

## Module 3. DEADLOCK AND STARVATION

### 3.1 PRINCIPLES OF DEADLOCK

Deadlock can be defined as the *permanent* blocking of a set of processes that either compete for system resources or communicate with each other. A set of processes is deadlocked when each process in the set is blocked awaiting an event (typically the freeing up of some requested resource) that can only be triggered by another blocked process in the set. Deadlock is permanent because none of the events is ever triggered.

Under the normal mode of operation, a process may utilize a resource in only the following sequence:

1. **Request:** If the request cannot be granted immediately (for example, the resource is being used by another process), then the requesting process must wait until it can acquire the resource.
2. **Use:** The process can operate on the resource.
3. **Release:** The process releases the resource.

But, if every process in the set is waiting for other processes to release the resource, then the deadlock happens. Two general categories of resources:

- **Reusable:** can be safely used by only one process at a time and *is not depleted* (that is not reduced) by that use. For example, Processors, I/O channels, main and secondary memory, devices, and data structures such as files, databases, and semaphores. Here, deadlock occurs if each process holds one resource and requests the other.
- **Consumable:** one that can be created (*produced*) and destroyed (*consumed*). For example, interrupts, signals, messages, and information in I/O buffers. Here, deadlock may occur if a Receive message is blocking.

#### 3.1.1 Conditions (or Characterization) for Deadlock

Following are the conditions for deadlock to present:

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

#### 3.1.2 Resource Allocation Graphs

The resource allocation graph is a directed graph that depicts a state of the system of resources and processes with each process and each resource represented by a node. It is a graph consisting of a set of vertices  $V$  and a set of edges  $E$  with following notations:

- $V$  is partitioned into two types:

- $P = \{P_1, P_2, \dots, P_n\}$ , the set consisting of all the processes in the system.
- $R = \{R_1, R_2, \dots, R_m\}$ , the set consisting of all resource types in the system.
- **Request edge:** A directed edge  $P_i \rightarrow R_j$  indicates that the process  $P_i$  has requested for an instance of the resource  $R_j$  and is currently waiting for that resource.
- **Assignment edge:** A directed edge  $R_j \rightarrow P_i$  indicates that an instance of the resource  $R_j$  has been allocated to the process  $P_i$ .

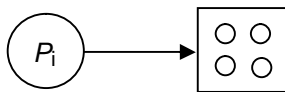
The following symbols are used while creating resource allocation graph:



A Process



A resource with 4 instances

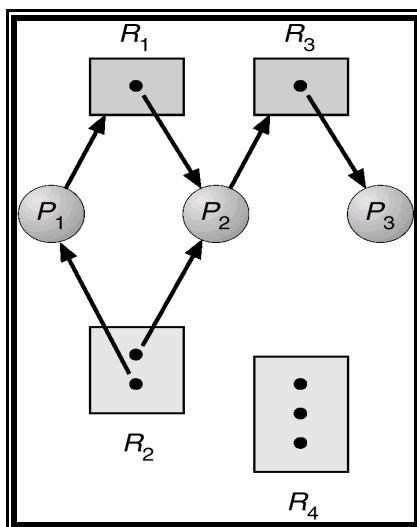


Process  $P_i$  requests for  $R_j$

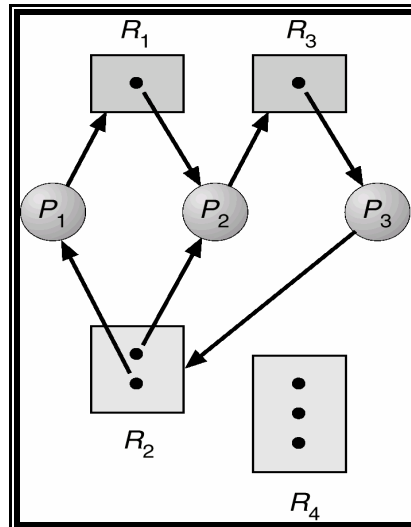


Process  $P_i$  is holding an instance of  $R_j$

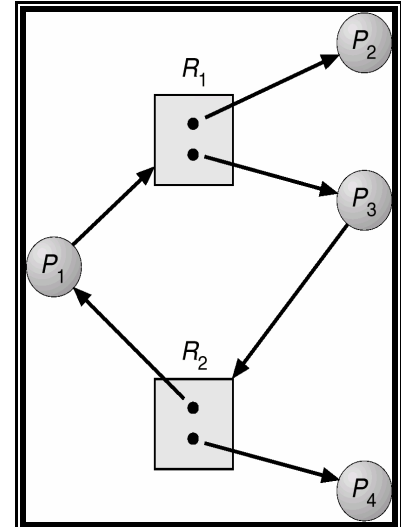
Examples of resource allocation graph are shown in Figure 3.1. Note that, in Figure 3.1(c), the processes  $P_2$  and  $P_4$  are not depending on any other resources. And, they will give up the resources  $R_1$  and  $R_2$  once they complete the execution. Hence, there will not be any deadlock.



(a) Resource allocation Graph



(b) With a deadlock



(c) with cycle but no deadlock

Figure 3.1 Resource allocation graphs

Given the definition of resource allocation graph, we can understand that, if there is no cycle in the graph, then there will not be a deadlock. If there is a cycle, there is a chance of deadlock.

There are three general approaches exist for dealing with deadlock.

- Prevent deadlock: Ensure that the system will *never* enter a deadlock state.
- Avoid deadlock: Make appropriate dynamic choices based on the current state of resource allocation.
- Detect Deadlock: Allow the system to enter a deadlock state and then recover.

### 3.2 DEADLOCK PREVENTION

The strategy of deadlock prevention is to design a system in such way that the possibility of deadlock is excluded. This is possible if we ensure that one of the four conditions (mutual exclusion, hold & wait, No preemption and circular wait) cannot hold. We will examine these conditions now.

- **Mutual Exclusion:** The mutual-exclusion condition must hold for nonsharable resources. For example, a printer cannot be simultaneously shared by several processes. On the other hand, sharable resources do not require mutually exclusive access, and thus cannot be involved in a deadlock. For example, simultaneous access can be granted for read-only file. A process never needs to wait for a sharable resource. In general, we cannot prevent deadlocks by denying the mutual-exclusion condition: Some resources are intrinsically nonsharable.
- **Hold and Wait:** To ensure that the hold-and-wait condition never occurs in the system, we must guarantee that, whenever a process requests a resource, it does not hold any other resources. That is, a process should request all its required resources at one time and block that process until all its requests can be granted simultaneously. This approach has certain problems:
  - A process may be held up for a long time waiting for all its required resources. Actually, it would have proceeded with only few of the processes.
  - Resources allocated to one process may remain unused for some time, during which they are denied to other processes.
  - A process may not know in advance all the resources that it requires.
- **No Preemption:** This condition can be prevented in this way: If a process is holding some resources and requests another resource that cannot be immediately allocated to it (that is, the process must wait), then all resources currently being held are preempted. In other words, these resources are implicitly released. The preempted resources are added to the list of resources for which the process is waiting. The process will be restarted only when it can regain its old resources, as well as the new ones that it is requesting.
- **Circular Wait:** The circular-wait condition can be prevented by defining a linear ordering of resources. That is, every process must request the resources in an increasing order. Let there be two processes P1 and P2. Let there are resources R<sub>i</sub> and R<sub>j</sub> such that  $i < j$ . Now, P1 has acquired R<sub>i</sub> and requested R<sub>j</sub>. And, P2 has acquired R<sub>j</sub>

and requesting  $R_j$ . This condition is impossible, because it implies  $i < j$  and  $j < i$ . But, here also, the problems seen in hold-and-wait prevention are seen.

### 3.3 DEADLOCK AVOIDANCE

Deadlock prevention methods seen in the previous section requires at least one condition should not hold. But, these methods result in low device utilization and reduced system throughput.

An alternative method for avoiding deadlocks is to require additional information about how resources are to be requested. For this purpose, a simplest and most useful model is designed: each process must declare the maximum number of resources of each type that it may need. Given a priori information about the maximum number of resources of each type that may be requested for each process, it is possible to construct an algorithm that ensures that the system will never enter a deadlock state. This algorithm defines the deadlock-avoidance approach. A deadlock-avoidance algorithm dynamically examines the resource-allocation state to ensure that a circular-wait condition can never exist. The resource-allocation state is defined by the number of available and allocated resources, and the maximum demands of the processes.

#### 3.3.1 Safe State

A state is **safe** if the system can allocate resources to each process in some order and still avoid a deadlock. More formally, a system is in a safe state only if there exists a **safe sequence**. A Sequence  $\langle P_1, P_2, \dots, P_n \rangle$  is safe if for each  $P_i$ , the resources that  $P_i$  can still request can be satisfied by currently available resources + resources held by all the  $P_j$ , with  $j < i$ .

- If the needs of  $P_i$  are not immediately available, then  $P_i$  can wait until all  $P_j$  have finished.
- When  $P_j$  is finished,  $P_i$  can obtain needed resources, execute, return allocated resources, and terminate.
- When  $P_i$  terminates,  $P_{i+1}$  can obtain its needed resources, and so on.

If no such sequence exists, then the system state is said to be unsafe.

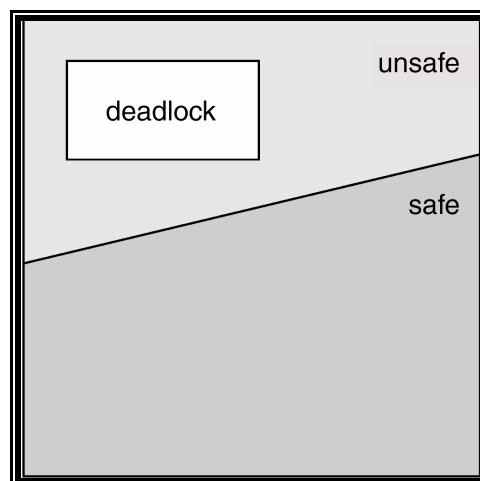


Figure 3.2 Safe, unsafe and deadlock state spaces

A safe state is not a deadlock state. A deadlock state is an unsafe state. But, not all unsafe states are deadlocks as shown in Figure 3.2. An unsafe state may lead to a deadlock. As long as the state is safe, the OS can avoid unsafe (and deadlock) states. In an unsafe state, the OS cannot prevent processes from requesting resources such that a deadlock occurs: The behavior of the processes controls unsafe states.

### 3.3.2 Resource Allocation Graph Algorithm

One of the techniques for avoiding a deadlock is using resource allocation graph with an additional edge called as **claim edge**. Claim edge  $P_i \rightarrow R_j$  indicated that process  $P_i$  may request resource  $R_j$  at some time in future. This edge is represented by a dashed line. The important steps are as below:

- When a process  $P_i$  requests a resource  $R_j$ , the claim edge  $P_i \rightarrow R_j$  is converted to a request edge.
- Similarly, when a resource  $R_j$  is released by the process  $P_i$ , the assignment edge  $R_j \rightarrow P_i$  is reconverted as claim edge  $P_i \rightarrow R_j$ .
- The request for  $R_j$  from  $P_i$  can be granted only if the converting request edge to assignment edge do not form a cycle in the resource allocation graph.

To apply this algorithm, each process  $P_i$  must know all its claims before it starts executing. If no cycle exists, then the allocation of the resource will leave the system in a safe state. If the cycle is found, system is put into unsafe state and may cause a deadlock.

An illustration: Consider a resource allocation graph shown in Figure 3.3(a). Suppose  $P_2$  requests  $R_2$ . Though  $R_2$  is currently free, we cannot allocate it to  $P_2$  as this action will create a cycle in the graph as shown in Figure 3.3(b). This cycle will indicate that the system is in unsafe state: because, if  $P_1$  requests  $R_2$  and  $P_2$  requests  $R_1$  later, a deadlock will occur.

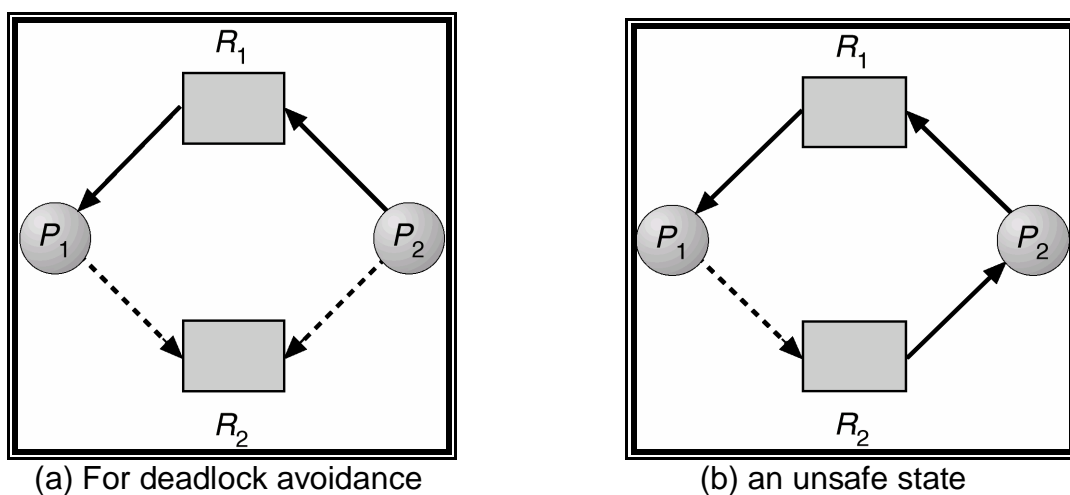


Figure 3.3 Resource Allocation graphs

### 3.3.3 Banker's Algorithm

The resource-allocation graph algorithm is not applicable when there are multiple instances for each resource. The banker's algorithm addresses this situation, but it is less efficient. The name was chosen because this algorithm could be used in a banking system to ensure that the bank never allocates its available cash such that it can no longer satisfy the needs of all its customers.

When a new process enters the system, it must declare the maximum number of instances of each resource type that it may need. This number may not exceed the total number of resources in the system. When a user requests a set of resources, the system must determine whether the allocation of these resources will leave the system in a safe state. If it will, the resources are allocated; otherwise, the process must wait until some other process releases enough resources.

Data structure for Banker's algorithms is as below –

Let  $n$  be the number of processes in the system and  $m$  be the number of resource types.

- **Available:** Vector of length  $m$  indicating number of available resources. If  $Available[j] = k$ , there are  $k$  instances of resource type  $R_j$  available.
- **Max:** An  $n \times m$  matrix defines the maximum demand of each process. If  $Max[i,j] = k$ , then process  $P_i$  may request at most  $k$  instances of resource type  $R_j$ .
- **Allocation:** An  $n \times m$  matrix defines the number of resources currently allocated to each process. If  $Allocation[i, j] = k$  then  $P_i$  is currently allocated  $k$  instances of  $R_j$ .
- **Need:** An  $n \times m$  matrix indicates remaining resource need of each process. If  $Need[i,j] = k$ , then  $P_i$  may need  $k$  more instances of  $R_j$  to complete its task. Note that,

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

The Banker's algorithm has two parts:

1. **Safety Algorithm:** It is for finding out whether a system is in safe state or not. The steps are as given below –
  1. Let  $Work$  and  $Finish$  be vectors of length  $m$  and  $n$ , respectively. Initialize:  
 $Work = Available$   
 $Finish [i] = false$  for  $i = 1, 2, 3, \dots, n$ .
  2. Find an  $i$  such that both:
    - (a)  $Finish [i] = false$
    - (b)  $Need_i \leq Work$If no such  $i$  exists, go to step 4.
  3.  $Work = Work + Allocation_i$   
 $Finish[i] = true$   
go to step 2.
  4. If  $Finish [i] == true$  for all  $i$ , then the system is in a safe state.
2. **Resource – Request Algorithm:** Let  $Request_i$  be the request vector for process  $P_i$ . If  $Request_i[j] = k$  then process  $P_i$  wants  $k$  instances of resource type  $R_j$ .

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

If the resulting resource allocation is safe, then the transaction is complete and the process  $P_i$  is allocated its resources. If the new state is unsafe, then  $P_i$  must wait for  $Request_i$ , and the old resource-allocation state is restored

### Example for Banker's algorithm:

Consider 5 processes  $P_0$  through  $P_4$  and 3 resources  $A$  (10 instances),  $B$  (5 instances), and  $C$  (7 instances). Snapshot at time  $T_0$  the snapshot of the system is as given in Table 3.1.

**Table 3.1 Snapshot of the system at time  $T_0$**

Process	Allocation (A B C)	Max (A B C)	Available (A B C)
P0	(0 1 0)	(7 5 3)	(3 3 2)
P1	(2 0 0)	(3 2 2)	
P2	(3 0 2)	(9 0 2)	
P3	(2 1 1)	(2 2 2)	
P4	(0 0 2)	(4 3 3)	

The matrix  $Need = Max - Allocation$ . It is given by the Table 3.2.

**Table 3.2 Need Matrix**

Process	Need (A B C)
P0	(7 4 3)
P1	(1 2 2)
P2	(6 0 0)
P3	(0 1 1)
P4	(4 3 1)

**Table 3.3 New State**

Process	Allocation (A B C)	Need (A B C)	Available (A B C)
P0	(0 1 0)	(7 4 3)	(2 3 0)
P1	(3 0 2)	(0 2 0)	
P2	(3 0 2)	(6 0 0)	
P3	(2 1 1)	(0 1 1)	
P4	(0 0 2)	(4 3 1)	

We can apply Safety algorithm to check whether the system is safe. We can find that the sequence  $\langle P1, P3, P4, P2, P0 \rangle$  is one of the safety sequences.

Suppose, now the process  $P1$  makes a request (1, 0, 2). To check whether this request can be immediately granted, we can apply Resource-Request algorithm. If we assume that this request is fulfilled, the new state would be as shown in Table 3.3. Now, by checking using safety algorithm, we see that the sequence  $\langle P1, P3, P4, P0, P2 \rangle$  is in safe state. Hence, this request can be granted.



### 3.4 DEADLOCK DETECTION

If a system does not make use of either a 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

Note that a detection-and-recovery scheme has some system overhead and run-time cost is more.

#### 3.4.1 Single Instance of Each Resource Type

If all resources have only a single instance, then we can define a deadlock detection algorithm that is similar to resource-allocation graph, called a **wait-for graph**. We obtain this graph from the resource-allocation graph by removing the resource-nodes and collapsing the appropriate edges. That is, an edge from  $P_i$  to  $P_j$  in a wait-for graph implies that process  $P_i$  is waiting for process  $P_j$  to release a resource that  $P_i$  needs.

Consider Figure 3.4 showing a resource allocation graph and its respective wait-for graph.

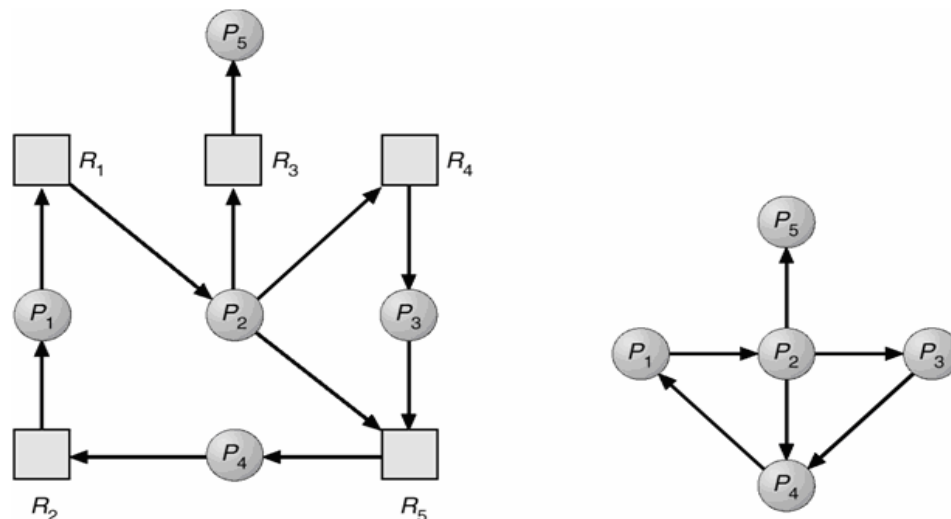


Figure 3.4 (a) Resource Allocation graph

(b) corresponding wait-for graph

A deadlock exists in the system if and only if the wait-for graph contains a cycle. To detect deadlocks, the system needs to maintain the wait-for graph and periodically to invoke an algorithm that searches for a cycle in the graph.

#### 3.4.2 Several Instances of a Resource Type

The wait-for graph scheme is not applicable to a resource-allocation system with multiple instances of each resource type. If resources have many instances, we use another algorithm which is similar to Banker's algorithm. The data structures used are:



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

The **detection algorithm** given below investigates every possible allocation sequence for the processes that remain to be completed.

1. Let  $Work$  and  $Finish$  be vectors of length  $m$  and  $n$ , respectively. Initialize:
  - (a)  $Work = Available$
  - (b) For  $i = 1, 2, \dots, n$ 
    - if  $Allocation_i \neq 0$ , then
    - $Finish[i] = false$ ;
    - Else
    - $Finish[i] = true$ .
2. Find an index  $i$  such that both:
  - (a)  $Finish[i] = false$
  - (b)  $Request_i \leq Work$
 If no such  $i$  exists, go to step 4.
3.  $Work = Work + Allocation_i$   
 $Finish[i] = true$   
 go to step 2.
4. If  $Finish[i] == false$ , for some  $i$ ,  $1 \leq i \leq n$ , then the system is in deadlock state. Moreover, if  $Finish[i] == false$ , then  $P_i$  is deadlocked.

### 3.4.3 Detection Algorithm Usage

The detection algorithm should be invoked based on following factors:

- How often is a deadlock likely to occur?
- How many processes will be affected by deadlock when it happens?

If deadlocks occur frequently, then the detection algorithm should be invoked frequently. Resources allocated to deadlocked processes will be idle until the deadlock can be broken.

Deadlocks occur only when some process makes a request that cannot be granted immediately. So, we could invoke the deadlock detection algorithm every time a request for allocation cannot be granted immediately. In this case, we can identify not only the set of processes that is deadlocked, but also the specific process that "caused" the deadlock. Another alternative is to invoke the algorithm in periodic intervals, say, once in an hour or whenever CPU utilization drops below certain level.

## 3.5 RECOVERY FROM DEADLOCK

When a detection algorithm determines that a deadlock exists, several alternatives are available. One possibility is to inform the operator that a deadlock has occurred, and to let the operator deal with the deadlock manually. The other possibility is to let the system recover from the deadlock automatically. There are two options for breaking a deadlock. One solution is simply to abort one or more processes to break the circular wait. The second option is to preempt some resources from one or more of the deadlocked processes.

### 3.5.1 Process Termination

Processes can be aborted for eliminating deadlock in two different ways:

- **Abort all deadlocked processes:** This method clearly will break the deadlock cycle. But, these processes may have computed for a long time, and the results of these partial computations must be discarded and probably recomputed later.
- **Abort one process at a time until the deadlock cycle is eliminated:** This method incurs considerable 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 may not be easy. If the process was in the midst of updating a file, terminating it will leave that file in an incorrect state. Similarly, if the process was in the midst of printing data on the printer, the system must reset the printer to a correct state before printing the next job. Many factors may determine which process is chosen to abort:

- Priority of the process.
- How long process has computed, and how much longer to completion.
- How many and what type of resources the process has used.
- How many more resources process needs to complete.
- How many processes will need to be terminated?
- Is process interactive or batch?

### 3.5.2 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. Following issues have to be considered:

- **Selecting a victim:** Which resources and which processes are to be preempted? We must determine the order of pre-emption to minimize cost. Cost factors may include parameters like the number of resources a deadlock process is holding, and the amount of time a deadlocked process has thus far consumed during its execution.
- **Rollback:** If we preempt a resource from a process, it cannot continue its normal execution and hence we must roll back the process as far as necessary to break the deadlock. This method requires the system to keep more information about the state of all the running processes.
- **Starvation:** In a system where victim selection is based primarily on cost factors, it may happen that the same process is always picked as a victim. As a result, this process

never completes its designated task, a starvation situation that needs to be dealt with in any practical system. Clearly, we must ensure that a process can be picked as a victim only a (small) finite number of times.

### 3.6 AN INTEGRATED DEADLOCK STRATEGY

There are strengths and weaknesses to all of the strategies for dealing with deadlock. Rather than attempting to design an OS facility that employs only one of these strategies, it might be more efficient to use different strategies in different situations. One of the approaches could be –

- Group resources into a number of different resource classes.
- Use the linear ordering strategy for the prevention of circular wait to prevent deadlocks between resource classes.
- Within a resource class, use the algorithm that is most appropriate for that class.

The resource classes can be

- **Swappable space:** Blocks of memory on secondary storage for use in swapping processes
- **Process resources:** Assignable devices, such as tape drives, and files
- **Main memory:** Assignable to processes in pages or segments
- **Internal resources:** Such as I/O channels

### 3.7 DINING PHILOSOPHERS PROBLEM

The dining philosopher's problem can be explained as below: Five philosophers live in a house, where a table is laid for them. The life of each philosopher consists of thinking and eating. The food they eat is spaghetti. Each philosopher requires two forks to eat spaghetti. The eating arrangements are simple as shown in Figure 3.5. There is a round table on which is set a large serving bowl of spaghetti, five plates, one for each philosopher, and five forks. A philosopher wishing to eat goes to his or her assigned place at the table and, using the two forks on either side of the plate, takes and eats some spaghetti.

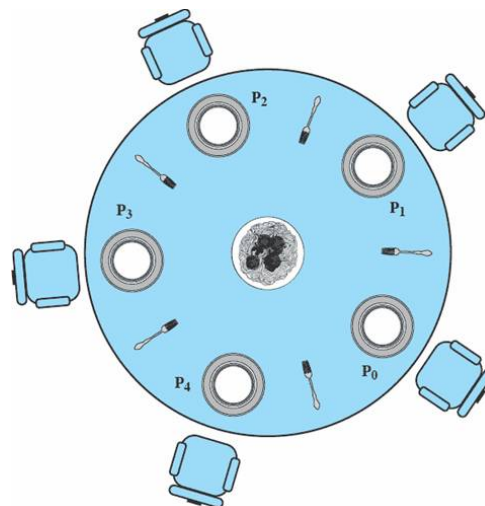


Figure 3.5 Dining arrangement for philosophers

The problem: Devise an algorithm that will allow the philosophers to eat. The algorithm must satisfy mutual exclusion (no two philosophers can use the same fork at the same time) while avoiding deadlock and starvation.

### Key points of dining philosopher's problem:

- illustrates basic problems in deadlock and starvation
- reveal many of the difficulties in concurrent programming
- deals with the coordination of shared resources, which may occur when an application includes concurrent threads of execution
- it is a standard test case for evaluating approaches to synchronization

### 3.7.1 Solution using Semaphores

The dining philosopher's problem can be solved using semaphores as shown in the code (Figure 3.6).

Each philosopher picks up the fork on the left side first and then the fork on the right. After the philosopher is finished eating, the two forks are replaced on the table. This solution leads to deadlock: If all of the philosophers are hungry at the same time, they all sit down, they all pick up the fork on their left, and they all reach out for the other fork, which is not there. In this undignified position, all philosophers starve.

To overcome the risk of deadlock, we could buy five additional forks or teach the philosophers to eat spaghetti with just one fork. As another approach, we could consider adding an attendant who only allows four philosophers at a time into the dining room. With at most four seated philosophers, at least one philosopher will have access to two forks.

```
/* program    diningphilosophers */
semaphore fork [5] = {1};
int i;
void philosopher (int i)
{
    while (true) {
        think();
        wait (fork[i]);
        wait (fork [(i+1) mod 5]);
        eat();
        signal(fork [(i+1) mod 5]);
        signal(fork[i]);
    }
}
void main()
{
    parbegin (philosopher (0), philosopher (1), philosopher
(2),
            philosopher (3), philosopher (4));
}
```

Figure 3.6 Solution to Dining Philosopher's problem

### 3.7.2 Solution using Monitors

Here, some condition variables are set to see that each philosopher must wait for the availability of the fork. A function (for entering a monitor) has to be written so that, a philosopher must seize two forks on his left and right sides. If anyone fork is not available, he must wait. Now, another philosopher can enter the monitor and try his luck. Another function has to be written for releasing the forks when a philosopher finishes eating. In this solution, deadlock will not occur.

## MEMORY MANAGEMENT

### 3.8 BACKGROUND

It has been discussed earlier that CPU scheduling is used to increase the utilization of CPU and to improve computer's response to the user. For this, several processes must be kept in the memory. Here we will discuss how memory is managed to achieve this.

Memory consists of a large array of words or bytes, each with its own address. The CPU fetches instructions from memory according to the value of the program counter. These instructions may cause additional loading/storing from/to specific memory addresses. Due to several reasons, we can ignore *how* a memory address is generated by a program. But, we have to consider only the sequence of memory addresses generated by the running program.

#### 3.8.1 Address Binding

Usually, a program resides on a disk as a binary executable file. The program must be brought into memory and placed within a process for it to be executed.

Most systems allow a user process to reside in any part of the physical memory. Thus, although the address space of the computer starts at 00000, the first address of the user process need not be 00000. This arrangement affects the addresses that the user program can use. In most cases, a user program will go through several steps-some of which may be optional-before being executed as shown in Figure 3.7. Addresses may be represented in different ways during these steps. Addresses in the source program are generally symbolic such as *count*. A compiler will typically bind these symbolic addresses to relocatable addresses (such as "14 bytes from the beginning of this module"). The linkage editor or loader will in turn bind these relocatable addresses to absolute addresses (such as 74014). Each binding is a mapping from one address space to another.

The binding of instructions and data to memory addresses can be done in following ways:

- **Compile Time:** If we know at compile time where the process will reside in memory, then **absolute code** can be generated. After sometime, if this location changes, then it is 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 can be moved during its execution from one memory segment to another, then binding must be delayed until run time. Special hardware must be available for this scheme to work.

### 3.8.2 Logical v/s Physical Address Space

An address generated by the CPU is commonly referred to as a logical address, whereas an address seen by the memory unit—that is, the one loaded into the memory-address register of the memory—is referred to as a physical address. The compile-time and load-time address-binding methods generate identical logical and physical addresses. However, the execution-time address binding scheme results in different logical and physical addresses. In this case, we usually refer to the logical address as a **virtual address**.

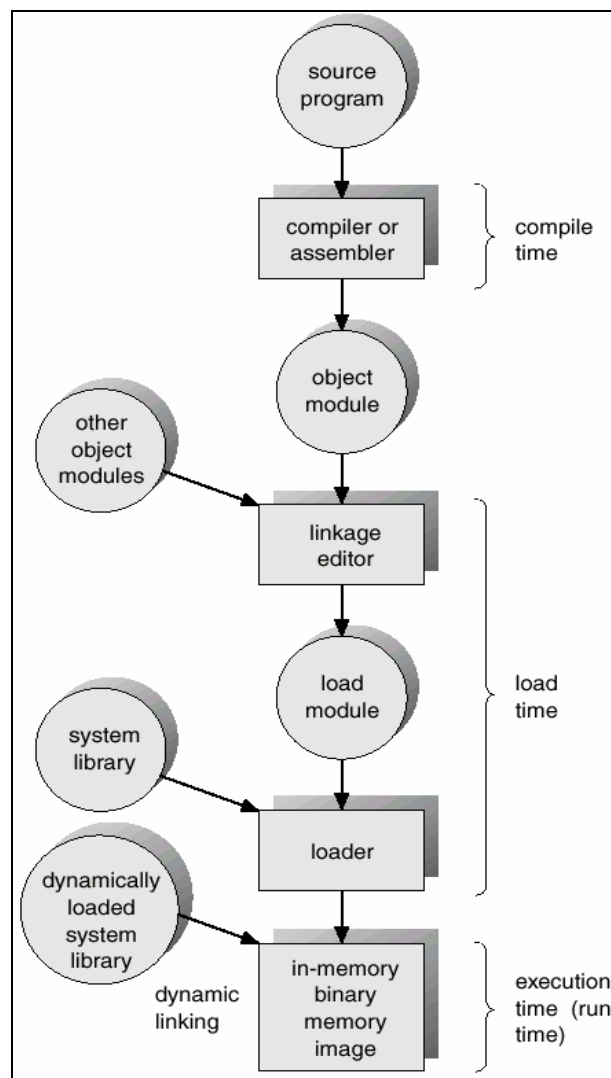


Figure 3.7 Multistep processing of a user program

The set of all logical addresses generated by a program is a logical-address space; the set of all physical addresses corresponding to these logical addresses is a physical-address



space. The run-time mapping from virtual to physical addresses is done by a hardware device called the **memory-management unit (MMU)**.

The value in the base register (also called a **relocation register**) is added to every address generated by a user process at the time it is sent to memory. 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 (Figure 3.8). The user program never sees the real physical addresses. The program can create a pointer to location 346, store it in memory, manipulate it, compare it to other addresses—all as the number 346. Only when it is used as a memory address (in an indirect load or store, perhaps) is it relocated relative to the base register. The user program deals with logical addresses. The memory-mapping hardware converts logical addresses into physical addresses. The final location of a referenced memory address is not determined until the reference is made. We now have two different types of addresses: logical addresses (in the range 0 to  $max$ ) and physical addresses (in the range  $R + 0$  to  $R + max$  for a base value  $R$ ). The user generates only logical addresses and thinks that the process runs in locations 0 to  $max$ .

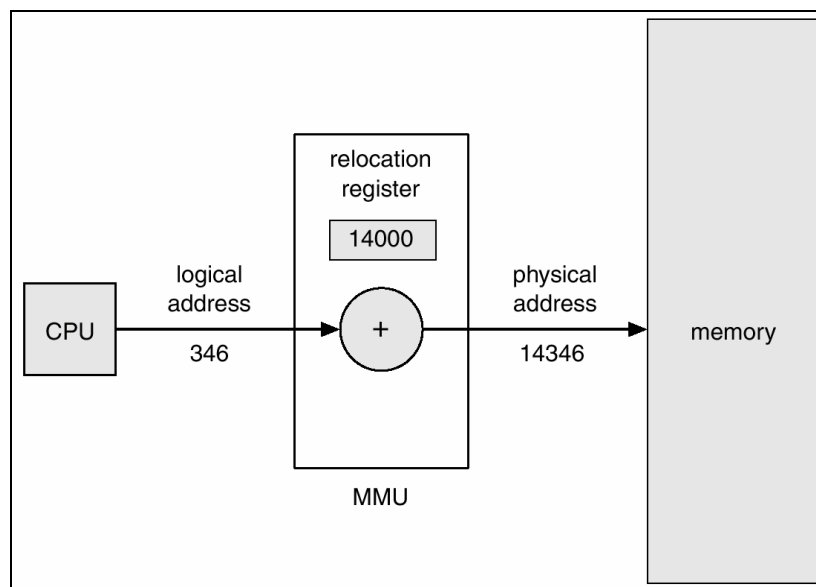


Figure 3.8 Dynamic relocation using a relocation register

### 3.8.3 Dynamic Loading

We have discussed earlier that the entire program and data of a process must be in physical memory for the process to execute. The size of a process is limited to the size of physical memory. To obtain better memory-space utilization, we can use dynamic loading. With dynamic loading, a routine is not loaded until it is called. All routines are kept on disk in a relocatable load format. The main program is loaded into memory and is executed. When a routine needs to call another routine, the calling routine first checks to see whether the other routine has been loaded. If not, the relocatable linking loader is called to load the desired routine into memory and to update the program's address tables to reflect this change. Then, control is passed to the newly loaded routine.



The advantages of dynamic loading:

- an unused routine is never loaded. Hence, though the total program size is large, the portion used will be much smaller.
- it does not require special support from the operating system

### 3.8.4 Dynamic Linking and Shared Libraries

Some OS supports static linking, in which language libraries are combined by the loader into the binary program image. The dynamic linking is similar to dynamic loading. Here, linking is postponed until execution time. Small piece of code called **stub** is used to locate the appropriate memory-resident library routine. When this stub is executed, it checks to see whether the needed routine is already in memory. If not, the program loads the routine into memory. Thus, the next time that that code segment is reached, the library routine is executed directly, incurring no cost for dynamic linking. Under this scheme, all processes that use a language library execute only one copy of the library code.

A shared library is a file that is intended to be shared by executable files and further shared object files. Modules used by a program are loaded from individual shared objects into memory at load time or run time, rather than being copied by a linker.

## 3.9 SWAPPING

A process needs to be in memory to be executed. A process, however, can be swapped temporarily out of memory to a backing store, and then brought back into memory for continued execution as shown in Figure 3.9. Such swapping may be necessary in priority based scheduling or round-robin scheduling.

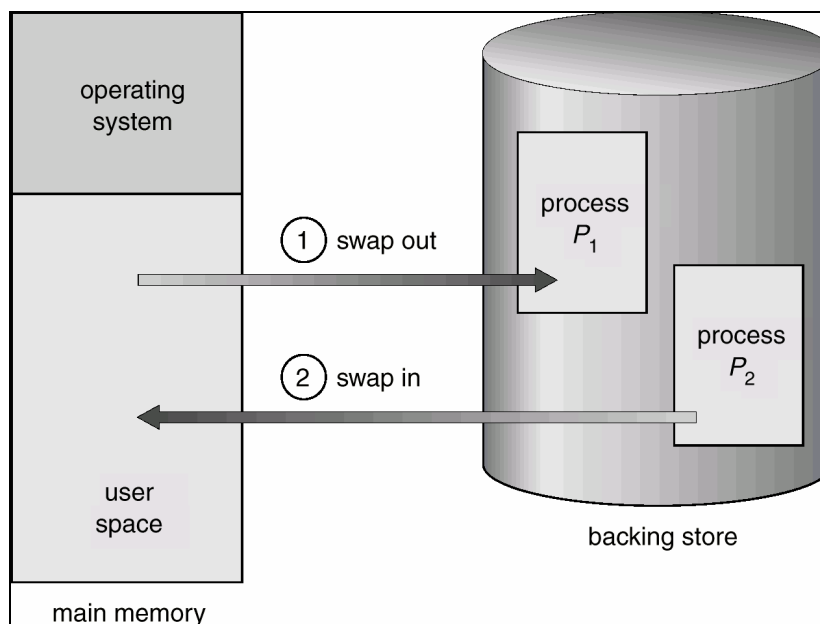


Figure 3.9 Swapping of two processes using a disk as a backing store

Normally a process that is swapped out will be swapped back into the same memory space that it occupied previously. This restriction is dictated by the method of address binding. If binding is done at assembly or load time, then the process cannot be moved to different locations. If execution-time binding is being used, then a process can be swapped into a different memory space, because the physical addresses are computed during execution time.

Swapping requires a backing store. The backing store is a fast disk. It must be large enough to accommodate copies of all memory images for all users, and it must provide direct access to these memory images. The system maintains a ready queue consisting of all processes whose memory images are on the backing store or in memory and are ready to run. Whenever the CPU scheduler decides to execute a process, it calls the dispatcher. The dispatcher checks to see whether the next process in the queue is in memory. If not, and there is no free memory region, the dispatcher swaps out a process currently in memory and swaps in the desired process. It then reloads registers as normal and transfers control to the selected process.

The context-switch time in such a swapping system is high. If we want to swap a process, it must be idle.

### 3.10 CONTIGUOUS MEMORY ALLOCATION

The main memory must accommodate both the operating system and the various user processes. We therefore need to allocate different parts of the main memory in the most efficient way possible. This is done using contiguous memory allocation. We usually want several user processes to reside in memory at the same time. Hence, we need to consider how to allocate available memory to the processes that are in the input queue waiting to be brought into memory. In this contiguous memory allocation, each process is contained in a single contiguous section of memory.

#### 3.10.1 Memory Protection

Before discussing memory allocation, we need to discuss the issue of memory protection. **Memory Protection** is to protect the OS from user processes, and to protect user processes from one another.

We can provide this protection by using following registers:

- **relocation register** : contains the value of the smallest physical address
- **limit register**: contains the range of logical addresses

For example, relocation = 100040 and limit = 74600. 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 memory as shown in Figure 3.10.

When the CPU scheduler selects a process for execution, the dispatcher loads the relocation and limit registers with the correct values as part of the context switch. Because

every address generated by the CPU is checked against these registers, we can protect OS and user programs and data from being modified by this running process.

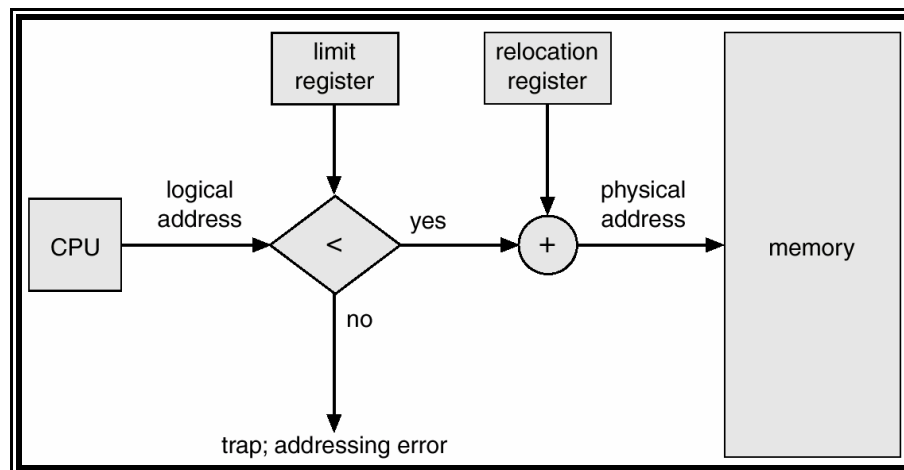


Figure 3.10 Hardware support for relocation and limit registers

### 3.10.2 Memory Allocation

One of the simplest methods for memory allocation is to divide memory into several fixed-sized **partitions**. Each partition may contain exactly one process. In this **multiple-partition method**, when a partition is free, a process is selected from the input queue and is loaded into the free partition. When the process terminates, the partition becomes available for another process. This method is no longer in use.

Another method is a generalization of the fixed-partition scheme. Here, the OS keeps a table indicating which parts of memory are available and which are occupied. Initially, all memory is available for user processes, and is considered as one large block of available memory, **a hole**. When a process arrives and needs memory, we search for a hole large enough for this process. If we find one, we allocate only as much memory as is needed, keeping the rest available to satisfy future requests. At any given time, we have a list of available block sizes and the input queue.

In general, a set of holes, of various sizes, is scattered throughout memory at any given time. When a process arrives and needs memory, the system searches this set for a hole that is large enough for this process. If the hole is too large, it is split into two: One part is allocated to the arriving process; the other is returned to the set of holes. When a process terminates, it releases its block of memory, which is then placed back in the set of holes. Two adjacent holes are merged to form one larger hole.

How to satisfy a request of size  $n$  from a list of free holes is a **dynamic storage-allocation problem**. There are following algorithms:

- **First Fit:** Allocate the first hole that is big enough.
- **Best Fit:** Allocate the smallest hole that is big enough. We must search the entire list, if the list is not ordered by size.
- **Worst Fit:** Allocate the largest hole.

Best fit and first fit are better than worst fit with respect to time and storage utilization. First fit is faster. All the three algorithms suffer from **external fragmentation**. As processes are loaded and removed from memory, there will be multiple holes. External fragmentation exists when enough total memory space is there to satisfy the request, but it is not contiguous.

### 3.10.3 Fragmentation

Sometimes, the size of the hole is much smaller than the overhead required to track it. Hence, normally, physical memory is divided into fixed-size block. Then memory is allocated in terms of units. So, sometimes, the memory allocated for a process may be slightly larger than what it requires. Such a difference is known as **internal fragmentation**. For example, if the physical memory is divided into block of 4 bytes, and the process requests for 6 bytes, then it will be allocated with 2 blocks. Hence, the total allocation is 8 bytes leading to 2 bytes of internal fragmentation.

One of the solutions to external fragmentation is **compaction**: shuffling of memory contents to place all free memory together into one large block. But, compaction is possible only when relocation is dynamic and is done during execution time. That is, if the relocation is static, we cannot apply this technique.

Another solution for external fragmentation is to permit logical-address space of a process to be non-contiguous. This allows a process to be allocated physical memory wherever it is available. To do so, two techniques are there:

- Paging
- Segmentation

These techniques are discussed in the following sections.

## 3.11 PAGING

Paging is a memory-management scheme that permits the physical-address space of a process to be noncontiguous. Paging avoids the problem of fitting the varying-sized memory chunks onto the backing store.

Paging is commonly used in most OS nowadays. It is implemented in hardware as well as in OS.

### 3.11.1 Basic Method

Physical memory is broken into fixed-sized blocks called **frames**. Logical memory is also broken into blocks of the same size called **pages**. When a process is to be executed, its pages are loaded into any available memory frames from the backing store. The backing store is divided into fixed-sized blocks that are of the same size as the memory frames. The hardware support for paging is illustrated in Figure 3.11.

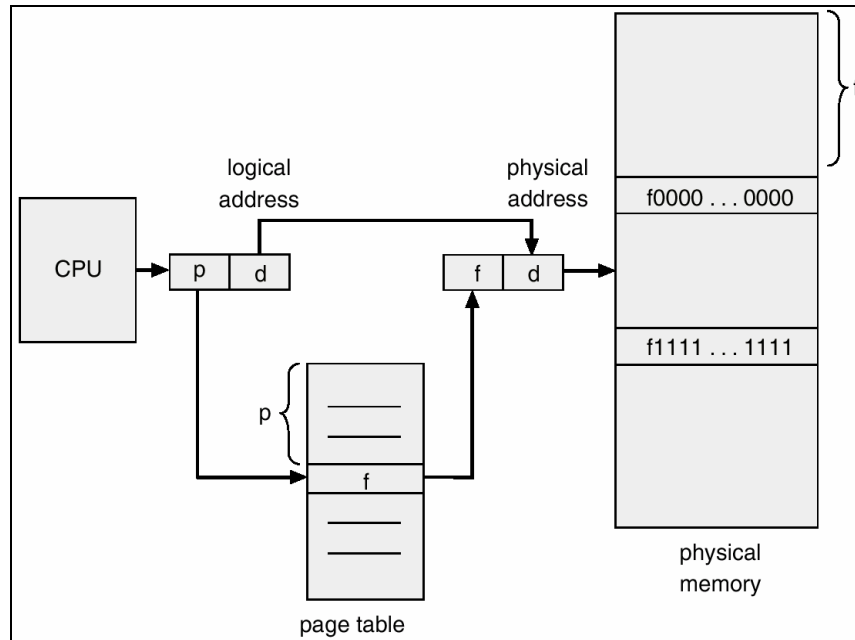


Figure 3.11 Paging hardware

Every address generated by the CPU is divided into two parts:

- **page number (p)**
- **page offset (d)**

The page number is used as an index into a page table. The page table contains the base address of each page in physical memory. This base address is combined with the page offset to define the physical memory address that is sent to the memory unit. The paging model of memory is shown in Figure 3.12.

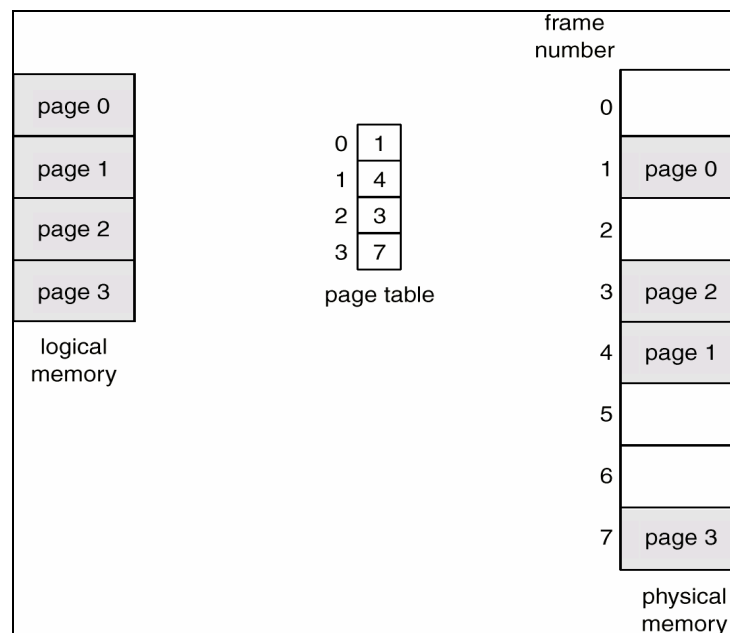


Figure 3.12 Paging model of logical and physical memory

When the CPU scheduler selects a process for execution, the dispatcher loads the relocation and limit registers with the correct values as part of the context switch. Because every address generated by the CPU is checked against these registers, we can protect OS and user programs and data from being modified by this running process.

**Example:**

Consider a page size of 4 bytes and a physical memory of 32 bytes (8 pages) as shown in Figure 3.13. We will see now, how the user's view of memory can be mapped into physical memory:

- Logical address 0 is page 0, offset 0. Indexing into the page table, we find that page 0 is in frame 5. Thus, logical address 0 maps to physical address 20 (= (5 x 4) + 0).
- Logical address 3 (page 0, offset 3) maps to physical address 23 (= (5 x 4) + 3).
- Logical address 4 (page 1, offset 0); according to the page table, page 1 is mapped to frame 6. Thus, logical address 4 maps to physical address 24 (= (6 x 4) + 0).
- Logical address 13 maps to physical address 9.

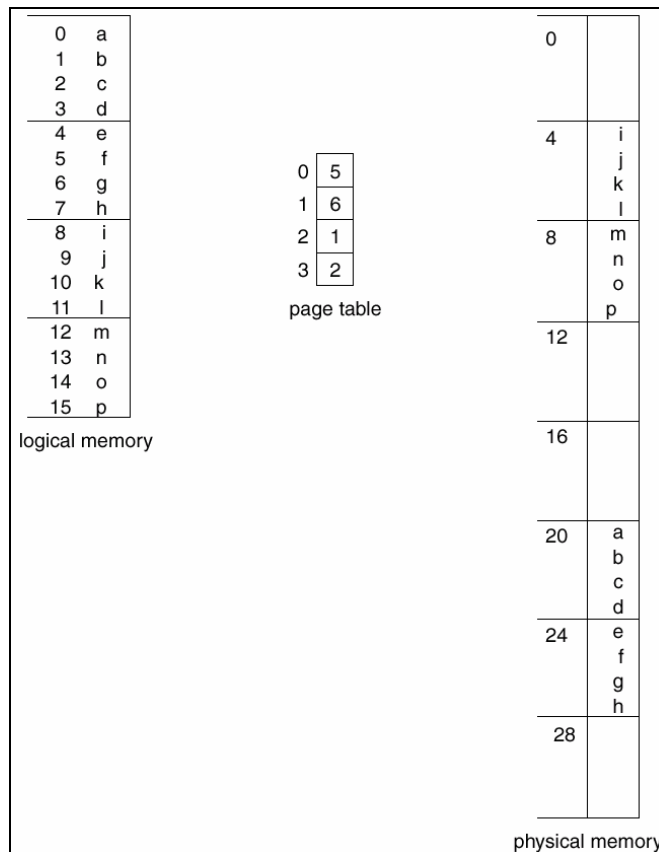


Figure 3.13 Paging Example

When paging is used, there will not be external fragmentation: Any free frame can be allocated to a process that needs it. However, some internal fragmentation will be there. To reduce it, small-sized pages are preferable. But, more number of pages will lead to many entries in page-table, and hence much overhead. So, optimum selection of size of each page is essential.

When a process needs to be executed, its size represented in terms of number of 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.

### 3.11.2 Hardware Support

Each operating system has its own methods for storing page tables. In the simplest case, the page table is implemented as a set of dedicated registers. Another method is: the page table is kept in main memory, and a **page-table base register (PTBR)** points to the page table. In this scheme every data/instruction access requires two memory accesses: One for the page table and one for the data/instruction.

This double access problem can be solved by the use of a special fast-lookup hardware cache called **associative memory** or **translation look-aside buffers (TLBs)**. The TLB is associative, high-speed memory. Each entry in the TLB consists of two parts: a key (or tag) and a value. The TLB is used with page tables in the following way:

- The TLB contains only a few of the page-table entries.
- When a logical address is generated by the CPU, its page number is presented to the TLB.
- If the page number is found, its frame number is immediately available and is used to access memory.
- If the page number is not in the TLB (known as a TLB miss), a memory reference to the page table must be made. When the frame number is obtained, we can use it to access memory.

This is shown in Figure 3.14.



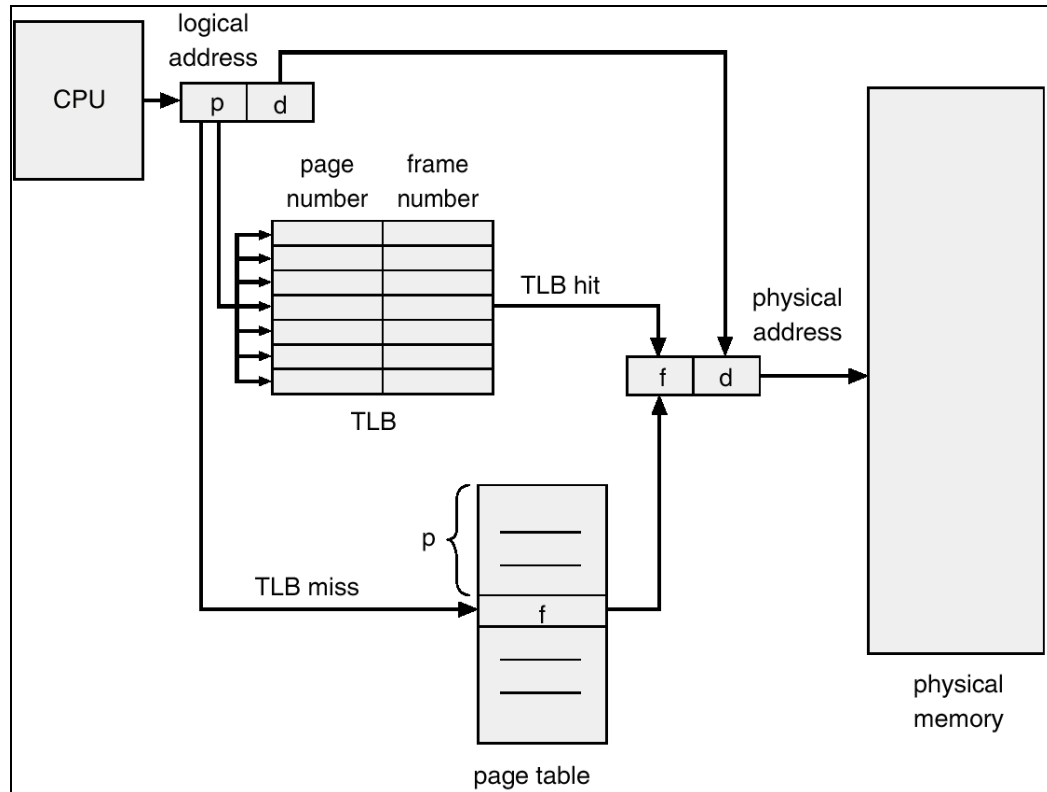


Figure 3.14 Paging hardware with TLB

### 3.11.3 Protection

In paging technique, protection bits are associated with each frame and are kept in the page table. This will help memory protection. One bit can define a page to be read-write or read-only. While referring the page table to get a frame number, the protection bits are also checked. An attempt to write a read-only page causes hardware trap to OS.

Another bit attached to page table is **valid – invalid bit**. When this bit is set to *valid*, it indicates that the respective page is in the logical address space of the process and hence it is a legal/valid page. If this bit is set to *invalid*, the page is not in logical address space of the process.

### 3.11.4 Structure of Page Table

The most common techniques for structuring the page table are:

- Hierarchical Paging
- Hashed Page Tables
- Inverted Page Tables

**Hierarchical Paging:** Most modern computer systems support a large logical-address space. In such an environment, the page table itself becomes excessively large. For example, consider a system with a 32-bit logical-address space. If the page size in such a system is 4 KB ( $2^{12}$ ), then a page table may consist of up to 1 million entries ( $2^{32}/2^{12}$ ). To avoid this problem, one of the solutions is to dividing the page table into smaller pieces. To

do so, **two-level paging algorithm** can be used. Here, the page table itself is paged as shown in Figure 3.15.

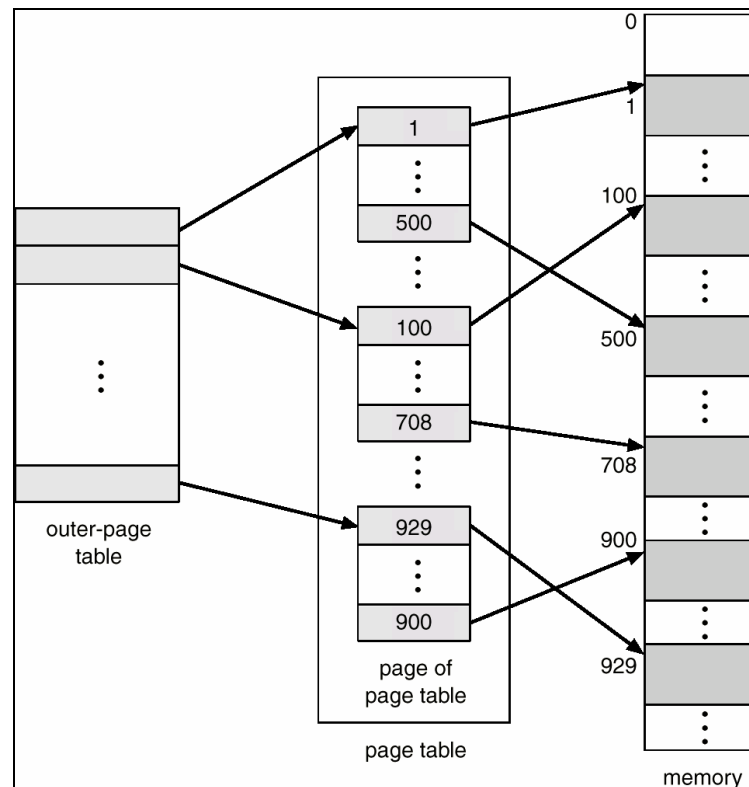


Figure 3.15 A two-level page-table scheme

**Hashed Page Tables:** A common approach for handling address spaces larger than 32 bits is to use a hashed page table, with the hash value being the virtual-page number. Each entry in the hash table contains an array of linked lists that hash to the same location (to handle collisions using open-hashing/separate chaining method). Each node in a linked list consists of three fields:

- The virtual page number
- the value of mapped page frame
- a pointer to the next element in the linked list.

The algorithm works as follows:

- i. The virtual page number in the virtual address is hashed into the hash table.
- ii. The virtual page number is compared to the first field of the first node in the linked list.
- iii. If there is a match, the corresponding page frame (second field of a node) is used to form the desired physical address.
- iv. If there is no match, subsequent nodes in the linked list are searched for a matching virtual page number.

This scheme is shown in Figure 3.16.

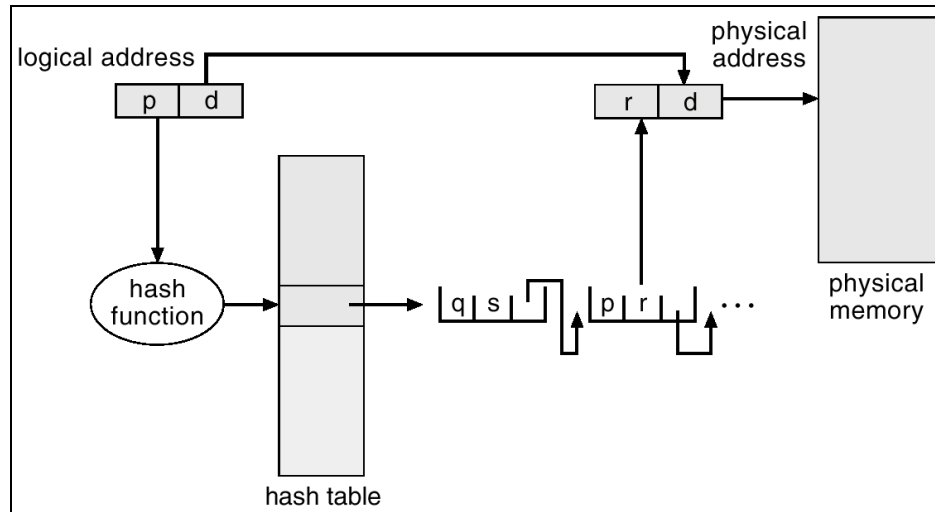


Figure 3.16 Hashed page table

**Inverted Page Table:** It is obvious that each process has its own page table and every page table has one entry for each page. We know that a page represents logical/virtual addresses. It is the OS which converts the logical addresses into physical addresses later. As, a page table contains millions of entries, it is a burden on OS to convert everything into physical address.

To avoid this problem, inverted page table is used. This table contains one entry for each real physical frame of memory rather than for logical memory. Hence, only one page table is sufficient for the system. As each process has its own address space, another entity specifying **process identifier (pid)** in its address – space is used as shown in Figure 3.17. The working is as explained below:

- The page table contains pid and page number.
- The CPU generates logical address containing pid, page number and offset.
- The pair pid and page number is searched for, in the page table.
- If the match is found at  $i^{\text{th}}$  entry, then the physical address will be  $\langle i, \text{offset} \rangle$
- If the match is not found, illegal address access error occurs.

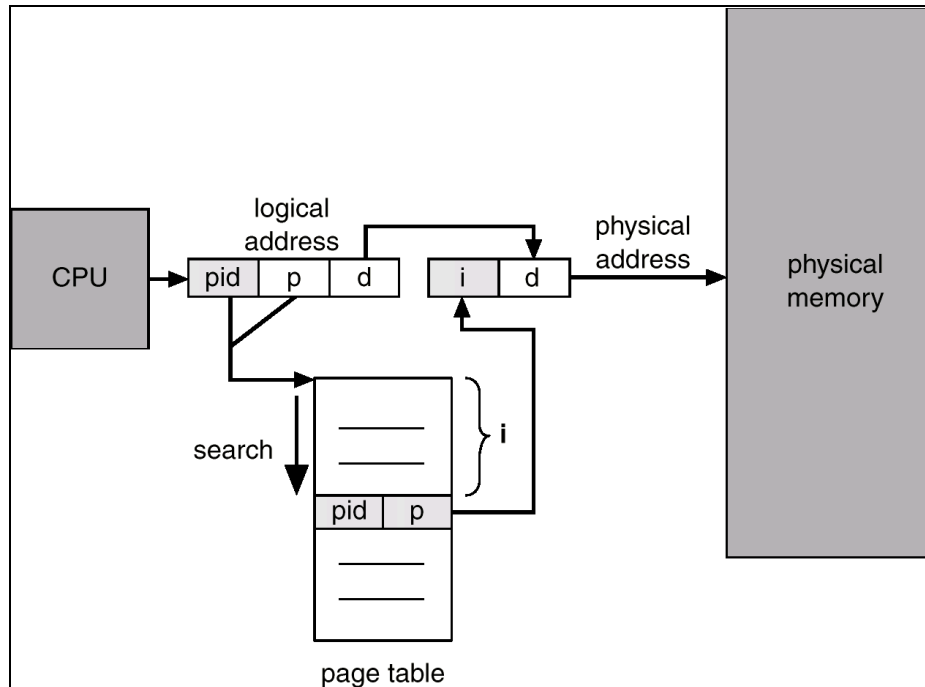


Figure 3.17 Inverted Page Table

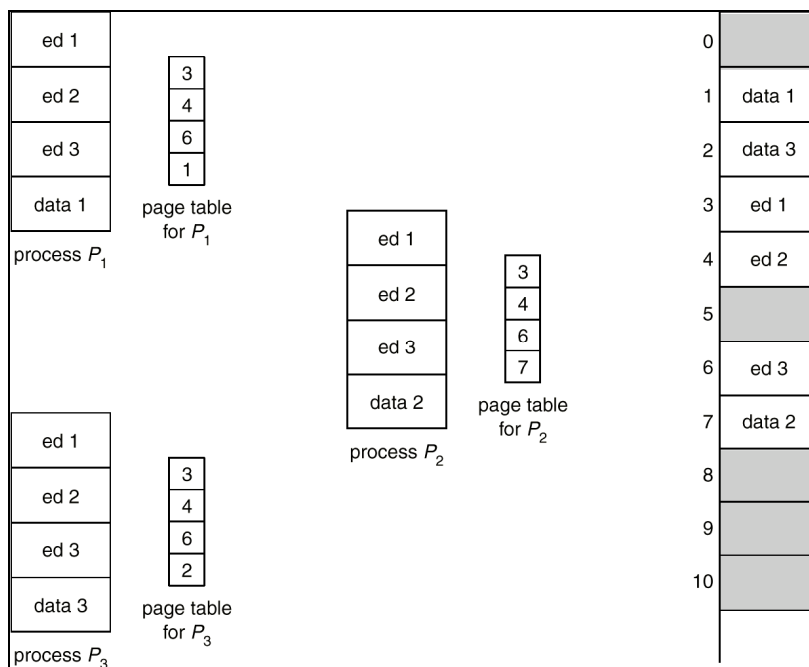


Figure 3.18 Sharing of code in Paging environment

### 3.11.5 Shared Pages

Paging technique is useful in time-sharing environment. It allows sharing of common code, if possible. Some of the processes have **re-entrant code**. That is, once they start executing, code will not change. In such situations, pages can be shared. Consider an

example: Assume that 40 users are using a text editor at a time. Let the text editor consists of 150KB of code and 50KB of data space. We need 8000KB for 40 users. As text-editor is reentrant, only one copy of the editor needs to be kept in the physical memory and it can be shared among all the users. But, the data space should be different for every user. This is depicted in Figure 3.18. In this example, size of each page (and also the frame) is taken as 50KB.

Such sharing technique can be applied for compilers, window systems, run-time libraries, database systems etc.

### 3.12 SEGMENTATION

In paging, we have observed that user view of the memory and the actual physical memory are different. User view of memory is mapped into physical memory.

#### 3.12.1 Basic Method

Users do not think memory as a linear array of bytes. Instead, users view memory as a collection of variable-sized segments, with no necessary ordering among segments as shown in Figure 3.19. For example, we think of a program as a main program with a set of subroutines, procedures, functions, or modules. There may also be various data structures: tables, arrays, stacks, variables, and so on. Each of these modules or data elements is referred to by name. We do not bother which module has is stored before/after other module.

**Segmentation** is a memory-management scheme that supports user view of memory. A logical-address space is a collection of segments. Each address is specified in terms of the segment number and the offset within the segment.

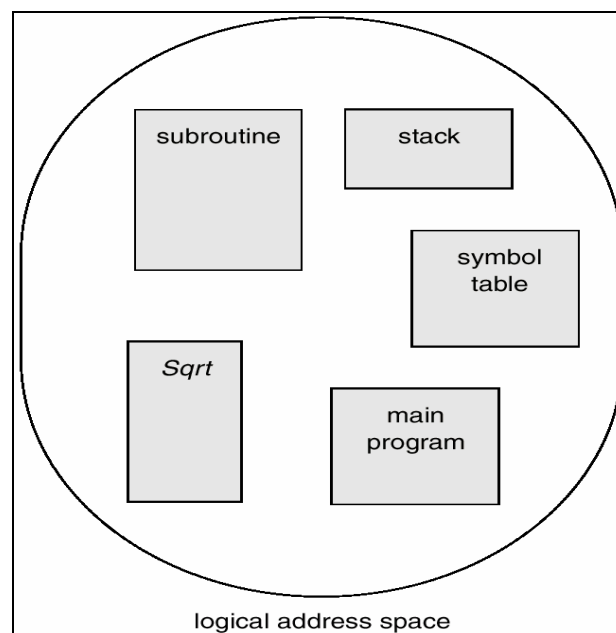


Figure 3.19 User view of a program

### 3.12.2 Hardware

Here, a two dimensional user-defined addresses are mapped into one-dimensional physical addresses. This mapping is affected by a segment table. Each entry of the segment table has

- *segment base* – contains the starting physical address where the segment resides in memory
- *segment limit* – specifies the length of the segment

A logical address consists of two parts: a segment number,  $s$  and an offset  $d$ . The segment number is used as an index into the segment table. The offset  $d$  of the logical address must be between 0 and the segment limit. If it is not, we trap to the OS (logical addressing attempt beyond end of segment). If this offset is legal, it is added to the segment base to produce the address in physical memory of the desired byte. The segment table is thus essentially an array of base-limit register pairs. The structure is as given in Figure 3.20.

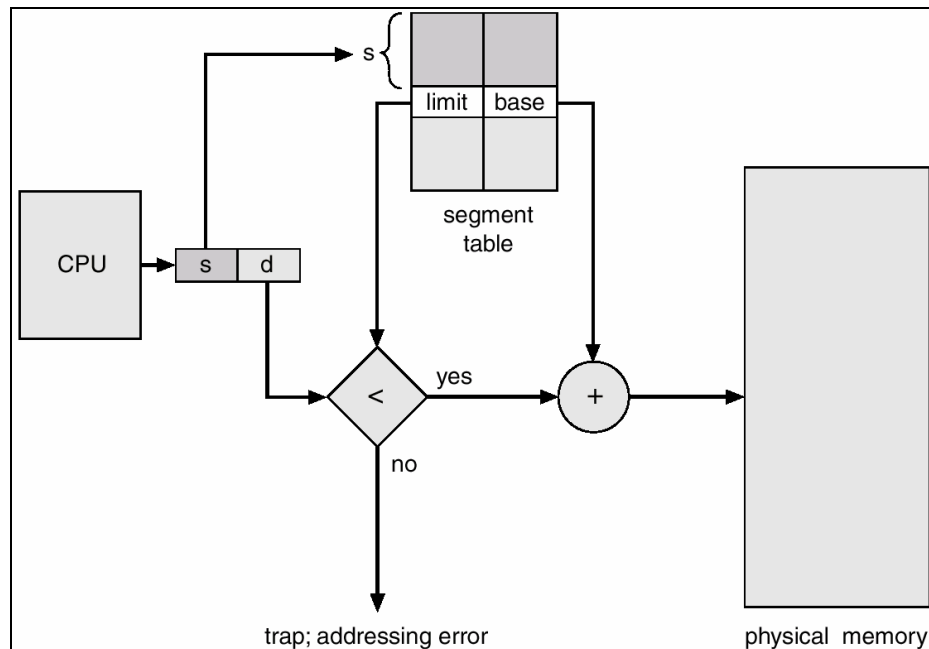


Figure 3.20 Segmentation Hardware

**Example:** Consider the situation shown in Figure 3.21. We have five segments numbered from 0 through 4. The segment table has a separate entry for each segment. It is giving the beginning address of the segment in physical memory (or base) and the length of that segment (or limit). For example, segment 2 is 400 bytes long and begins at the location 4300.

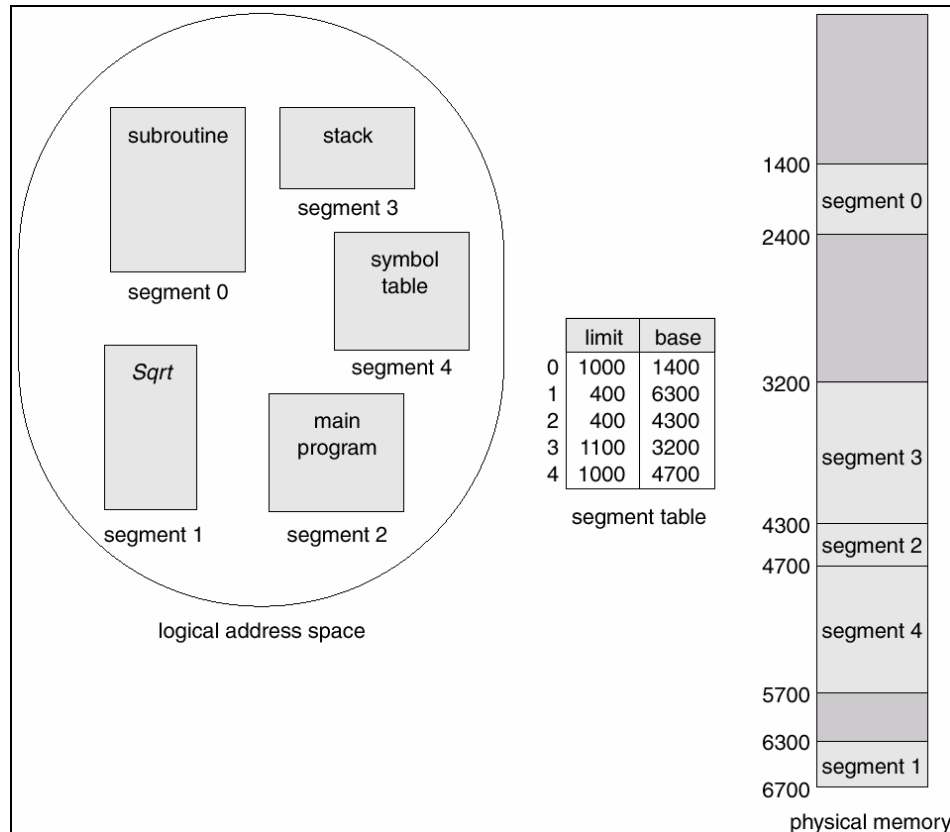


Figure 3.21 Example of Segmentation

### 3.12.3 Fragmentation

The long-term scheduler must find and allocate memory for all the segments of a user program. This situation is similar to paging *except* that the segments are of variable length; pages are all the same size. Thus, as with the variable-sized partition scheme, memory allocation is a dynamic storage-allocation problem, usually solved with a best-fit or first-fit algorithm. Segmentation may then cause external fragmentation, when all blocks of free memory are too small to accommodate a segment. In this case, the process may simply have to wait until more memory (or at least a larger hole) becomes available, or until compaction creates a larger hole. Because segmentation is by its nature a dynamic relocation algorithm, we can compact memory whenever we want. If the CPU scheduler must wait for one process, because of a memory allocation problem, it may (or may not) skip through the CPU queue looking for a smaller, lower-priority process to run.

### 3.13 SEGMENTATION WITH PAGING

Both paging and segmentation have advantages and disadvantages. Paging has advantages like –

- No external fragmentation
- Faster allocation

Segmentation has advantages like –

- Sharing
- Protection



These advantages are combined to get better results.

Segmentation with paging is explained as below: Instead of an actual memory location, the segment information includes the address of a page table for the segment. When a program references a memory location the offset is translated to a memory address using the page table. A segment can be extended simply by allocating another memory page and adding it to the segment's page table. Also, the logical address space of a process is divided into two partitions: one consisting of segments that are private to that process, and the other partition consists of segments that can be shared among all the processes.

### 3.14 DEMAND PAGING

Demand paging is one of the techniques of implementing virtual memory.

**(NOTE: Virtual Memory** is the separation of user logical memory from physical memory. It allows the execution of processes that may not be completely in memory. Hence, logical address space can therefore be much larger than physical address space. Moreover, it allows address spaces to be shared by several processes and also, allows for more efficient process creation.)

Demand paging is similar to paging system with swapping. Whenever process needs to be executed, only the required pages are swapped into memory. This is called as **lazy swapping**. As, the term *swapper* has a different meaning of 'swapping entire process into memory', another term **pager** is used in the discussion of demand paging.

When a process is to be swapped in, the pager guesses which pages will be used before the process is swapped out again. The pager brings only those necessary pages into memory. Hence, it decreases the swap time and the amount of physical memory needed.

In this scheme, we need to distinguish the pages which are in the memory and pages which are there on the disk. For this purpose, the valid – invalid bit is used. When this bit is set to *valid*, it indicates the page is in the memory. Whereas, the value of bit as *invalid* indicates page is on the disk.

If the process tries to access a page which is not in the memory (means, it is on the disk), **page fault** occurs. The paging hardware notices the *invalid* bit in the page table and cause a trap to the OS. This trap is the result of the failure of OS to bring the desired page into memory. This error has to be corrected. The procedure for handling this page fault is as shown in Figure 3.22. The steps are explained below:

1. We check an internal table (usually kept with the process control block) for this process, to determine whether the reference was a valid or invalid memory access.
2. If the reference was invalid, we terminate the process. If it was valid, but we have not yet brought in that page into memory, it is brought now.
3. We find a free frame.
4. We schedule a disk operation to read the desired page into the newly allocated frame.

5. When the disk read is complete, we modify the internal table kept with the process and the page table to indicate that the page is now in memory.
6. We restart the instruction that was interrupted by the illegal address trap. The process can now access the page as though it had always been in memory.

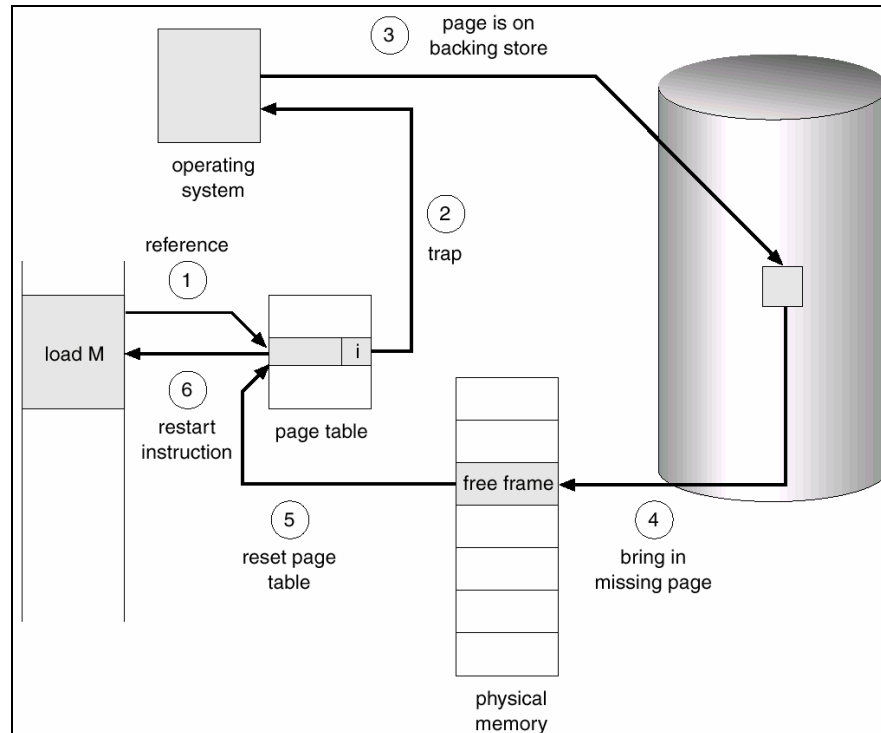


Figure 3.22 Steps in handling page fault

It is important to realize that, because we save the state (registers, condition code, instruction counter etc.) of the interrupted process when the page fault occurs, we can restart the process in exactly the same place and state

In an extreme situation, a process may start executing with no page in the memory. So, each time an instruction has to be executed, page fault occurs and the required page needs to be brought into the memory. This situation is called as **pure demand paging**. That is, no page is brought into the memory until it is required.

### 3.15 PROCESS CREATION

Virtual memory and paging helps in process creation. Two such techniques are discussed here.

#### 3.15.1 Copy-on-Write

Usually, a system call *fork()* can be used to create a child process as a duplicate of its parent. Alternatively, **copy-on-write** technique can be used. In this technique, initially, parent and child processes share the same pages. These shared pages are marked as copy-on-write pages. That is, if either process writes to a shared page, a copy of the shared page is created. For example, assume the child process attempts to modify a page

containing portions of the stack; the OS recognizes this as a copy-on-write page. The operating system will then create a copy of this page, mapping it to the address space of the child process. Therefore, the child process will modify its copied page and not the page belonging to the parent process. And, all non-modified pages may be shared by the parent and child processes. Pages that cannot be modified (i.e., pages containing executable code) may be shared by the parent and child. Copy-on-write is a common technique used by several OS when duplicating processes.

### 3.15.2 Memory-Mapped Files

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

## 3.16 PAGE REPLACEMENT

In a multiprogramming system, there are several processes running. Each of these processes might have been divided into several pages. Physical memory is divided into finite set of frames. Demand paging is applied for bringing the necessary pages from disk to memory. In such a set up, a situation may arise like this:

A process requests for a page which is not there in the memory. Hence, the page fault occurs. Now, this page has to be swapped into memory from the back store. But, at this moment, no free frames are available.

Now, OS has to make the frames free by swapping out some of the frames to disk and bring newly requested pages into memory. This procedure is called as **page replacement**.

### 3.16.1 Basic Scheme

Page replacement takes the following approach. If no frame is free, we find one page that is not currently being used and free it. Necessary changes have to be made in page table to indicate that this page is no longer in memory. We can now use the freed frame to hold the page for which the process faulted. The logic is shown in Figure 3.23. Modified page-fault service routine to include page replacement is as below:

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

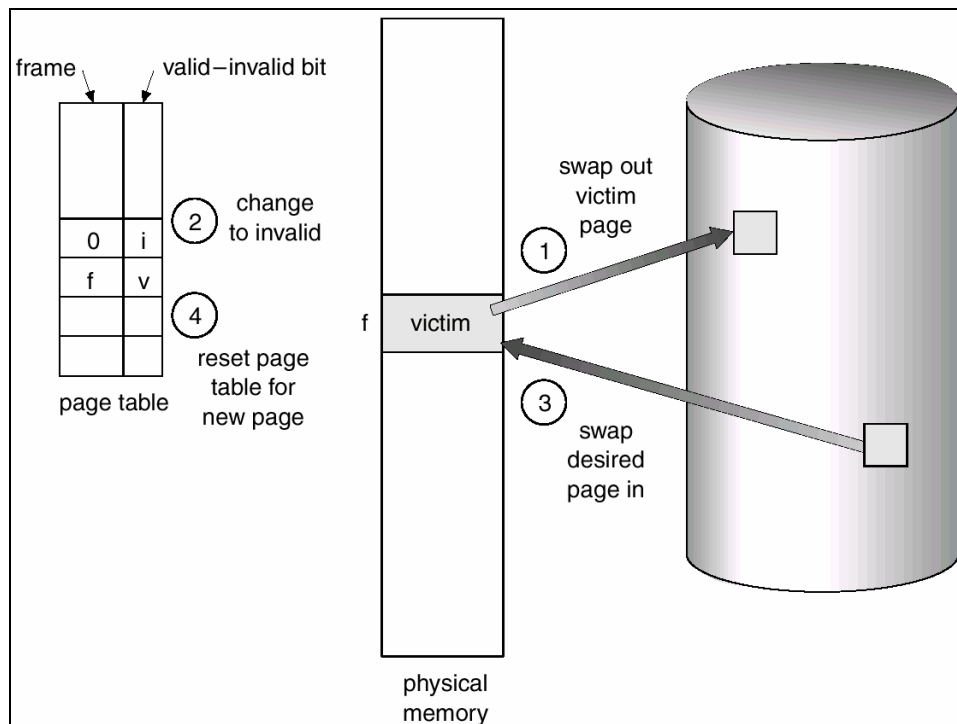


Figure 3.23 Page Replacement

**(NOTE:** Whenever a page has to be swapped out, the question arises: whether this page has to be written back to disk? Writing is necessary if the page has been modified during the process execution. Instead of writing every time a page is being swapped out, a bit called as **modify bit** (or **dirty bit**) is introduced in the hardware. If any modification is done, the bit is set, otherwise not. So, while swapping out a page, the value of this bit is checked and the page is written to disk only if the bit is set. Otherwise, as the copy of the page will always reside in the disk, it need not be written again.)

There are many page replacement algorithms. Every OS adopts its own page-replacement scheme. The algorithm with the lowest page-fault rate must be chosen. Page-replacement algorithm is evaluated by running it on a string of memory references (called as a **reference string**). Normally, a reference string consists of the numbers indicating the sequence of pages required by the process for execution. For applying any algorithm for page – replacement, we must know the number of frames available.

### 3.16.2 FIFO Page Replacement

It is the simplest page – replacement algorithm. As the name suggests, the first page which has been brought into memory will be replaced first when there no space for new page to arrive. Initially, we assume that no page is brought into memory. Hence, there will be few (that is equal to number of frames) page faults, initially. Then, whenever there is a request

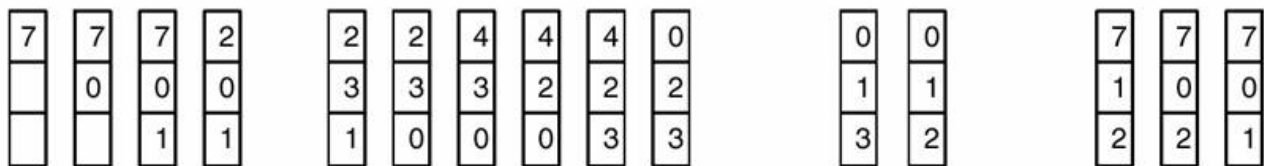
for a page, it is checked inside the frames. If that page is not available, page – replacement should take place.

**Example: Consider a reference string:** 7, 0, 1, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2, 1, 2, 0, 1, 7, 0, 1  
 Let the number of frames be 3.

Now, page – replacement according to FIFO algorithm is done as below –

reference string

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



page frames

In the above example, there are 15 page faults.

**NOTE** that, the performance of FIFO page replacement algorithm is not good. Because, the page that has been just swapped out may be needed immediately in the next step. This will cause more page faults.

Sometimes, one may feel that increase in number of frames may decrease the page faults. But this is not true always. For example, consider the reference string: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5. Now, the number of page faults for three frames will be 9 and number of page faults for four frames will be 10. Such an unexpected result is known as **Belady's anomaly**. It can be stated as: *the page fault rate may increase as the number of allocated frames increases.*

### 3.16.3 Optimal Page Replacement

An Optimal Page Replacement algorithm (also known as **OPT** or **MIN** algorithm) do not suffer from Belady's anomaly. It is stated as:

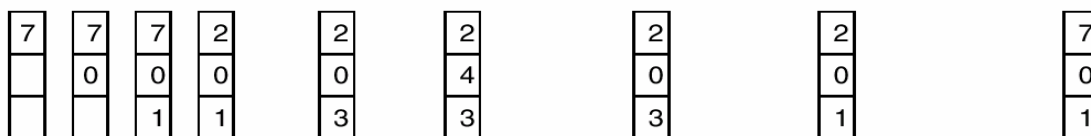
**Replace the page that will not be used for the longest period of time.**

This algorithm results in lowest page – faults.

**Example:**

reference string

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



page frames

Here, number of page faults = 9

The optimal page-replacement algorithm is difficult to implement, because it requires future knowledge of the reference string. So, the optimal algorithm is used mainly for comparison studies.

### 3.16.4 LRU Page Replacement

Least Recently Used page replacement algorithm states that:

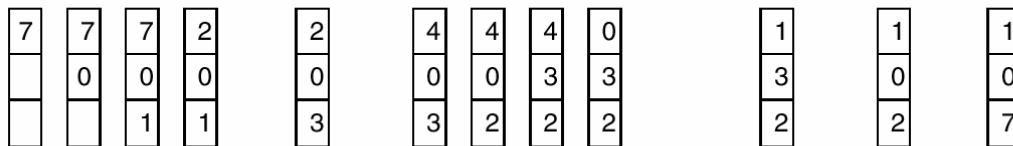
**Replace the page that has not been used for the longest period of time.**

This algorithm is better than FIFO.

#### Example:

reference string

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



page frames

Here, number of page faults = 12

**NOTE:** The major problem in LRU algorithm is hardware implementation. There are two approaches for hardware implementation:

- **Counters:** Each page table is associated with a *time-of-use* field and CPU is added with a logical clock or counter. When a page needs to be changed, look at the counters to determine which page has to be changed.
- **Stack:** Keep a stack of page numbers. Whenever a page is referenced, it is removed from the stack and put on the top. In this way, the top of the stack is always the most recently used page and the bottom is the least recently used page. To maintain this stack, a doubly linked list can be used.

## 3.17 ALLOCATION OF FRAMES

Allocation of limited set of free frames to all the processes is challenge. How many frames must be allocated to each process – is a question of interest.

In a single processor system, if only one process has to be executed, all available frames can be allocated to that process. If there are more pages than the frames, and the request is made for a new page, then page – replacement has to be done. But, in case of multiprogramming environment, there will be more than one process in memory at time. Now, the logic used in single-processor system doesn't work.

### 3.17.1 Minimum Number of Frames

Allocation of frames has various constraints:

- We cannot allocate more than the total number of available frames
- We must allocate at least a minimum number of frames to each process

In most of the cases, as the number of frames allocated to each process decreases, the page fault increases. This will slow down the process execution.

In general, each process should get a minimum number of frames. And this minimum number is defined by the instruction `-set architecture`. Remember that, when a page fault occurs before an executing instruction is complete, the instruction must be restarted. Consequently, we must have enough frames to hold all the different pages that any single instruction can reference. The maximum number of frames that can be allocated to a process depends on available physical memory.

### 3.17.2 Global v/s Local Allocation

As multiple processes competing for frames, the page – replacement algorithms can be classified into two categories as explained below:

- **Global Replacement:**
  - It allows a process to select a replacement frame from the set of all frames, including the frames that are being used by other processes.
  - High – priority processes can select a replacement from the lower priority processes.
  - Number of frames allocated for a process may increase/decrease in the course of execution.
  - A process cannot control its own page-fault rate.
  - As a process depends on paging behavior of other processes also, its performance may vary from one execution to other execution.
  - This method results in greater throughput.
- **Local Replacement:**
  - Here, a process can chose a replacement frame from only its own set of allocated frames.
  - Number of frames allocated for a process remains constant in the course of execution.
  - A process depends only on its paging behavior. Hence, the performance is constant.

## 3.18 THRASHING

If the number of frames allocated to a low-priority process falls below the minimum number required by the computer architecture, we must suspend the execution of that process. We should then page out its remaining pages, freeing all its allocated frames. This provision introduces a swap-in, swap-out level of intermediate CPU scheduling.

Whenever any process does not have enough frames, it will page-fault. At this point, it must replace some page. However, since all its pages are in active use, it must replace a page that will be needed again right away. Consequently, it quickly faults again, and again, and again. The process continues to fault, replacing pages for which it then faults and brings back in right away. This high paging activity is called **thrashing**. A process is thrashing if it is spending more time paging than executing.

Thrashing affects the performance of CPU as explained below:



If the CPU utilization is low, we normally increase the degree of multiprogramming by adding a new process to the system. A global page-replacement algorithm is used, and hence, the new process replaces the frames belonging to other processes as well. As the degree of multiprogramming increases, obviously there will be more page faults leading to thrashing. When every process starts waiting for paging rather than executing, the CPU utilization decreases. This problem is shown in Figure 3.24. The effects of thrashing can be limited by using local replacement algorithm.

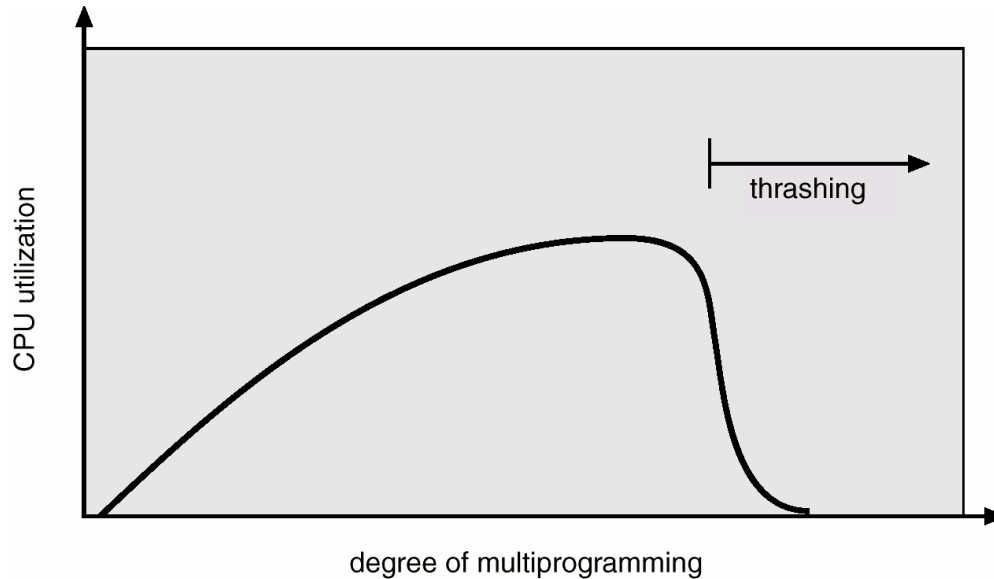


Figure 3.24 Thrashing