using Monitors**

Monitors make solving the producer-consumer a little easier. Mutual exclusion is achieved by placing the critical section of a program inside a monitor. In the code below, the critical sections of the producer and consumer are inside the monitor *Producer Consumer*. Once inside the monitor, a process is blocked by the **Wait** and **Signal** primitives if it cannot continue.

#### **Two condition variables**

- has\_empty: buffer has at least one empty slot
- has\_full: buffer has at least one full slot

**nfull:** number of filled slots need to do our own counting for condition variables

**monitor ProducerConsumer**

```
{
    int nfull = 0;
    condition has_empty, has_full;
    producer()
    {
        if (nfull == N)
            wait (has_empty); ... // fill a slot
            ++ nfull;
            signal (has_full);
    }
    consumer()
    {
        if (nfull == 0)
            wait (has_full); ... // empty a slot
            -- nfull;
            signal (has_empty);
    }
}
```