



[Scan or click here for more resources](#)

## UNIT-2

# Introduction of Process Management

### Program vs Process

A process is a program in execution. For example, when we write a program in C or C++ and compile it, the compiler creates binary code. The original code and binary code are both programs. When we actually run the binary code, it becomes a process.

A process is an 'active' entity, instead of a program, which is considered a 'passive' entity. A single program can create many processes when run multiple times; for example, when we open a .exe or binary file multiple times, multiple instances begin (multiple processes are created).

### Attributes or Characteristics of a Process

A process has the following attributes.

*div block*

1. Process Id: A unique identifier assigned by the operating system  
2. Process State: Can be ready, running, etc.  
3. CPU registers: Like the Program Counter (CPU registers must be saved and

restored when a process is swapped in and out of CPU)  
5.

Accounts information:  
6. I/O status information: For example, devices allocated to the process,

open files, etc.  
8. CPU scheduling information: For example, Priority (Different processes

may have different priorities, for example  
a shorter process assigned high priority  
in the shortest job first scheduling)

### States of Process:

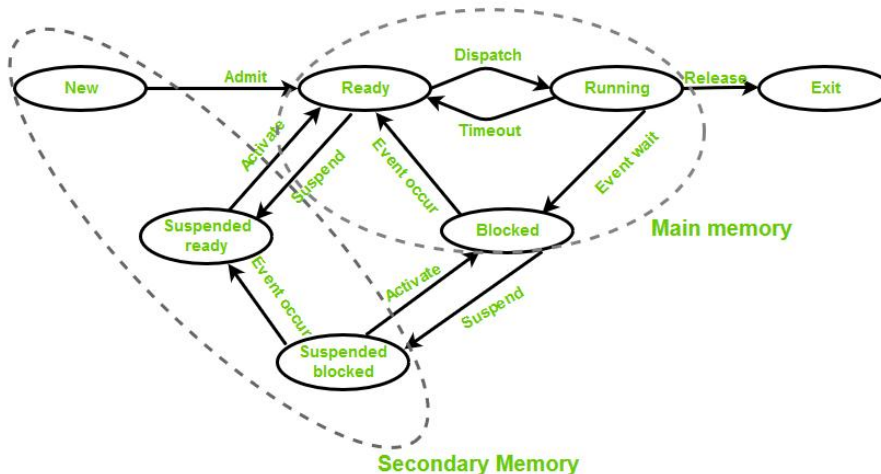
A process is in one of the following states:

1. **New:** Newly Created Process (or) being-created process.
2. **Ready:** After creation process moves to Ready state, i.e. the process is ready for execution.
3. **Run:** Currently running process in CPU (only one process at a time can be under execution in a single processor).
4. **Wait (or Block):** When a process requests I/O access.
5. **Complete (or Terminated):** The process completed its execution.
6. **Suspended Ready:** When the ready queue becomes full, some processes are moved to suspended ready state
7. **Suspended Block:** When waiting queue becomes full.

### Context Switching

The process of saving the context of one process and loading the context of another process is known as Context Switching. In simple terms, it is like loading and unloading the process from the running state to the ready state.

## States of a Process in Operating Systems



- **New (Create)** – In this step, the process is about to be created but not yet created, it is the program which is present in secondary memory that will be picked up by OS to create the process.
- **Ready** – New → Ready to run. After the creation of a process, the process enters the ready state i.e. the process is loaded into the main memory. The process here is ready to run and is waiting to get the CPU time for its execution. Processes that are ready for execution by the CPU are maintained in a queue for ready processes.
- **Run** – The process is chosen by CPU for execution and the instructions within the process are executed by any one of the available CPU cores.

- **Blocked or wait** – Whenever the process requests access to I/O or needs input from the user or needs access to a critical region (the lock for which is already acquired) it enters the blocked or wait state. The process continues to wait in the main memory and does not require CPU. Once the I/O operation is completed the process goes to the ready state.
- **Terminated or completed** – Process is killed as well as PCB is deleted.
- **Suspend ready** – Process that was initially in the ready state but were swapped out of main memory (refer Virtual Memory topic) and placed onto external storage by scheduler are said to be in suspend ready state. The process will transition back to ready state whenever the process is again brought onto the main memory.
- **Suspend wait or suspend blocked** – Similar to suspend ready but uses the process which was performing I/O operation and lack of main memory caused them to move to secondary memory. When work is finished it may go to suspend ready.

#### CPU and IO Bound Processes:

If the process is intensive in terms of CPU operations then it is called CPU bound process. Similarly, If the process is intensive in terms of I/O operations then it is called IO bound process.

#### Types of schedulers:

1. **Long term – performance** – Makes a decision about how many processes should be made to stay in the ready state, this decides the degree of multiprogramming. Once a decision is taken it lasts for a long time hence called long term scheduler.
2. **Short term – Context switching time** – Short term scheduler will decide which process to be executed next and then it will call dispatcher. A dispatcher is a software that moves process from ready to run and vice versa. In other words, it is context switching.
3. **Medium term – Swapping time** – Suspension decision is taken by medium term scheduler. Medium term scheduler is used for swapping that is moving the process from main memory to secondary and vice versa.

## Process Schedulers in Operating System

The process scheduling is the activity of the process manager that handles the removal of the running process from the CPU and the selection of another process on the basis of a particular strategy.

There are three types of process scheduler.

1. **Long Term or job scheduler :**  
It brings the new process to the 'Ready State'. It controls *Degree of Multi-programming*, i.e., number of process present in ready state at any point of time. It is important that the long-term scheduler make a careful selection of both IO and CPU bound process. IO bound tasks are which use much of their time in input and output operations while CPU bound processes are which spend their time on CPU. The job scheduler increases efficiency by maintaining a balance between the two.
2. **Short term or CPU scheduler :**  
It is responsible for selecting one process from ready state for scheduling it on the running state.

Note: Short-term scheduler only selects the process to schedule it doesn't load the process on running. Here is when all the scheduling algorithms are used. The CPU scheduler is responsible for ensuring there is no starvation owing to high burst time processes.

**Dispatcher** is responsible for loading the process selected by Short-term scheduler on the CPU (Ready to Running State) Context switching is done by dispatcher only. A dispatcher does the following:

1. Switching context.
2. Switching to user mode.
3. Jumping to the proper location in the newly loaded program.

3. **Medium-term scheduler :**

It is responsible for suspending and resuming the process. It mainly does swapping (moving processes from main memory to disk and vice versa). Swapping may be necessary to improve the process mix or because a change in memory requirements has overcommitted available memory, requiring memory to be freed up. It is helpful in maintaining a perfect balance between the I/O bound and the CPU bound. It reduces the degree of multiprogramming.

## Process Table and Process Control Block (PCB)

While creating a process the operating system performs several operations. To identify the processes, it assigns a process identification number (PID) to each process. As the operating system supports multi-programming, it needs to keep track of all the processes. For this task, the process control block (PCB) is used to track the process's execution status. Each block of memory contains information about the process state, program counter, stack pointer, status of opened files, scheduling algorithms, etc. All these information is required and must be saved when the process is switched from one state to another. When the process makes a transition from one state to another, the operating system must update information in the process's PCB.

A process control block (PCB) contains information about the process, i.e. registers, quantum, priority, etc. The process table is an array of PCB's, that means logically contains a PCB for all of the current processes in the system.

|                                     |
|-------------------------------------|
| Pointer                             |
| Process State                       |
| Process Number                      |
| Program Counter                     |
| Registers                           |
| Memory Limits                       |
| Open File Lists                     |
| Misc. Accounting<br>and Status Data |

Process Control Block

- **Pointer** – It is a stack pointer which is required to be saved when the process is switched from one state to another to retain the current position of the process.
- **Process state** – It stores the respective state of the process.

- **Process number** – Every process is assigned with a unique id known as process ID or PID which stores the process identifier.
- **Program counter** – It stores the counter which contains the address of the next instruction that is to be executed for the process.
- **Register** – These are the CPU registers which includes: accumulator, base, registers and general purpose registers.
- **Memory limits** – This field contains the information about memory management system used by operating system. This may include the page tables, segment tables etc.
- **Open files list** – This information includes the list of files opened for a process.

### **What is a Thread?**

A thread is a path of execution within a process. A process can contain multiple threads.

### **Why Multithreading?**

A thread is also known as lightweight process. The idea is to achieve parallelism by dividing a process into multiple threads. For example, in a browser, multiple tabs can be different threads. MS Word uses multiple threads: one thread to format the text, another thread to process inputs, etc. More advantages of multithreading are discussed below

### **Process vs Thread?**

The primary difference is that threads within the same process run in a shared memory space, while processes run in separate memory spaces.

Threads are not independent of one another like processes are, and as a result threads share with other threads their code section, data section, and OS resources (like open files and signals). But, like process, a thread has its own program counter (PC), register set, and stack space.

### **Difference between Process and Thread:**

| S.NO | Process   | Thread   |
|------|---|--|
| 1.   | Process means any program is in execution.          | Thread means segment of a process.                 |
| 2.   | Process takes more time to terminate.               | Thread takes less time to terminate.               |
| 3.   | It takes more time for creation.                    | It takes less time for creation.                   |
| 4.   | It also takes more time for context switching.      | It takes less time for context switching.          |
| 5.   | Process is less efficient in term of communication. | Thread is more efficient in term of communication. |

| S.NO | Process  | Thread   |
|------|--|--|
| 6.   | Multi programming holds the concepts of multi process.   | We don' t need multi programs in action for multiple threads because a single process consists of multiple threads.  |
| 7.   | Process is isolated.   | Threads share memory.  |
| 8.   | Process is called heavy weight process.  | A Thread is lightweight as each thread in a process shares code, data and resources.   |
| 9.   | Process switching uses interface in operating system.  | Thread switching does not require to call a operating system and cause an interrupt to the kernel.   |
| 10.  | If one process is blocked then it will not effect the execution of other process                   | Second thread in the same task could not run, while one server thread is blocked.  |
| 11.  | Process has its own Process Control Block, Stack and Address Space.                                | Thread has Parents' PCB, its own Thread Control Block and Stack and common Address space.  |
| 12.  | If one process is blocked, then no other process can execute until the first process is unblocked. | While one thread is blocked and waiting, a second thread in the same task can run.   |
| 13.  | Changes to the parent process does not affect child processes.                                     | Since all threads of the same process share address space and other resources so any changes to the main thread may affect the behavior of the other threads of the process. |

## Inter Process Communication (IPC)

*A process can be of two types:*

- *Independent process.*
- *Co-operating process.*

*An independent process is not affected by the execution of other processes while a co-operating process can be affected by other executing processes. Though one can think that those processes, which are running independently, will execute very efficiently, in reality, there are many situations when co-operative nature can be utilized for increasing computational speed, convenience, and modularity . Inter-process communication (IPC) is a mechanism that allows processes to communicate with each other and synchronize their actions. The communication between these*



processes can be seen as a method of co-operation between them. Processes can communicate with each other through both:

1. Shared Memory
2. Message passing

Figure 1 below shows a basic structure of communication between processes via the shared memory method and via the message passing method.

An operating system can implement both methods of communication. First, we will discuss the shared memory methods of communication and then message passing. Communication between processes using shared memory requires processes to share some variable, and it completely depends on how the programmer will implement it. One way of communication using shared memory can be imagined like this: Suppose process1 and process2 are executing simultaneously, and they share some resources or use some information from another process. Process1 generates information about certain computations or resources being used and keeps it as a record in shared memory. When process2 needs to use the shared information, it will check in the record stored in shared memory and take note of the information generated by process1 and act accordingly. Processes can use shared memory for extracting information as a record from another process as well as for delivering any specific information to other processes.

Let's discuss an example of communication between processes using the shared memory method.

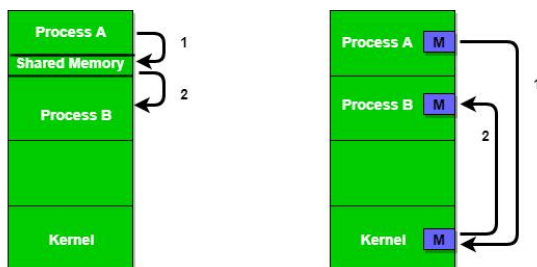


Figure 1 - Shared Memory and Message Passing

### i) Shared Memory Method

#### Ex: Producer-Consumer problem

There are two processes: Producer and Consumer. The producer produces some items and the Consumer consumes that item. The two processes share a common space or memory location known as a buffer where the item produced by the Producer is stored and from which the Consumer consumes the item if needed. There are two versions of this problem:

### ii) Messaging Passing Method

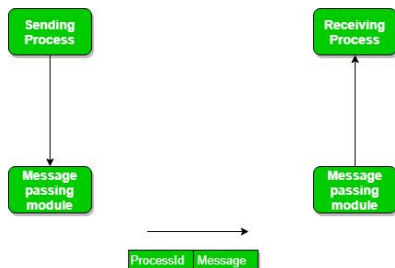
Now, We will start our discussion of the communication between processes via message passing. In this method, processes communicate with each other without using any kind of shared memory. If

two processes p1 and p2 want to communicate with each other, they proceed as follows:

- Establish a communication link (if a link already exists, no need to establish it again.)
- Start exchanging messages using basic primitives.

We need at least two primitives:

- `send(message, destination)` or `send(message)`
- `receive(message, host)` or `receive(message)`



The message size can be of fixed size or of variable size. If it is of fixed size, it is easy for an OS designer but complicated for a programmer and if it is of variable size then it is easy for a programmer but complicated for the OS designer. A standard message can have two parts: **header** and **body**.

## Introduction of Process Synchronization

On the basis of synchronization, processes are categorized as one of the following two types:

- **Independent Process** : Execution of one process does not affects the execution of other processes.
- **Cooperative Process** : Execution of one process affects the execution of other processes.

Process synchronization problem arises in the case of Cooperative process also because resources are shared in Cooperative processes.

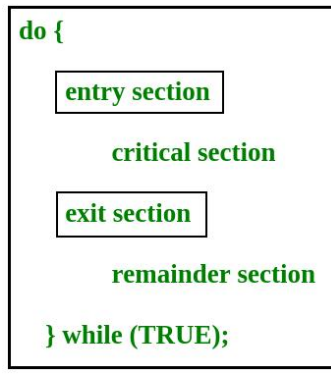
### Critical Section Problem

When more than one processes access a same code segment that segment is known as critical section. Critical section contains shared variables or resources which are needed to be synchronized to maintain consistency of data variable.

In simple terms a critical section is group of instructions/statements or region of code that need to be executed atomically ([read this post](#) for atomicity), such as accessing a resource (file, input or output port, global data, etc.).



Critical section is a code segment that can be accessed by only one process at a time. Critical section contains shared variables which need to be synchronized to maintain consistency of data variables.



In the entry section, the process requests for entry in the **Critical Section**.

Any solution to the critical section problem must satisfy three requirements:

- **Mutual Exclusion** : If a process is executing in its critical section, then no other process is allowed to execute in the critical section.
- **Progress** : If no process is executing in the critical section and other processes are waiting outside the critical section, then only those processes that are not executing in their remainder section can participate in deciding which will enter in the critical section next, and the selection can not be postponed indefinitely.
- **Bounded Waiting** : A bound must exist on the number of times that other processes are allowed to enter their critical sections after a process has made a request to enter its critical section and before that request is granted.

**Semaphores:-**

A semaphore is a signaling mechanism and a thread that is waiting on a semaphore can be signaled by another thread. This is different than a mutex as the mutex can be signaled only by the thread that called the wait function.

A semaphore uses two atomic operations, wait and signal for process synchronization.

A Semaphore is an integer variable, which can be accessed only through two operations wait() and signal().

There are two types of semaphores: **Binary Semaphores** and **Counting Semaphores**

1. **Binary Semaphore –**

This is also known as mutex lock. It can have only two values – 0 and 1. Its value is initialized to 1. It is used to implement the solution of critical section problems with multiple processes.

2. **Counting Semaphore –**

Its value can range over an unrestricted domain. It is used to control access to a resource that has multiple instances.

## CPU Scheduling in Operating Systems

Scheduling of processes/work is done to finish the work on time.

Below are different time with respect to a process.

*Arrival Time: Time at which the process arrives in the ready queue.*

*Completion Time: Time at which process completes its execution.*

*Burst Time: Time required by a process for CPU execution.*

*Turn Around Time: Time Difference between completion time and arrival time.*

*Turn Around Time = Completion Time – Arrival Time*

*Waiting Time(W.T): Time Difference between turn around time and burst time.*

*Waiting Time = Turn Around Time – Burst Time*

### Why do we need scheduling?

A typical process involves both I/O time and CPU time. In a uni programming system like MS-DOS, time spent waiting for I/O is wasted and CPU is free during this time. In multi programming systems, one process can use CPU while another is waiting for I/O. This is possible only with process scheduling.

### Objectives of Process Scheduling Algorithm

*Max CPU utilization [Keep CPU as busy as possible]*

*Fair allocation of CPU.*

*Max throughput [Number of processes that complete their execution per time unit]*

*Min turnaround time [Time taken by a process to finish execution]*

*Min waiting time [Time a process waits in ready queue]*

*Min response time [Time when a process produces first response]*

## Different Scheduling Algorithms

First Come First Serve (FCFS): Simplest scheduling algorithm that schedules according to arrival times of processes. First come first serve scheduling algorithm states that the process that requests the CPU first is allocated the CPU first. It is implemented by using the 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. FCFS is a non-preemptive scheduling algorithm.

**Note:**First come first serve suffers from convoy effect.

Shortest Job First (SJF): Process which have the shortest burst time are scheduled first. If two processes have the same burst time then FCFS is used to break the tie. It is a non-preemptive scheduling algorithm.

Round Robin Scheduling: Each process is assigned a fixed time (Time Quantum/Time Slice) in cyclic way. It is designed especially for the time-sharing system. 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 Round Robin 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.

Priority Based scheduling (Non-Preemptive): In this scheduling, processes are scheduled according to their priorities, i.e., highest priority process is scheduled first. If priorities of two processes match, then schedule according to arrival time. Here starvation of process is possible.

Multilevel Queue Scheduling: According to the priority of process, processes are placed in the different queues. Generally high priority process are placed in the top level queue. Only after completion of processes from top level queue, lower level queued processes are scheduled. It can suffer from starvation.

Multi level Feedback Queue Scheduling: It allows the process to move in 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 is moved to a lower-priority queue.

## Preemptive and Non-Preemptive Scheduling

### 1. Preemptive Scheduling:

Preemptive scheduling is used when a process switches from running state to ready state or from the waiting state to ready state. The resources (mainly CPU cycles) are allocated to the process for a limited amount of time and then taken away, and the process is again placed back in the ready queue if that process still has CPU burst time remaining. That process stays in the ready queue till it gets its next chance to execute.

Algorithms based on preemptive scheduling are: Round Robin (RR), Shortest Remaining Time First (SRTF), Priority (preemptive version), etc.

| Process | Arrival Time | CPU Burst Time (in millisec.) |
|---------|--------------|-------------------------------|
| P0      | 3            | 2                             |
| P1      | 2            | 4                             |
| P2      | 0            | 6                             |
| P3      | 1            | 4                             |



**Preemptive Scheduling**

### 2. Non-Preemptive Scheduling:

Non-preemptive Scheduling is used when a process terminates, or a process switches from running to the waiting state. In this scheduling, once the resources (CPU cycles) are allocated to a process, the process holds the CPU till it gets terminated or

reaches a waiting state. In the case of non-preemptive scheduling does not interrupt a process running CPU in the middle of the execution. Instead, it waits till the process completes its CPU burst time, and then it can allocate the CPU to another process.

Algorithms based on non-preemptive scheduling are: Shortest Job First (SJF basically non preemptive) and Priority (non preemptive version), etc.

| Process | Arrival Time | CPU Burst Time (in millisec.) |
|---------|--------------|-------------------------------|
| P0      | 3            | 2                             |
| P1      | 2            | 4                             |
| P2      | 0            | 6                             |
| P3      | 1            | 4                             |



**Non-Preemptive Scheduling**

### Comparison Chart:

| Parameter       | PREEMPTIVE SCHEDULING   | NON-PREEMPTIVE SCHEDULING   |
|-----------------|---|---|
| Basic           | In this resources(CPU Cycle) are allocated to a process for a limited time.                                 | Once resources(CPU Cycle) are allocated to a process, the process holds it till it completes its burst time or switches to waiting state. |
| Interrupt       | Process can be interrupted in between.  | Process can not be interrupted until it terminates itself or its time is up.  |
| Starvation      | If a process having high priority frequently arrives in the ready queue, a low priority process may starve. | If a process with a long burst time is running CPU, then later coming process with less CPU burst time may starve.                        |
| Overhead        | It has overheads of scheduling the processes.   | It does not have overheads.   |
| Flexibility     | flexible  | rigid   |
| Cost            | cost associated   | no cost associated  |
| CPU Utilization | In preemptive scheduling, CPU utilization is high.  | It is low in non preemptive scheduling.   |
| Examples        | Examples of preemptive scheduling are Round Robin and Shortest Remaining Time First.                        | Examples of non-preemptive scheduling are First Come First Serve and Shortest Job First.  |