

# CPU Scheduling- Intro

---





# Chapter 6: CPU Scheduling

---

- ❑ Basic Concepts
- ❑ Scheduling Criteria
- ❑ Scheduling Algorithms
- ❑ Thread Scheduling
- ❑ Multiple-Processor Scheduling
- ❑ Real-Time CPU Scheduling
- ❑ Operating Systems Examples
- ❑ Algorithm Evaluation





# Objectives

---

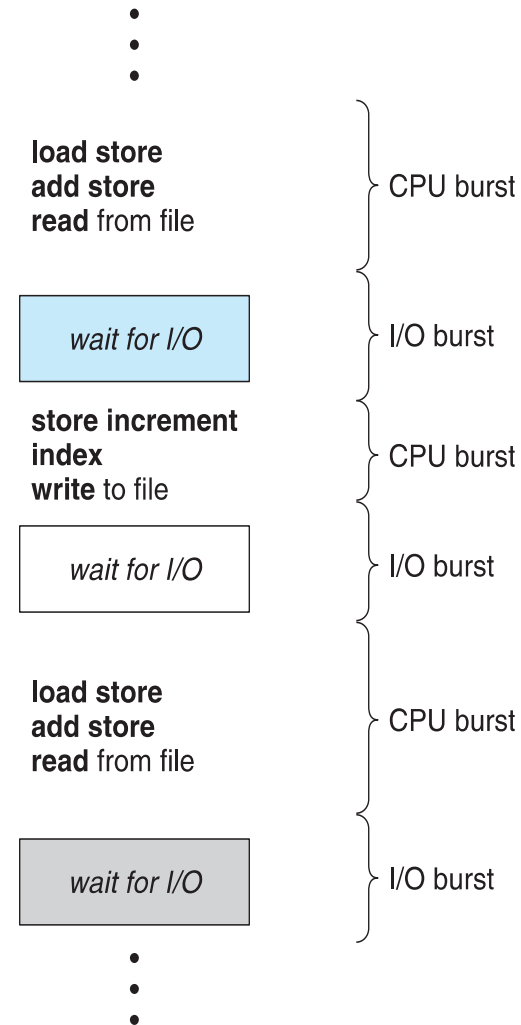
- To introduce CPU scheduling, which is the basis for multiprogrammed operating systems
- To describe various CPU-scheduling algorithms
- To discuss evaluation criteria for selecting a CPU-scheduling algorithm for a particular system
- To examine the scheduling algorithms of several operating systems





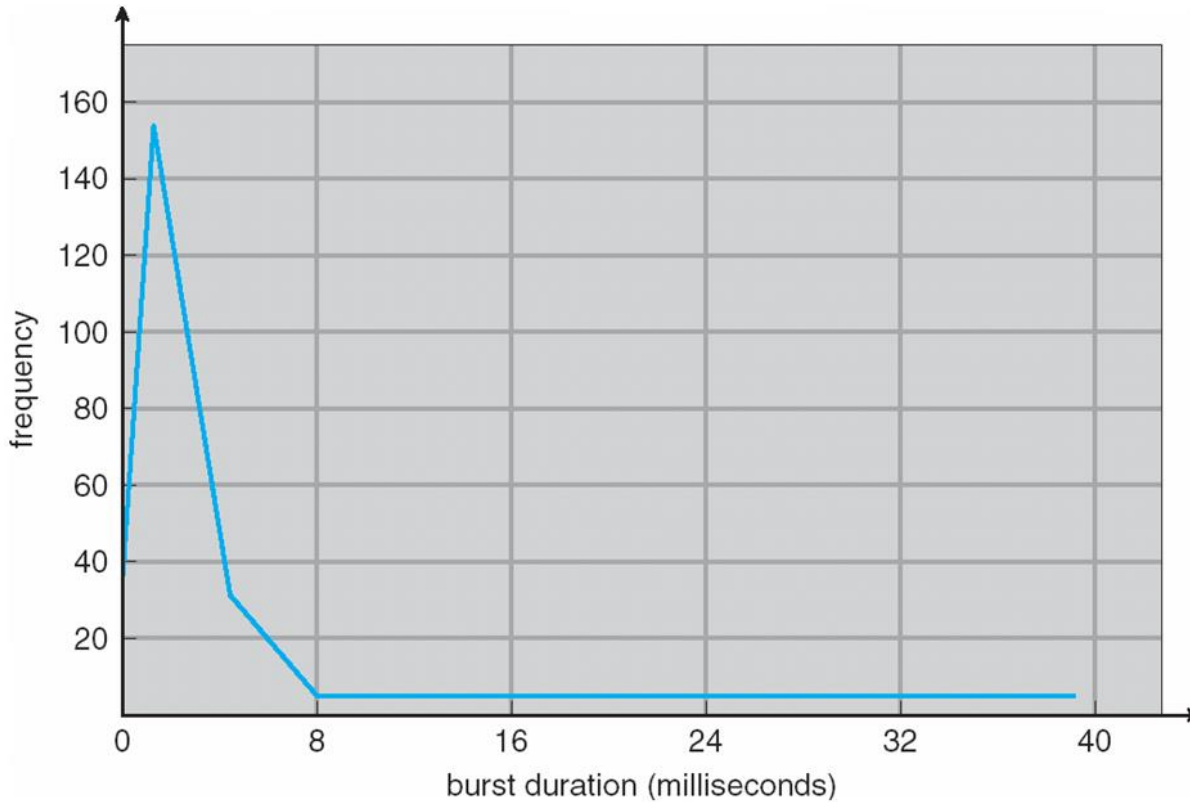
# Basic Concepts

- ❑ Maximum CPU utilization obtained with multiprogramming
- ❑ CPU–I/O Burst Cycle – Process execution consists of a **cycle** of CPU execution and I/O wait
- ❑ **CPU burst** followed by **I/O burst**
- ❑ CPU burst distribution is of main concern





# Histogram of CPU-burst Times





# CPU Scheduler

- **Short-term scheduler** selects from among the processes in ready queue, and allocates the CPU to one of them
  - Queue may be ordered in various ways
- CPU scheduling decisions may take place when a process:
  1. Switches from running to waiting state
  2. Switches from running to ready state
  3. Switches from waiting to ready
  4. Terminates
- Scheduling under 1 and 4 is **nonpreemptive**
- All other scheduling is **preemptive**
  - Consider access to shared data
  - Consider preemption while in kernel mode
  - Consider interrupts occurring during crucial OS activities





# Dispatcher

---

- ❑ Dispatcher module gives control of the CPU to the process selected by the short-term scheduler; this involves:
  - ❑ switching context
  - ❑ switching to user mode
  - ❑ jumping to the proper location in the user program to restart that program
- ❑ **Dispatch latency** – time it takes for the dispatcher to stop one process and start another running





# Scheduling Criteria

---

- ❑ **CPU utilization** – keep the CPU as busy as possible
- ❑ **Throughput** – # of processes that complete their execution per time unit
- ❑ **Turnaround time** – amount of time to execute a particular process
- ❑ **Waiting time** – amount of time a process has been waiting in the ready queue
- ❑ **Response time** – amount of time it takes from when a request was submitted until the first response is produced, not output (for time-sharing environment)







# Scheduling Algorithm Optimization Criteria

---

- ❑ Max CPU utilization
- ❑ Max throughput
- ❑ Min turnaround time
- ❑ Min waiting time
- ❑ Min response time





# Basic Concepts

- ❑ Maximum CPU utilization obtained with multiprogramming
- ❑ A CPU bursts when it is executing instructions; an I/O system bursts when it services requests to fetch information.
- ❑ **CPU burst** distribution
- ❑ Turnaround Time= Turnaround time (TAT) is the time interval from the time of submission of a process to the time of the completion of the process
- ❑ **TURNAROUND TIME=WAITING TIME + SERVICE TIME**
- ❑ **Burst Time** = The slice that it gets, is called the **CPU burst**. In simple terms, the duration for which a process gets control of the CPU is the **CPU burst time**, and the concept of gaining control of the CPU is the **CPU burst**.
- ❑ **Waiting Time** = *Starting Time - Arrival Time*





# First Come First Serve Scheduling

- ❑ In the "First come first serve" scheduling algorithm, as the name suggests, the process which arrives first, gets executed first, or we can say that the process which requests the CPU first, gets the CPU allocated first.
- ❑ First Come First Serve, is just like **FIFO**(First in First out) Queue data structure, where the data element which is added to the queue first, is the one who leaves the queue first.
- ❑ This is used in Batch Systems.
- ❑ It's **easy to understand and implement** programmatically, using a Queue data structure, where a new process enters through the **tail** of the queue, and the scheduler selects process from the **head** of the queue.
- ❑ A perfect real life example of FCFS scheduling is **buying tickets at ticket counter**.



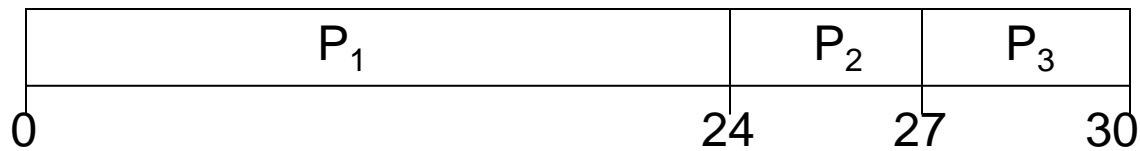


# First-Come, First-Served (FCFS) Scheduling

<u>Process</u>	<u>Burst Time</u>
$P_1$	24
$P_2$	3
$P_3$	3

- Suppose that the processes arrive in the order:  $P_1, P_2, P_3$   
The Gantt Chart for the schedule is:

□



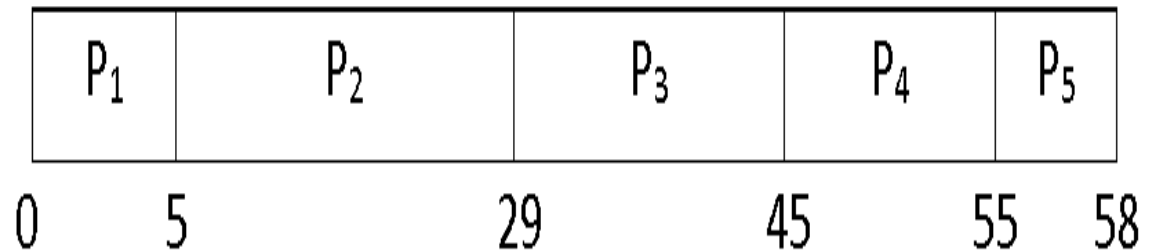
- Waiting time for  $P_1 = 0$ ;  $P_2 = 24$ ;  $P_3 = 27$
- Average waiting time:  $(0 + 24 + 27)/3 = 17$





# FCFS Algorithm

Process	Burst Time(ms)
P <sub>1</sub>	5
P <sub>2</sub>	24
P <sub>3</sub>	16
P <sub>4</sub>	10
P <sub>5</sub>	3





- Consider the above set of processes that arrive at time zero. The length of the CPU **burst time** given in millisecond. Now we calculate the average waiting time, average turnaround time.





## □ Average Waiting Time and Turnaround Time

### □ Average Waiting Time

- First of all, we have to calculate the waiting time of each process.

*Waiting Time = Starting Time - Arrival Time*

Waiting time of

$$P1 = 0$$

$$P2 = 5 - 0 = 5 \text{ ms}$$

$$P3 = 29 - 0 = 29 \text{ ms}$$

$$P4 = 45 - 0 = 45 \text{ ms}$$

$$P5 = 55 - 0 = 55 \text{ ms}$$

*Average Waiting Time = Waiting Time of all Processes / Total Number of Process*

Therefore, average waiting time =  $(0 + 5 + 29 + 45 + 55) / 5 = 25 \text{ ms}$





## □ Average Turnaround Time

- *Turnaround Time = Waiting time in the ready queue + executing time + waiting time in waiting-queue for I/O*

Turnaround time of

$$P1 = 0 + 5 + 0 = 5\text{ms}$$

$$P2 = 5 + 24 + 0 = 29\text{ms}$$

$$P3 = 29 + 16 + 0 = 45\text{ms}$$

$$P4 = 45 + 10 + 0 = 55\text{ms}$$

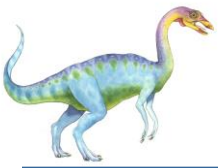
$$P5 = 55 + 3 + 0 = 58\text{ms}$$

$$\text{Total Turnaround Time} = (5 + 29 + 45 + 55 + 58)\text{ms} = 192\text{ms}$$

$$\text{Average Turnaround Time} = (\text{Total Turnaround Time} / \text{Total Number of Process}) = (192 / 5)\text{ms} = 38.4\text{ms}$$







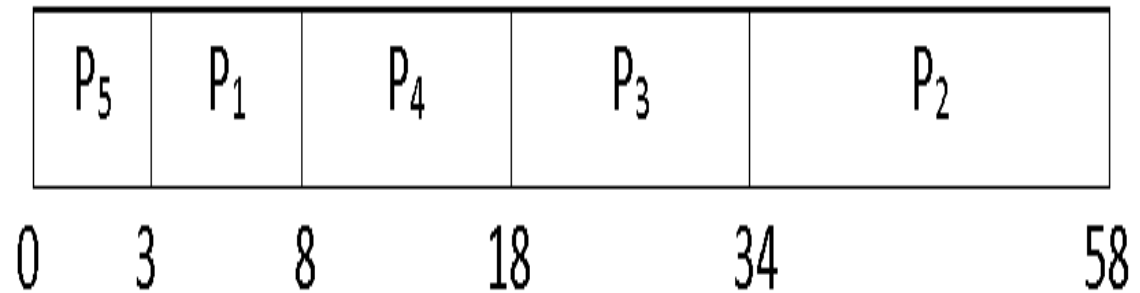
- ❑ Draw back
- ❑ What is Convoy Effect?
- ❑ Convoy Effect is a situation where many processes, who need to use a resource for short time are blocked by one process holding that resource for a long time.
- ❑ This essentially leads to poor utilization of resources and hence poor performance.





# SJF (Shortest Job First) Scheduling

Process	Burst Time(ms)
P <sub>1</sub>	5
P <sub>2</sub>	24
P <sub>3</sub>	16
P <sub>4</sub>	10
P <sub>5</sub>	3





## □ Average Waiting Time

- We will apply the same formula to find average waiting time in this problem. Here arrival time is common to all processes(i.e., zero).

Waiting Time for

$$P1 = 3 - 0 = 3\text{ms}$$

$$P2 = 34 - 0 = 34\text{ms}$$

$$P3 = 18 - 0 = 18\text{ms}$$

$$P4 = 8 - 0 = 8\text{ms}$$

$$P5 = 0\text{ms}$$

$$\text{Now, Average Waiting Time} = (3 + 34 + 18 + 8 + 0) / 5 = 12.6\text{ms}$$





## □ **Average Turnaround Time**

- According to the SJF Gantt chart and the turnaround time formulae,  
Turnaround Time of

$$P1 = 3 + 5 = 8\text{ms}$$

$$P2 = 34 + 24 = 58\text{ms}$$

$$P3 = 18 + 16 = 34\text{ms}$$

$$P4 = 8 + 10 = 18\text{ms}$$

$$P5 = 0 + 3 = 3\text{ms}$$

Therefore, Average Turnaround Time =  $(8 + 58 + 34 + 18 + 3) / 5 = 24.2\text{ms}$

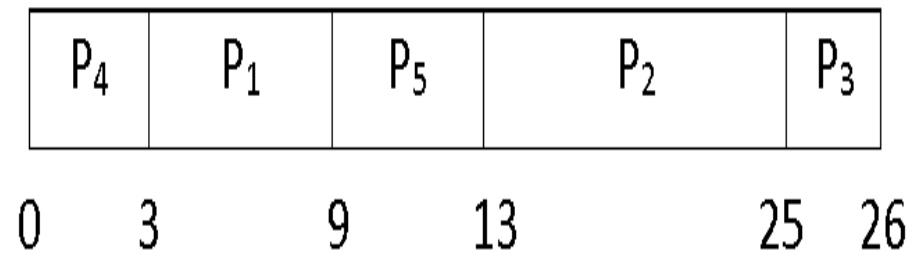


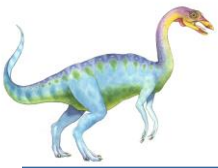


# Priority Scheduling Example

Process	CPU Burst Time	Priority
P <sub>1</sub>	6	2
P <sub>2</sub>	12	4
P <sub>3</sub>	1	5
P <sub>4</sub>	3	1
P <sub>5</sub>	4	3

Gantt Chart





- Processes with same priority are executed in FCFS manner.
- **Average Waiting Time**
- First of all, we have to find out the waiting time of each process.  
Waiting Time of process  
P1 = 3ms  
P2 = 13ms  
P3 = 25ms  
P4 = 0ms  
P5 = 9ms  
Therefore, Average Waiting Time =  $(3 + 13 + 25 + 0 + 9) / 5 = 10\text{ms}$





## □ **Average Turnaround Time**

- First finding Turnaround Time of each process.

Turnaround Time of process

$$P1 = (3 + 6) = 9\text{ms}$$

$$P2 = (13 + 12) = 25\text{ms}$$

$$P3 = (25 + 1) = 26\text{ms}$$

$$P4 = (0 + 3) = 3\text{ms}$$

$$P5 = (9 + 4) = 13\text{ms}$$

Therefore, Average Turnaround Time =  $(9 + 25 + 26 + 3 + 13) / 5 = 15.2\text{ms}$





# Round Robin (RR)

---

- Each process gets a small unit of CPU time (**time quantum**  $q$ ), usually 10-100 milliseconds. After this time has elapsed, the process is preempted and added to the end of the ready queue.
- Timer interrupts every quantum to schedule next process
- Performance
  - $q$  large  $\Rightarrow$  FIFO
  - $q$  small  $\Rightarrow q$  must be large with respect to context switch, otherwise overhead is too high

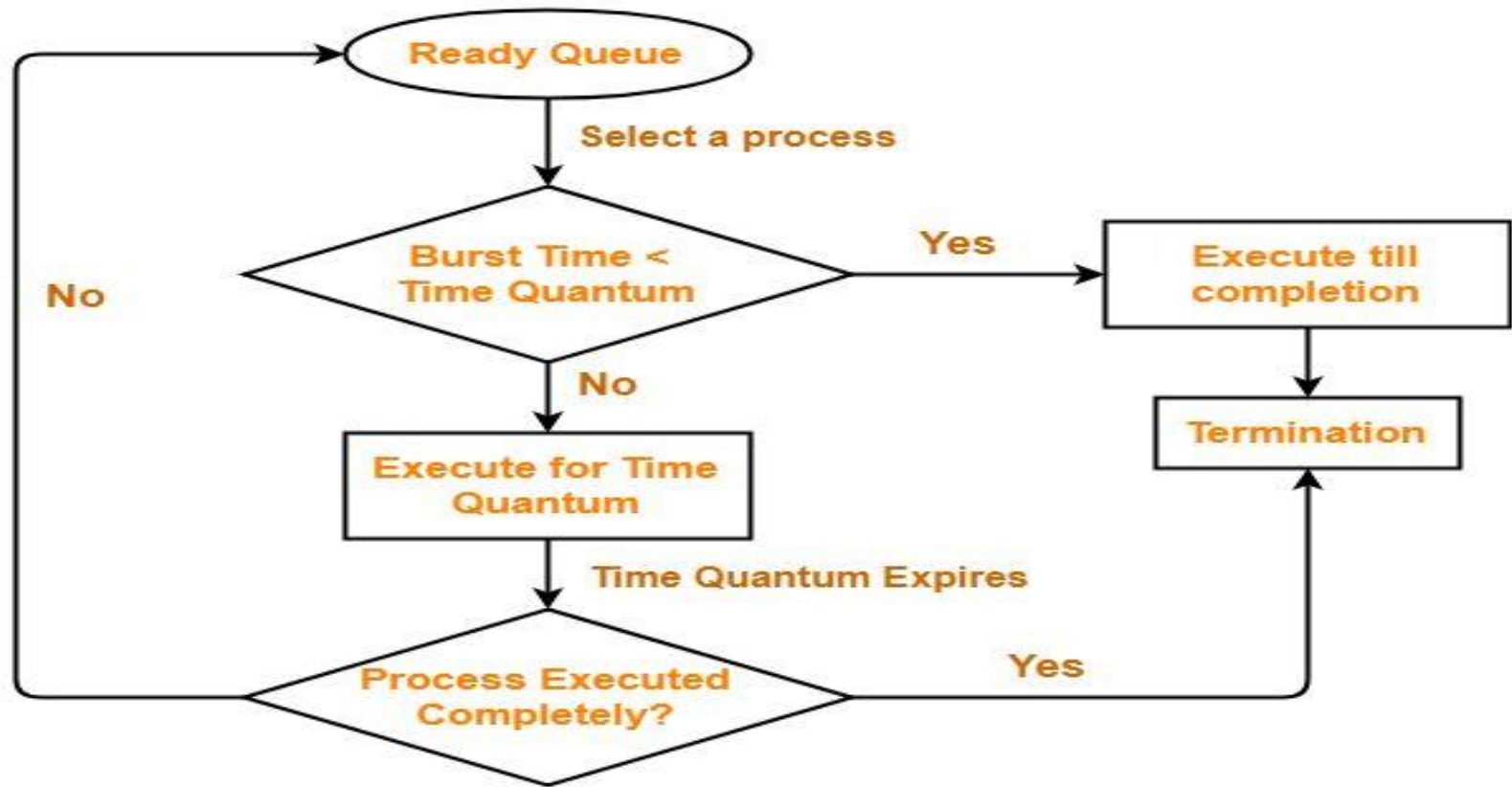






- ❑ CPU is assigned to the process on the basis of FCFS for a fixed amount of time.
- ❑ This fixed amount of time is called as **time quantum** or **time slice**.
- ❑ After the time quantum expires, the running process is preempted and sent to the ready queue.
- ❑ Then, the processor is assigned to the next arrived process.
- ❑ It is always preemptive in nature.
- ❑ Round Robin Scheduling is FCFS Scheduling with preemptive mode





### Round Robin Scheduling

