# RAJIV GANDHI PROUDYOGIKI VISHWAVIDYALAYA, BHOPAL

## New Scheme Based On AICTE Flexible Curricula Computer Science & Engineering, VI-Semester CSE-405 [Operating System]
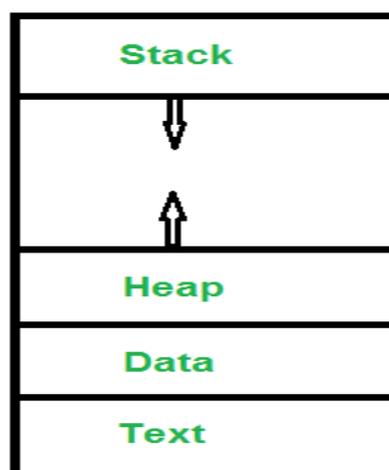
## Topic Covered

## Process Management

## 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, as opposed to a program, which is considered to be 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).

### What does a process look like in memory?

*Text Section*: A Process, sometimes known as the Text Section, also includes the current activity represented by the value of the *Program Counter*.

*Stack*: The Stack contains the temporary data, such as function parameters, returns addresses, and local variables.

*Data Section*: Contains the global variable.

*Heap Section*: Dynamically allocated memory to process during its run time.

Refer this for more details on sections.

## Attributes or Characteristics of a Process

A process has following attributes.

**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

**7. CPU scheduling information:** For example, Priority (Different processes may have different priorities, for example a short process may be assigned a low priority in the shortest job first scheduling)

All of the above attributes of a process are also known as the *context of the process*. Every process has its own program control block(PCB), i.e each process will have a unique PCB. All of the above attributes are part of the PCB.
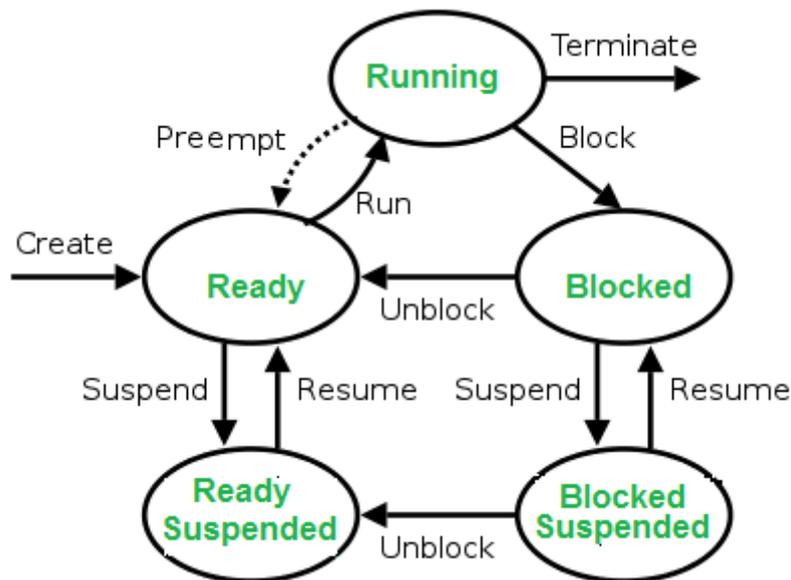
## States of Process:

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.



## Process Scheduling

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.

Process scheduling is an essential part of a Multiprogramming operating systems. Such operating systems allow more than one process to be loaded into the executable memory at a time and the loaded process shares the CPU using time multiplexing.

## Process Scheduling Queues

The OS maintains all PCBs in Process Scheduling Queues. The OS maintains a separate queue for each of the process states and PCBs of all processes in the same execution state are placed in the same queue. When the state of a process is changed, its PCB is unlinked from its current queue and moved to its new state queue.

The Operating System maintains the following important process scheduling queues −

- **Job queue** − This queue keeps all the processes in the system.

- **Ready queue** − This queue keeps a set of all processes residing in main memory, ready and waiting to execute. A new process is always put in this queue.

- **Device queues** − The processes which are blocked due to unavailability of an I/O device constitute this queue.

The OS can use different policies to manage each queue (FIFO, Round Robin, Priority, etc.). The OS scheduler determines how to move processes between the ready and run queues which can only have one entry per processor core on the system; in the above diagram, it has been merged with the CPU.

Two-State Process Model

Two-state process model refers to running and non-running states which are described below −

| S.N. | State & Description |
|------|---------------------|
| 1 | **Running** <br><br> When a new process is created, it enters into the system as in the running state. |
| 2 | **Not Running** <br><br> Processes that are not running are kept in queue, waiting for their turn to execute. Each entry in the queue is a pointer to a particular process. Queue is implemented by using linked list. Use of dispatcher is as follows. When a process is interrupted, that process is transferred in the waiting queue. If the process has completed or aborted, the process is discarded. In either case, the dispatcher then selects a process from the queue to execute. |

## Schedulers

Schedulers are special system software which handle process scheduling in various ways. Their main task is to select the jobs to be submitted into the system and to decide which process to run. Schedulers are of three types −

- Long-Term Scheduler
- Short-Term Scheduler
- Medium-Term Scheduler

**Long Term Scheduler**

It is also called a **job scheduler**. A long-term scheduler determines which programs are admitted to the system for processing. It selects processes from the queue and loads them into memory for execution. Process loads into the memory for CPU scheduling.

The primary objective of the job scheduler is to provide a balanced mix of jobs, such as I/O bound and processor bound. It also controls the degree of multiprogramming. If the degree of multiprogramming is stable, then the average rate of process creation must be equal to the average departure rate of processes leaving the system.

On some systems, the long-term scheduler may not be available or minimal. Time-sharing operating systems have no long term scheduler. When a process changes the state from new to ready, then there is use of long-term scheduler.

**Short Term Scheduler**

It is also called as **CPU scheduler**. Its main objective is to increase system performance in accordance with the chosen set of criteria. It is the change of ready state to running state of the process. CPU scheduler selects a process among the processes that are ready to execute and allocates CPU to one of them.

Short-term schedulers, also known as dispatchers, make the decision of which process to execute next. Short-term schedulers are faster than long-term schedulers.

**Medium Term Scheduler**

Medium-term scheduling is a part of **swapping**. It removes the processes from the memory. It reduces the degree of multiprogramming. The medium-term scheduler is in-charge of handling the swapped out-processes.

A running process may become suspended if it makes an I/O request. A suspended processes cannot make any progress towards completion. In this condition, to remove the process from memory and make space for other processes, the suspended process is moved to the

secondary storage. This process is called **swapping**, and the process is said to be swapped out or rolled out. Swapping may be necessary to improve the process mix.

## Comparison among Scheduler

| S.N. | Long-Term Scheduler | Short-Term Scheduler | Medium-Term Scheduler |
|------|---------------------|----------------------|------------------------|
| 1 | It is a job scheduler | It is a CPU scheduler | It is a process swapping scheduler. |
| 2 | Speed is lesser than short term scheduler | Speed is fastest among other two | Speed is in between both short and long term scheduler. |
| 3 | It controls the degree of multiprogramming | It provides lesser control over degree of multiprogramming | It reduces the degree of multiprogramming. |
| 4 | It is almost absent or minimal in time sharing system | It is also minimal in time sharing system | It is a part of Time sharing systems. |
| 5 | It selects processes from pool and loads them into memory for execution | It selects those processes which are ready to execute | It can re-introduce the process into memory and execution can be continued. |

**Context Switch**

A context switch is the mechanism to store and restore the state or context of a CPU in Process Control block so that a process execution can be resumed from the same point at a later time. Using this technique, a context switcher enables multiple processes to share a

single CPU. Context switching is an essential part of a multitasking operating system features.

When the scheduler switches the CPU from executing one process to execute another, the state from the current running process is stored into the process control block. After this, the state for the process to run next is loaded from its own PCB and used to set the PC, registers,

etc. At that point, the second process can start executing.

Context switches are computationally intensive since register and memory state must be saved and restored. To avoid the amount of context switching time, some hardware systems employ two or more sets of processor registers. When the process is switched, the following information is stored for later use.

- Program Counter

- Scheduling information

- Base and limit register value

- Currently used register

- Changed State

- I/O State information

- Accounting information

## Operating System Scheduling algorithms

A Process Scheduler schedules different processes to be assigned to the CPU based on particular scheduling algorithms. There are six popular process scheduling algorithms which we are going to discuss in this chapter −
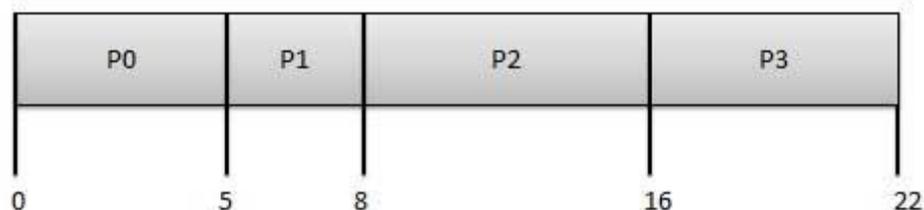
- First-Come, First-Served (FCFS) Scheduling

- Shortest-Job-Next (SJN) Scheduling

- Priority Scheduling

- Shortest Remaining Time

- Round Robin(RR) Scheduling

- Multiple-Level Queues Scheduling

These algorithms are either **non-preemptive or preemptive**. Non-preemptive algorithms are designed so that once a process enters the running state, it cannot be preempted until it completes its allotted time, whereas the preemptive scheduling is based on priority where a scheduler may preempt a low priority running process anytime when a high priority process enters into a ready state.

## First Come First Serve (FCFS)

- Jobs are executed on first come, first serve basis.

- It is a non-preemptive, pre-emptive scheduling algorithm.

- Easy to understand and implement.

- Its implementation is based on FIFO queue.

- Poor in performance as average wait time is high.

| Process | Arrival Time | Execute Time | Service Time |
|---------|--------------|--------------|--------------|
| P0 | 0 | 5 | 0 |
| P1 | 1 | 3 | 5 |
| P2 | 2 | 8 | 8 |
| P3 | 3 | 6 | 16 |

| P0 | P1 | P2 | P3 |
|----|----|----|----|

0   5   8        16        22

**Wait time** of each process is as follows −

| Process | Wait Time : Service Time - Arrival Time |
|---------|-----------------------------------------|
| P0 | 0 - 0 = 0 |
| P1 | 5 - 1 = 4 |

| | |
|---|---|
| P2 | 8 - 2 = 6 |
| P3 | 16 - 3 = 13 |

Average Wait Time: (0+4+6+13) / 4 = 5.75

## Shortest Job Next (SJN)

- This is also known as **shortest job first**, or SJF

- This is a non-preemptive, pre-emptive scheduling algorithm.

- Best approach to minimize waiting time.

- Easy to implement in Batch systems where required CPU time is known in advance.

- Impossible to implement in interactive systems where required CPU time is not known.

- The processer should know in advance how much time process will take.

Given: Table of processes, and their Arrival time, Execution time

| Process | Arrival Time | Execution Time | Service Time |
|---------|-------------|----------------|--------------|
| P0 | 0 | 5 | 0 |
| P1 | 1 | 3 | 5 |
| P2 | 2 | 8 | 14 |
| P3 | 3 | 6 | 8 |

**Waiting time** of each process is as follows −

| Process | Waiting Time |
|---------|--------------|
| P0 | 0 - 0 = 0 |
| P1 | 5 - 1 = 4 |
| P2 | 14 - 2 = 12 |
| P3 | 8 - 3 = 5 |

Average Wait Time: (0 + 4 + 12 + 5)/4 = 21 / 4 = 5.25

## Priority Based Scheduling

- Priority scheduling is a non-preemptive algorithm and one of the most common scheduling algorithms in batch systems.

- Each process is assigned a priority. Process with highest priority is to be executed first and so on.

- Processes with same priority are executed on first come first served basis.

- Priority can be decided based on memory requirements, time requirements or any other resource requirement.

Given: Table of processes, and their Arrival time, Execution time, and priority. Here we are considering 1 is the lowest priority.

| Process | Arrival Time | Execution Time | Priority | Service Time |
|---------|--------------|----------------|----------|--------------|
| P0 | 0 | 5 | 1 | 0 |
| P1 | 1 | 3 | 2 | 11 |

| | | | | |
|----|----|----|----|----|
| P2 | 2 | 8 | 1 | 14 |
| P3 | 3 | 6 | 3 | 5 |

**Waiting time** of each process is as follows −

| Process | Waiting Time |
|---------|--------------|
| P0 | 0 - 0 = 0 |
| P1 | 11 - 1 = 10 |
| P2 | 14 - 2 = 12 |
| P3 | 5 - 3 = 2 |

Average Wait Time: (0 + 10 + 12 + 2)/4 = 24 / 4 = 6
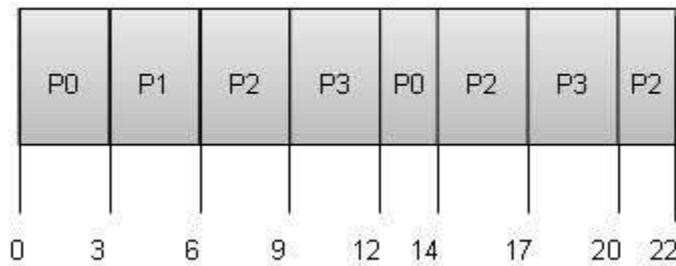
## Shortest Remaining Time

- Shortest remaining time (SRT) is the preemptive version of the SJN algorithm.

- The processor is allocated to the job closest to completion but it can be preempted by a newer ready job with shorter time to completion.

- Impossible to implement in interactive systems where required CPU time is not known.

- It is often used in batch environments where short jobs need to give preference.

## Round Robin Scheduling

- Round Robin is the preemptive process scheduling algorithm.

- Each process is provided a fix time to execute, it is called a **quantum**.

- Once a process is executed for a given time period, it is preempted and other process executes for a given time period.

- Context switching is used to save states of preempted processes.

Quantum = 3



**Wait time** of each process is as follows −

| Process | Wait Time : Service Time - Arrival Time |
|---------|------------------------------------------|
| P0 | (0 - 0) + (12 - 3) = 9 |
| P1 | (3 - 1) = 2 |
| P2 | (6 - 2) + (14 - 9) + (20 - 17) = 12 |
| P3 | (9 - 3) + (17 - 12) = 11 |

Average Wait Time: (9+2+12+11) / 4 = 8.5

## Multiple-Level Queues Scheduling

Multiple-level queues are not an independent scheduling algorithm. They make use of other existing algorithms to group and schedule jobs with common characteristics.

- Multiple queues are maintained for processes with common characteristics.
- Each queue can have its own scheduling algorithms.
- Priorities are assigned to each queue.
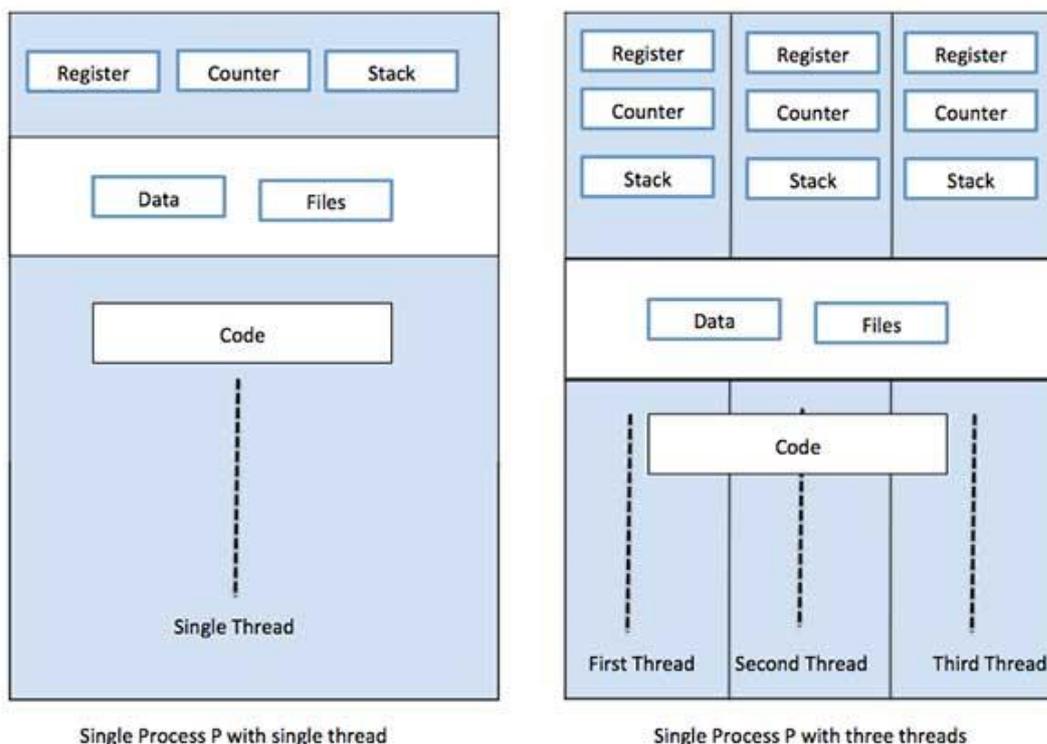
## Operating System - Multi-Threading

Thread

A thread is a flow of execution through the process code, with its own program counter that keeps track of which instruction to execute next, system registers which hold its current working variables, and a stack which contains the execution history.

A thread shares with its peer threads few information like code segment, data segment and open files. When one thread alters a code segment memory item, all other threads see that.

A thread is also called a **lightweight process**. Threads provide a way to improve application performance through parallelism. Threads represent a software approach to improving performance of operating system by reducing the overhead thread is equivalent to a classical process.

Each thread belongs to exactly one process and no thread can exist outside a process. Each thread represents a separate flow of control. Threads have been successfully used in implementing network servers and web server. They also provide a suitable foundation for parallel execution of applications on shared memory multiprocessors. The following figure shows the working of a single-threaded and a multithreaded process.



Single Process P with single thread                Single Process P with three threads

## Difference between Process and Thread

| S.N. | Process | Thread |
|------|---------|--------|
| 1 | Process is heavy weight or resource intensive. | Thread is light weight, taking lesser resources than a process. |
| 2 | Process switching needs interaction with operating system. | Thread switching does not need to interact with operating system. |
| 3 | In multiple processing environments, each process executes the same code but has its own memory and file resources. | All threads can share same set of open files, child processes. |
| 4 | 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. |
| 5 | Multiple processes without using threads use more resources. | Multiple threaded processes use fewer resources. |
| 6 | In multiple processes each process operates independently of the others. | One thread can read, write or change another thread's data. |

## Advantages of Thread

- Threads minimize the context switching time.
- Use of threads provides concurrency within a process.
- Efficient communication.
- It is more economical to create and context switch threads.

- Threads allow utilization of multiprocessor architectures to a greater scale and efficiency.
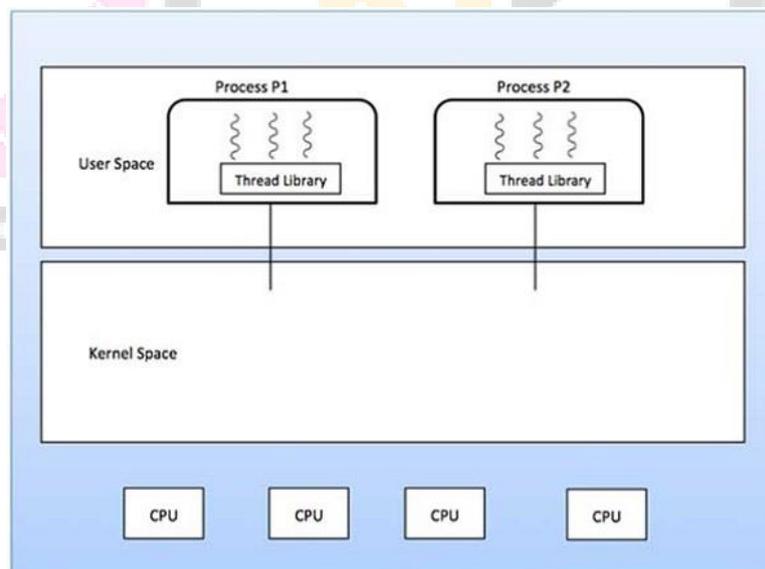
## Types of Thread

Threads are implemented in following two ways −

- **User Level Threads** − User managed threads.

- **Kernel Level Threads** − Operating System managed threads acting on kernel, an operating system core.

## User Level Threads

In this case, the thread management kernel is not aware of the existence of threads. The thread library contains code for creating and destroying threads, for passing message and data between threads, for scheduling thread execution and for saving and restoring thread contexts. The application starts with a single thread.



## Advantages

- Thread switching does not require Kernel mode privileges.

- User level thread can run on any operating system.

- Scheduling can be application specific in the user level thread.

- User level threads are fast to create and manage.

Disadvantages

- In a typical operating system, most system calls are blocking.

- Multithreaded application cannot take advantage of multiprocessing.

## Kernel Level Threads

In this case, thread management is done by the Kernel. There is no thread management code in the application area. Kernel threads are supported directly by the operating system. Any application can be programmed to be multithreaded. All of the threads within an application are supported within a single process.

The Kernel maintains context information for the process as a whole and for individuals threads within the process. Scheduling by the Kernel is done on a thread basis. The Kernel performs thread creation, scheduling and management in Kernel space. Kernel threads are generally slower to create and manage than the user threads.

Advantages

- Kernel can simultaneously schedule multiple threads from the same process on multiple processes.

- If one thread in a process is blocked, the Kernel can schedule another thread of the same process.

- Kernel routines themselves can be multithreaded.

Disadvantages

- Kernel threads are generally slower to create and manage than the user threads.

- Transfer of control from one thread to another within the same process requires a mode switch to the Kernel.

## Multithreading Models

Some operating system provide a combined user level thread and Kernel level thread facility. Solaris is a good example of this combined approach. In a combined system,

multiple threads within the same application can run in parallel on multiple processors and a blocking system call need not block the entire process.

**Difference between User-Level & Kernel-Level Thread**

| S.N. | User-Level Threads | Kernel-Level Thread |
|------|--------------------|---------------------|
| 1 | User-level threads are faster to create and manage. | Kernel-level threads are slower to create and manage. |
| 2 | Implementation is by a thread library at the user level. | Operating system supports creation of Kernel threads. |
| 3 | User-level thread is generic and can run on any operating system. | Kernel-level thread is specific to the operating system. |
| 4 | Multi-threaded applications cannot take advantage of multiprocessing. | Kernel routines themselves can be multithreaded. |

## References

- https://www.tutorialspoint.com/operating_system
- https://www.javatpoint.com/os-tutorial
- https://www.studytonight.com/operating-system/
- https://www.geeksforgeeks.org/operating-systems/
- Silberschatz, Galvin, Gagne, "Operating System Concepts'', Wiley, 9/E
- William Stalling, "Operating Systems", Pearson Education
- Andrew S. Tanenbaum, "Modern Operating Systems", 3/e, Prentice Hall