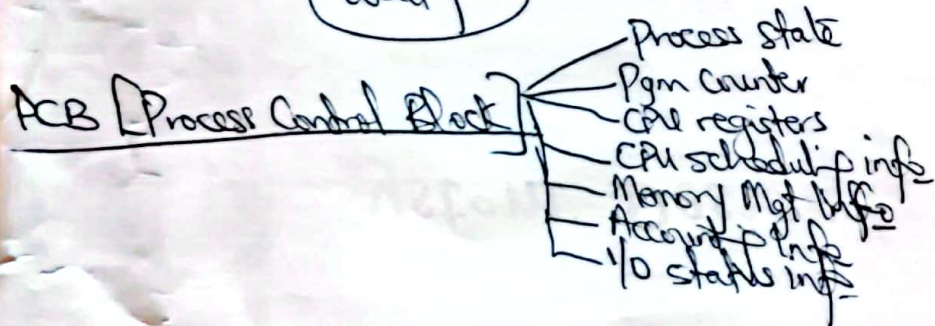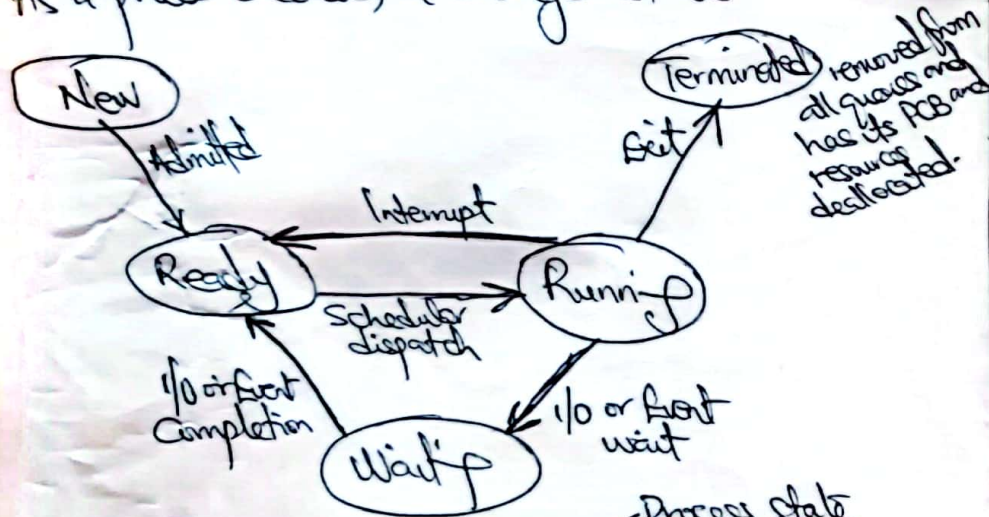# Process

- A program in execution
- While a pgm is just a passive collection of instructions. A process is the actual execution of those instructions.

- A process consists of the ffg resources:
  i) Memory
  ii) Operating system descriptors of resources that are allocated to the process
  iii) Security attributes
  iv) Processor state

# Process States

As a process executes, it changes states:



- Terminated: removed from all queues and has its PCB and resources deallocated.

# PCB [Process Control Block]

- Process state
- Pgm Counter
- CPU registers
- CPU scheduling info
- Memory Mgt Info
- Accounting info
- I/O status info

# Scheduling Queues

- As processes enter the system, they are put into a job queue
- Ready processes waiting to be executed are kept in a list called the ready queue, stored as linked list.
- The list of processes waiting for a particular I/O device is called a device queue.

# Schedulers

- Schedule processes into queues
- Long Term Scheduler (Job scheduler) — selects processes from a pool of processes and loads them into memory for execution.
- The Short-term Scheduler (CPU scheduler) selects from among the processes that are ready to execute, and allocates the CPU to one of them.
- The primary distinction b/w these 2 schedulers is the frequency of their execution.
- STS schedules more frequently and is faster while LTS executes much less frequently and controls the degree of multipgming i.e. the no. of processes in memory.

# CPU - I/O Burst Cycle

- Process execution consists of a cycle of CPU execution and I/O wait.
- Processes alternate b/w these 2 states.
- Process execution begins with a CPU burst followed by an I/O burst, and so on. The last CPU burst will end c a request to terminate execution.
- CPU burst can be extensively measured.
- An I/O-bound pgm ⇒ has many very short CPU bursts
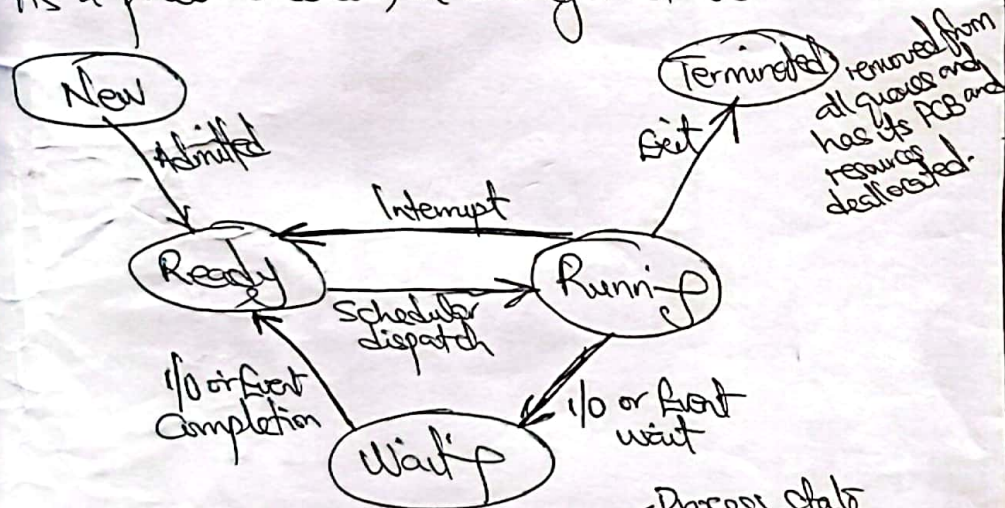- CPU-bound pgm ⇒ has a few very long CPU bursts

## Process

- A program in execution
- While a pgm is just a passive collection of instructions, A process is the actual execution of those instructions.
- A process consists of the ffg resources:
  (i) Memory
  (ii) Operating System descriptors of resources that are allocated to the process
  (iii) Security attributes
  (iv) Processor state

## Process States

As a process executes, it changes states:



New → Admitted → Ready
Ready → Scheduler dispatch → Running
Running → Interrupt → Ready
Running → Exit → Terminated (removed from all queues and has its PCB and resources deallocated.)
Running → I/O or event wait → Wait
Wait → I/O or event Completion → Ready

## PCB [Process Control Block]

- Process state
- Pgm Counter
- CPU registers
- CPU scheduling info
- Memory Mgt info
- Accounting info
- I/O status info

## Scheduling Queues

- As processes enter the system, they are put into a job queue
- Ready processes waiting to be executed are kept in a list called the ready queue, stored as linked list.
- The list of processes waiting for a particular I/O device is called a device queue.

## Schedulers

- Schedules processes into queues
- Long Term Scheduler (Job scheduler) selects processes from a pool of processes and loads them into memory for execution.
- The Short-term Scheduler (CPU scheduler) selects from among the processes that are ready to execute, and allocates the CPU to one of them.
- The primary distinction b/w these 2 schedulers is the frequency of their execution.
- STS schedules more frequently and is faster while LTS executes much less frequently and controls the degree of multipgming i.e. the no of processes in memory.

## CPU-I/O Burst Cycle

- Process execution consists of a cycle of CPU execution and I/O wait.
- Processes alternate b/w these 2 states.
- Process execution begins with a CPU burst, followed by an I/O burst, an so on. The last CPU burst will end c a request to terminate execution.
- CPU burst can be extensively measured.
- An I/O bound pgm ⇒ has many very short CPU bursts
- CPU-bound pgm ⇒ has a few very long CPU bursts

Preemptive Scheduling & Non-Preemptive scheduling

- Preemptive scheduling: Once CPU has been allocated to a process the process keeps the CPU until it releases the CPU either by terminating or by switching to the wait state.

1) Preemptive scheduling the CPU is forcefully taken away a process while it is running & state the of the process switches from the wait state (or when interrupts occurs) to ready state.
2) When a process switches from the wait state (completion of I/O) ready state (completion of I/O).
→ Scheduling of Preemptive scheduling.

Scheduling Criteria,
① CPU utilization
② Throughput
③ Turnaround Time → Time interval from submission to completion.
④ Waiting Time → The amount of time a process spends in the ready queue.
⑤ Response Time → Amount of time it takes to start responding but not the time it takes to output the response from the submission of a request.
⑥ Schedule the throughput. Deals with the problem of deciding which process is to be allocated the CPU.
FCFS scheduling → Avg. process that requests the CPU first is allocated the CPU first.

---

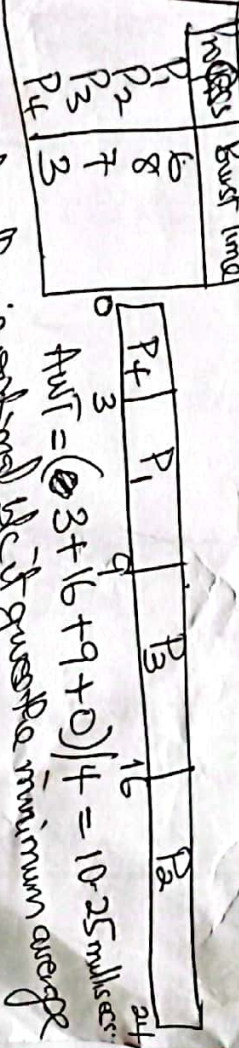- early managed with a FIFO queue. The Throughput & time under this policy is quite long.

| Process | Burst time |
|---------|-----------|
| P1      | 24        |
| P2      | 3         |
| P3      | 3         |



$AWT = (0 + 24 + 27)/3 = 17ms.$

Assume the processes arrive in the order $P_2, P_3, P_1$

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



$AWT = (6 + 0 + 3)/3 = 3 ms$

② SJF scheduling. with each process the length of the CPU burst. This algorithm associates with each process the length of the next CPU burst. When the CPU is available, it is assigned to the process that has the smallest next CPU burst, then use FCFS.

| Process | Burst time |
|---------|-----------|
| P1      | 6         |
| P2      | 8         |
| P3      | 7         |
| P4      | 3         |

| P4 | P1 | P3 | P2 |
|----|----|----|----|



$AWT = (3+16+9+0)/4 = 10.25 milliseconds$

- SJF algorithm is optimal it gives the minimum average waiting time for a given set of processes.
- It cannot be implemented at the level of short-term CPU scheduling. There is no way to know the length of the next CPU burst.
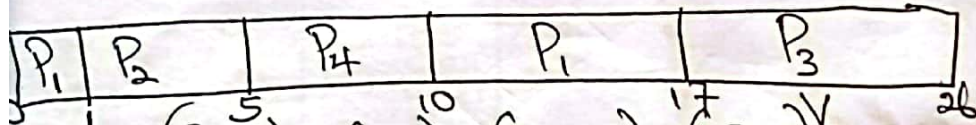- SJF scheduling is either preemptive or non-preemptive. Preemptive SJF scheduling is sometimes called Shortest remaining time-first scheduling.

Consider the ffg 4 processes, with length of the CPU-burst given in Milliseconds.

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

If the processes arrive at the ready queue at the times shown, then the Preemptive SJF is shown below in Gantt chart below :-

| $P_1$ | $P_2$ | $P_4$ | $P_1$ | $P_3$ |
|---|---|---|---|---|

1   5   10   17   26

$$AWT = \left[(10-1) + (1-1) + (17-2) + (5-3)\right]/4$$
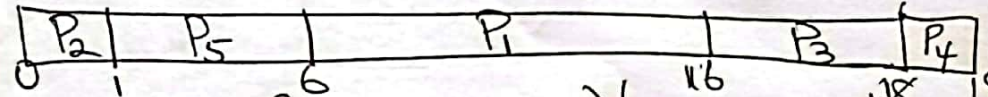
$$= 26/4 = 6.5 \text{ millisecs.}$$

A Non-Preemptive SJF Scheduling (Assignment)
(7.75 millisecs)

3) **Priority Scheduling**
- CPU is allocated to the process with the highest priority.
- Equal priority processes are scheduled in FCFS.
- The higher the CPU burst, the lower the priority & Vice versa.
- Low nos $\Rightarrow$ High priority.

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

Usip priority scheduling, the processes will be scheduled thus :

| $P_2$ | $P_5$ | $P_1$ | $P_3$ | $P_4$ |
|---|---|---|---|---|

0   1   6   16   18   19

$$AWT = (6+0+16+18+1)/5 = 8.2 \text{ milli seconds}$$

Disadvantage $\Rightarrow$ Indefinite block-p or starvation.
Soln $\Rightarrow$ Aging (a technique of gradually increas-p the priority of processes that wait in the system for a long time).
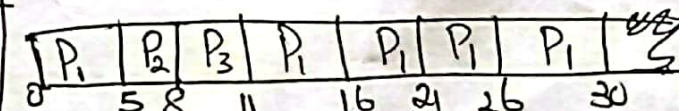
④ **Round-Robin Scheduling**
- Simple & easy to implement
- Assigns time slices to each process in equal portions and in circular order, handl-p all processes without priority.
- Starvation free.
- Similar to FCFS Scheduling, but preemption is added to switch b/w processes.
- A small unit of time called a time quantum / time slice is defined. Generally from 10-100 milliseconds.
- To implement this, we keep the ready queue as a FIFO queue of processes.
- The Average wait-p time is quite long.

Eg. Consider the ffg set of processes, assumed to have arrived at time 0, in the order $P_1, P_2, P_3$, with the length of the CPU burst given in milliseconds.

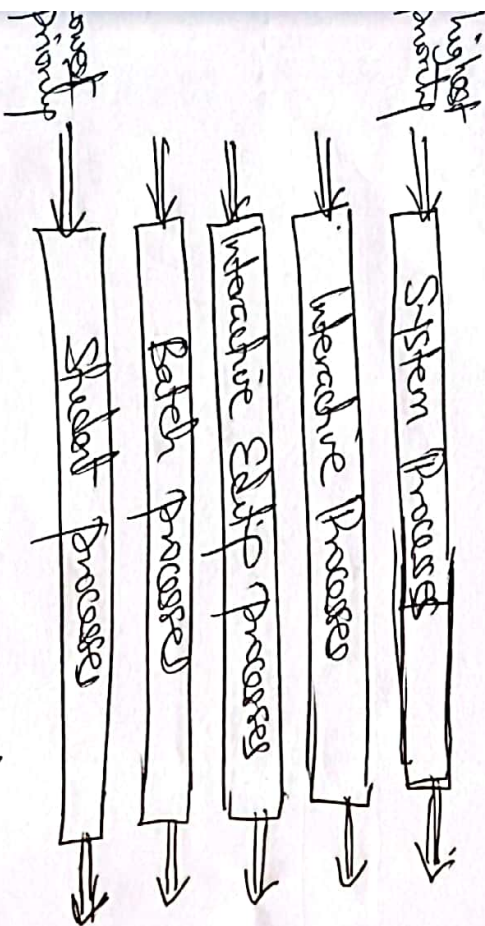| Process | Burst Time |
|---|---|
| $P_1$ | 24 |
| $P_2$ | 3 |
| $P_3$ | 3 |

Time Quantum = 5 ms

| $P_1$ | $P_2$ | $P_3$ | $P_1$ | $P_1$ | $P_1$ | $P_1$ |
|---|---|---|---|---|---|---|

0   5   8   11   16   21   26   30

$$AWT = \left[(11-5) + 5 + 8\right]/3$$
$$= (6+5+8)/3 = 19/3 = 6$$
$$= 6.3 \text{ millisecs}$$

# 5) Multilevel Queue (MLQ) Scheduling

Processes are classified into different groups. The processes are permanently assigned to one queue, generally based on some property of the process. Generally based on some property of the process, process priority etc. Each queue has its own scheduling algorithm. The foreground might use RR while background might use FCFS.

highest priority →

System Process

Interactive Process

Interactive Editing Process

Batch Process

Student Process

lowest priority

## Systematic Approach to Performance Evaluation:

1) State Goals & Define the System
2) List Services and Outcomes
3) Select Metrics
4) List Parameters
5) Select Factors to study
6) Select Evaluation Technique
7) Select Workload
8) Design Experiments
9) Analyze & Interpret Data

# 6) Present Result = Commonly Used Performance Metrics

1) Response time — The total amount of time it takes to respond to a request for service. Response time is the sum of the service time and wait time. Increases as the load on the system increases for FCFS.

2) Throughput → the maximum rate of production. Throughput increases as load. In communication networks, network throughput or the rate of successful message delivery over a communication channel. Throughput is measured in jobs per second for batch systems. Throughput is measured in requests per second for interactive systems. It is measured in Millions of Instructions per second (MIPS) or Millions of Floating Point Operations per second (MFLOPS) or Packets per second for networks. It is measured in bits per second (bps) for transaction processing systems. It is measured in transactions per second (TPS).

System throughput generally increases as the load increases. The maximum achievable throughput under ideal workload conditions is called Nominal Capacity of the system. Nominal Capacity == Bandwidth, expressed in bits per second. Nominal capacity of Network.

③ **Efficiency** → This is the ratio of maximum achievable throughput to nominal capacity.
- the ratio of the performance of an $n$-processor system to the performance of a uni-processor system is the efficiency. usually measured in terms of MIPS or MFLOPS.

④ **Utilization** → if a resource is measured as the fraction of time the resource's busy during a request.

⑤ **Reliability** → if a system is usually measured as the probability of errors or by the mean time b/w errors.
- If the system is not available is called **Downtime**. The time during which the system is available is called **uptime**.
- Mean uptime is also known as the Mean time to failure (MTTF)

⑥ **Availability**: is a ratio of a system's uptime to the total of operational performance period.

⑦ **Transmission time** - gets added to response time when host request over a network and it can be very significant.
- Two types of propagation delays due to latency — transmission errors, and data communication.

---

bandwidth limits

## Amdahl's Law
- is an expression used to find the maximum expected improvement to an overall system when only part of the system is improved. Computer to predict the theoretical maximum speedup using multiple processors.

Amdahl's law - states that the performance improvement to be gained from using some faster mode of execution is limited by the fraction of time the faster mode can be used.

Amdahl's Law is formulated thus:

$$\text{Speedup}(S) = \frac{1}{(1-P) + \frac{P}{S}}$$

where Speedup(S) is the theoretical speedup of the execution of the whole task.
- S is the speedup of the part of the task that benefits from improved system resources
- P is the proportion of execution time that the part benefiting from improved resources

Furthermore, $\text{Speedup}(S) \leq \dfrac{1}{1-P}$

$$\lim_{S \to \infty} \text{Speedup}(S) = \frac{1}{1-P}$$

# Processor Performance

$$\text{CPU time} = \frac{\text{CPU clock cycles for a program} \times}{\text{Clock cycle time}}$$

$$\text{CPU time} = \frac{\text{CPU clock cycles for a program}}{\text{Clock rate}}$$

CPI (cycles per instruction) $= \dfrac{\text{CPU clock cycles for a program}}{\text{Instruction Count}}$

$$\text{CPU time} = \text{Instruction Count} \times \text{Cycles per instruction} \times \text{Clock cycle time}$$

$$\text{CPU Time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{clock cycle}} = \frac{\text{Seconds}}{\text{Program}}$$

$$\text{Performance}(X) = \frac{1}{\text{Execution time}(X)}$$

"X is n times faster than Y" means

$$n = \frac{\text{Performance}(X)}{\text{Performance}(Y)} = \frac{\text{Execution time}(Y)}{\text{Execution time}(X)}$$

Performance for entire task using the enhancement when possible

① Speedup = $\dfrac{\text{Performance for entire task without using the enhancement}}{\text{Performance for entire task using the enhancement when possible}}$

② Speedup = $\dfrac{\text{Execution time b/f improvement}}{\text{Execution time after improvement}}$

---

* If a CPU of A runs a pgm in 10 sec and of B runs the same pgm in 15 seconds, is CPU A faster? and How much faster?

Ans  A is n times faster than B $\Rightarrow$

$$\frac{\text{Performance of A}}{\text{Performance of B}} = \frac{\text{Execution time of B}}{\text{Execution time of A}} = n$$

The performance ratio n is $\dfrac{15}{10} = 1.5$

So A is 1.5 times (50%) faster than B.