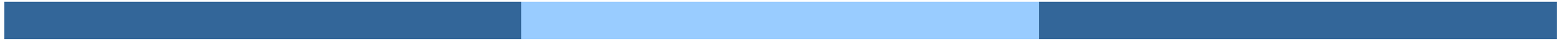


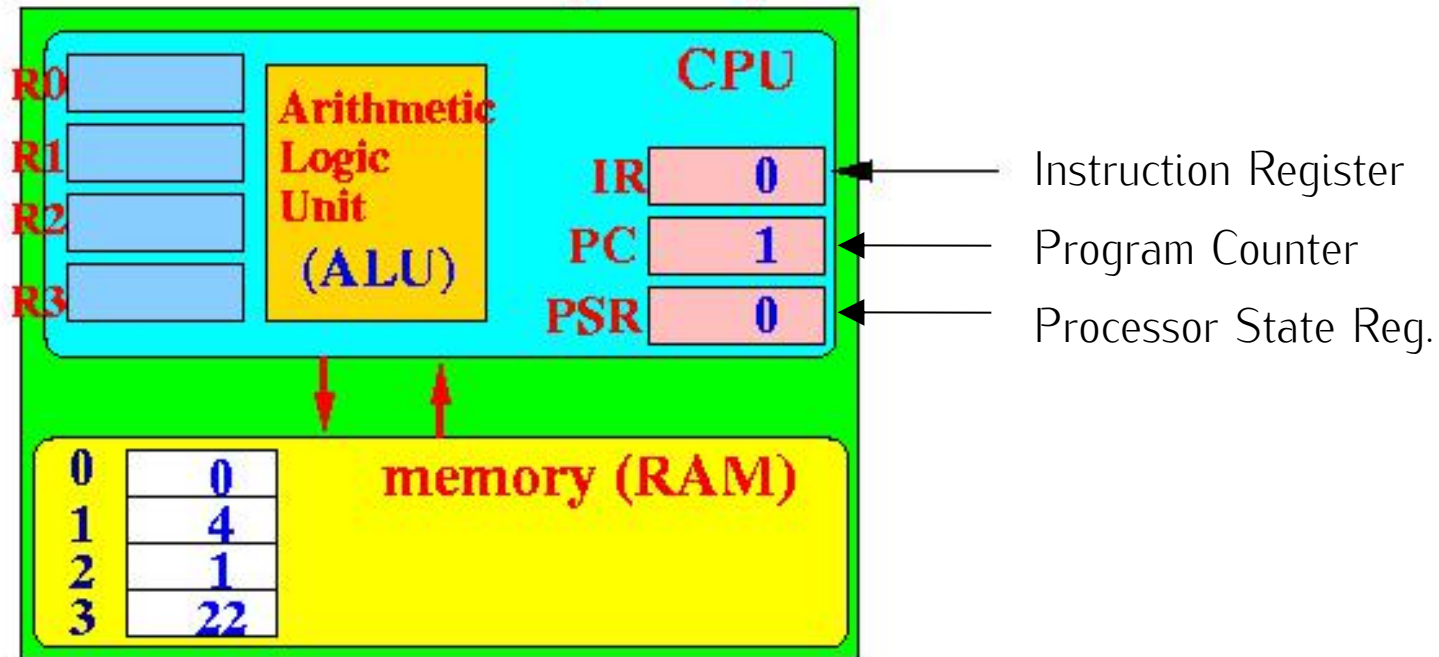
CPU Scheduling



Review



Architecture of CPU



Process in Memory

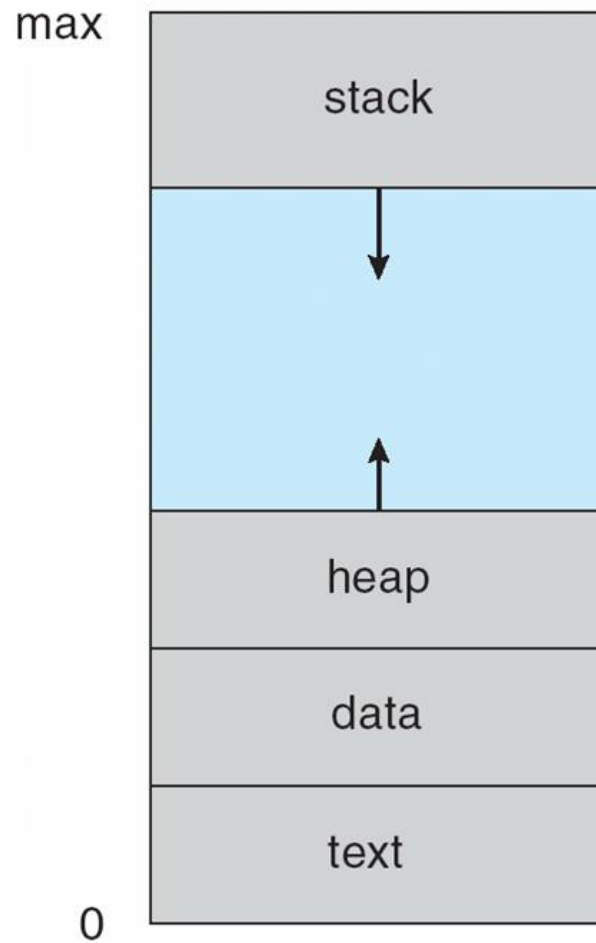
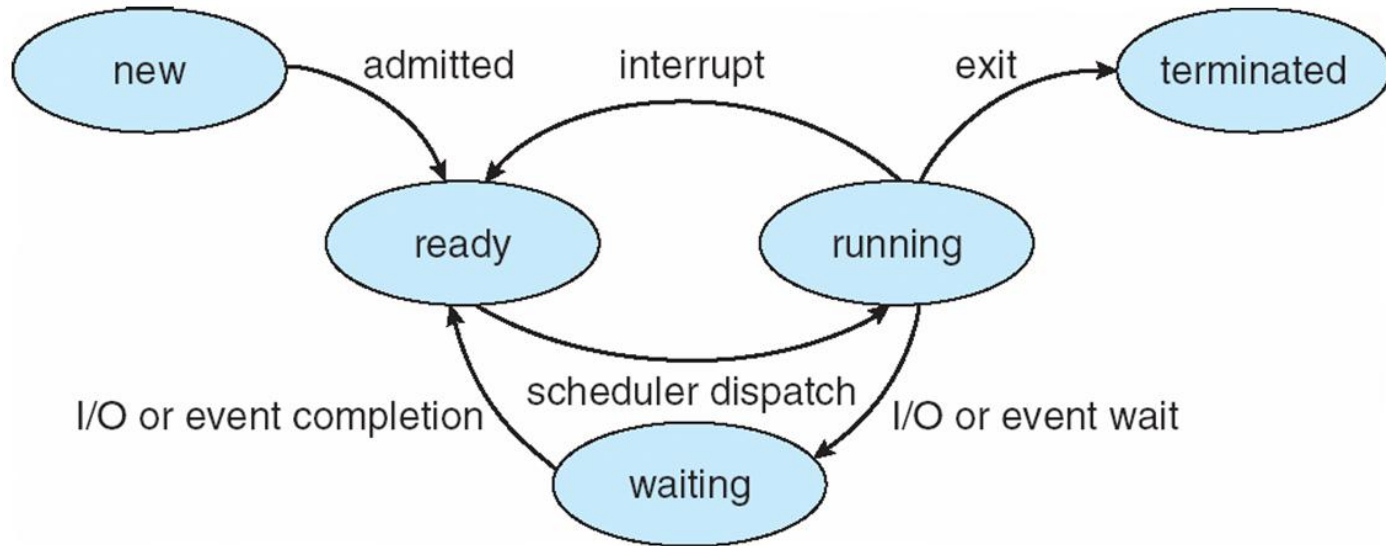


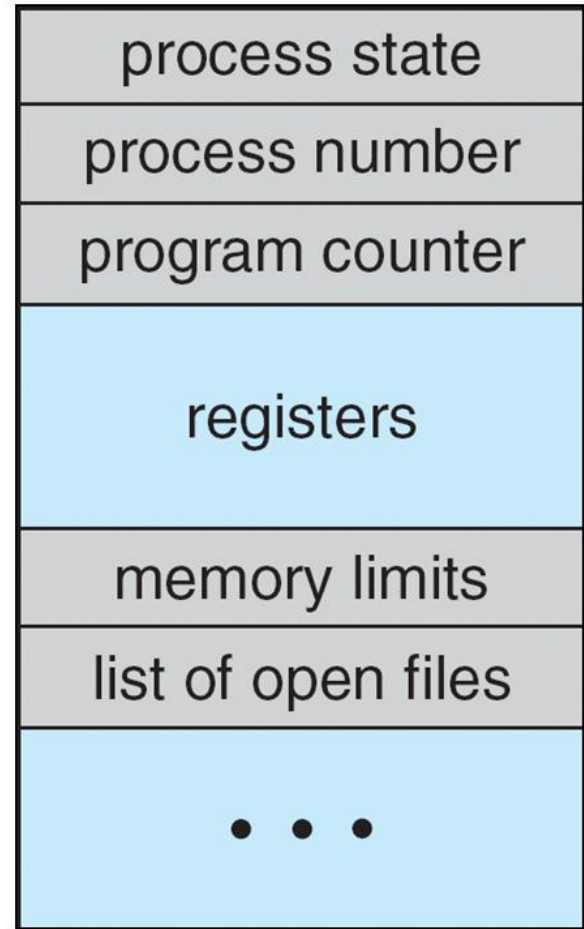
Diagram of Process State



Process Control Block (PCB)

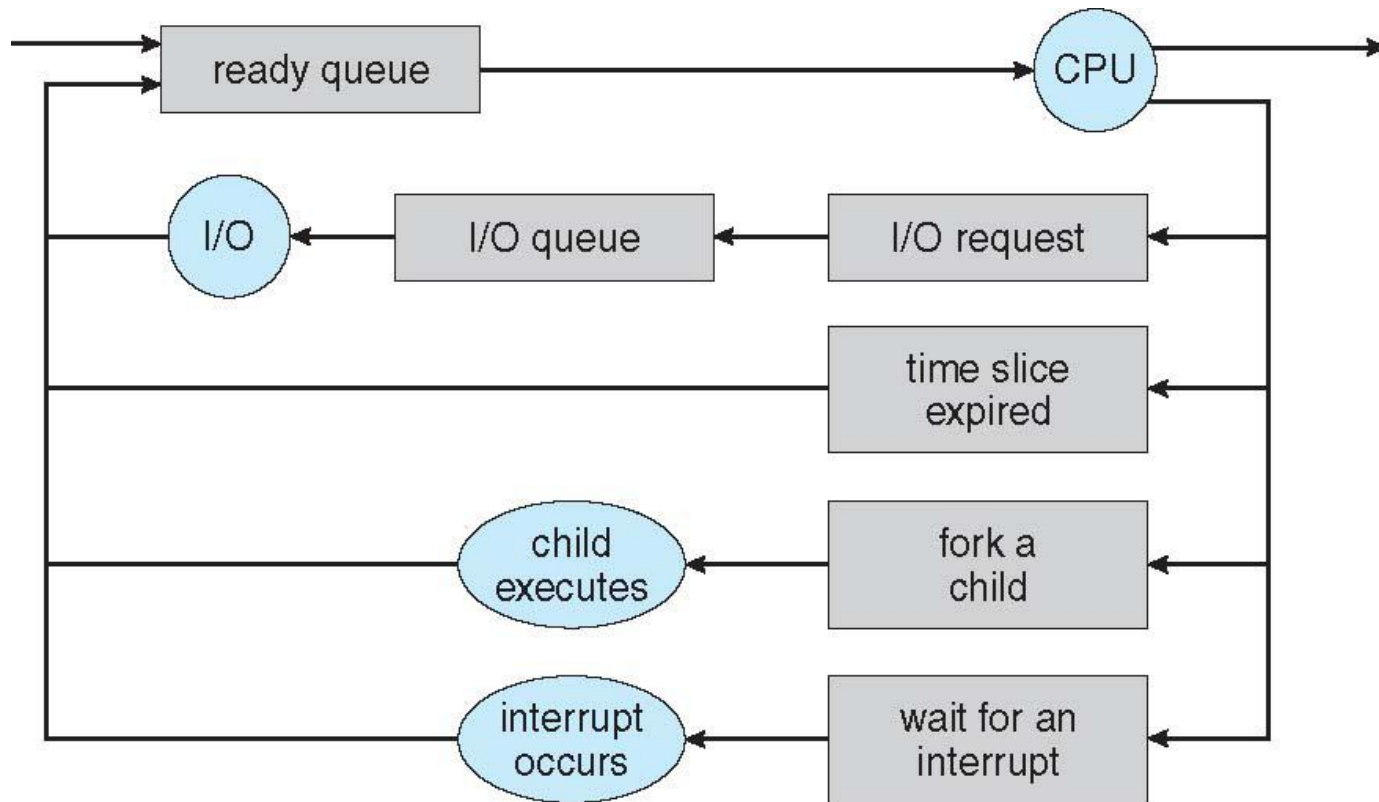
Information associated with each process
(also called **task control block**)

- Process state – running, waiting, etc
- Program counter – location of instruction to next execute
- CPU registers – contents of all process-centric registers
- CPU scheduling information– priorities, scheduling queue pointers
- Memory-management information – memory allocated to the process
- Accounting information – CPU used, clock time elapsed since start, time limits
- I/O status information – I/O devices allocated to process, list of open files



Representation of Process Scheduling

- n Queueing diagram represents queues, resources, flows

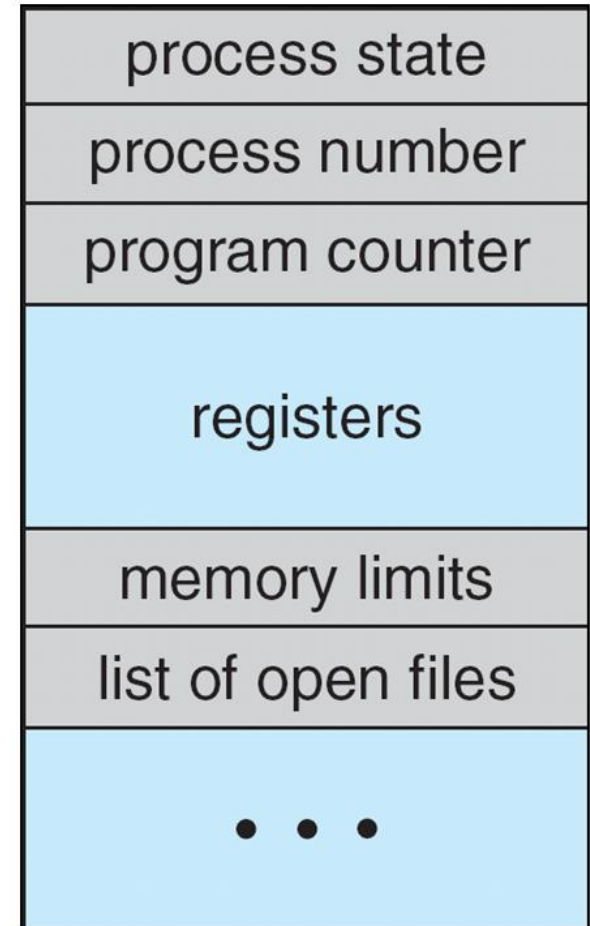


Schedulers

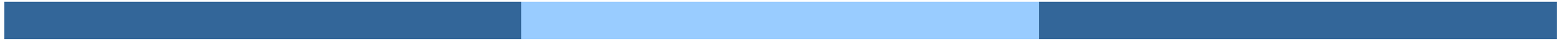
- **Short-term scheduler** (or **CPU scheduler**) – selects which process should be executed next and allocates CPU
 - l Sometimes the only scheduler in a system
 - l Short-term scheduler is invoked frequently (milliseconds) c (must be fast)
- **Long-term scheduler** (or **job scheduler**) – selects which processes should be brought into the ready queue
 - l Long-term scheduler is invoked infrequently (seconds, minutes) c (may be slow)
 - l The long-term scheduler controls the **degree of multiprogramming**
- Processes can be described as either:
 - l **I/O-bound process** – spends more time doing I/O than computations, many short CPU bursts
 - l **CPU-bound process** – spends more time doing computations; few very long CPU bursts
- Long-term scheduler strives for good *process mix*

Context Switch

- When CPU switches to another process, the system must **save the state** of the old process and load the **saved state** for the new process via a **context switch**
- **Context** of a process represented in the PCB
- Context-switch time is overhead; the system does no useful work while switching
 - └ The more complex the OS and the PCB the longer the context switch
- Time dependent on hardware support
 - └ Some hardware provides multiple sets of registers per CPU multiple contexts loaded at once

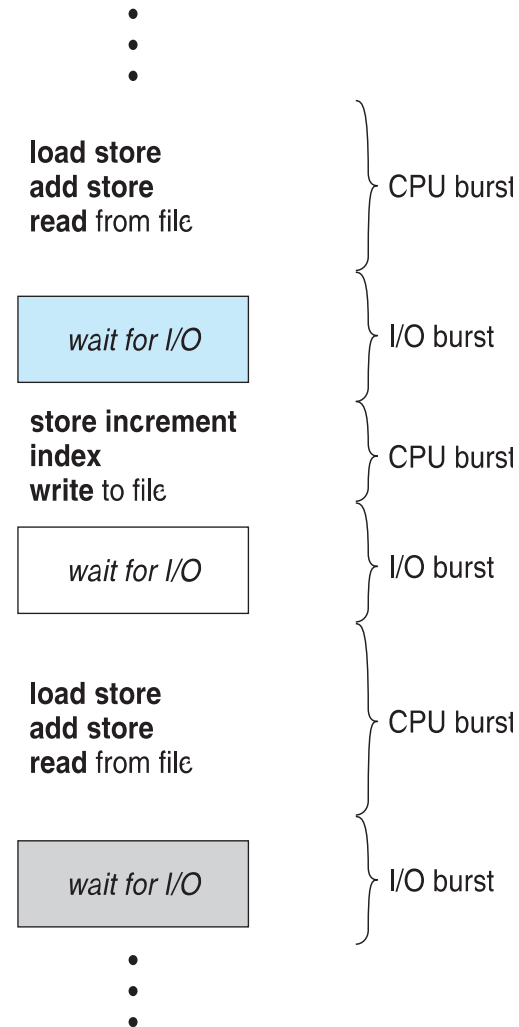


Scheduling



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



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 is the cause of following state changes:
 1. Switches from running to waiting state
 2. Switches from running to ready state
 3. Switches from waiting to ready
 4. Terminates
- For 1 and 4 **non-preemptive** scheduler is used
- 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:
 - l switching context
 - l switching to user mode
 - l 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

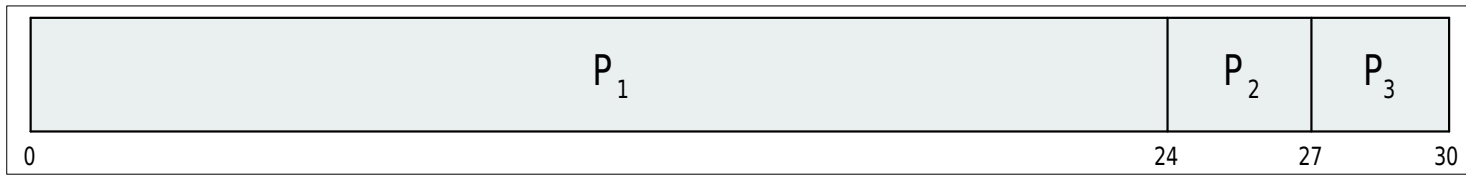
- **CPU utilization** – keep the CPU as busy as possible
- **Throughput** – # of processes that complete their execution per time unit
- **Turnaround time** – amount of time to execute a particular process
- **Waiting time** – amount of time a process has been waiting in the ready queue
- **Response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

First- Come, First-Served (FCFS) Scheduling

| Process | Burst Time |
|---------|------------|
| 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:



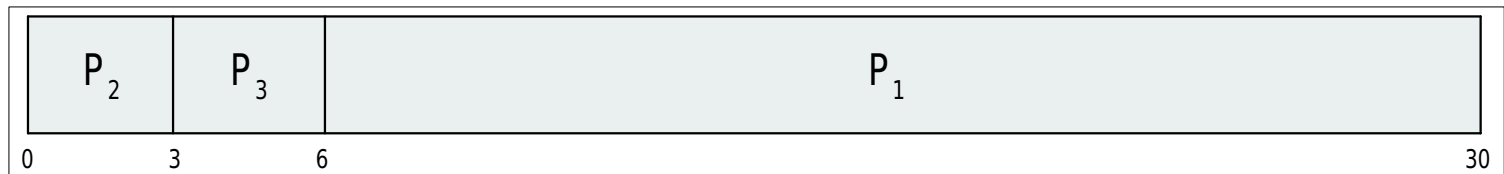
- Waiting time for $P_1 = 0$; $P_2 = 24$; $P_3 = 27$
- Average waiting time: $(0 + 24 + 27)/3 = 17$

FCFS Scheduling (Cont.)

Suppose that the processes arrive in the order:

$$P_2, P_3, P_1$$

- The Gantt chart for the schedule is:



- Waiting time for $P_1 = 6$; $P_2 = 0$; $P_3 = 3$
- Average waiting time: $(6 + 0 + 3)/3 = 3$
- Much better than previous case
- **Convoy effect** – short process behind long process
 - └ Consider one CPU-bound and many I/O-bound processes

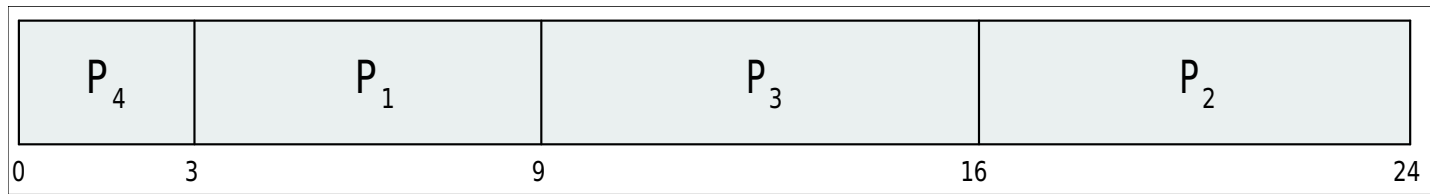
Shortest-Job-First (SJF) Scheduling

- Associate with each process the length of its next CPU burst
 - | Use these lengths to schedule the process with the shortest time
- SJF is optimal – gives minimum average waiting time for a given set of processes
 - | The difficulty is knowing the length of the next CPU request
 - | Could ask the user

Example of SJF

| Process | Burst Time |
|---------|------------|
| P_1 | 6 |
| P_2 | 8 |
| P_3 | 7 |
| P_4 | 3 |

- SJF scheduling chart

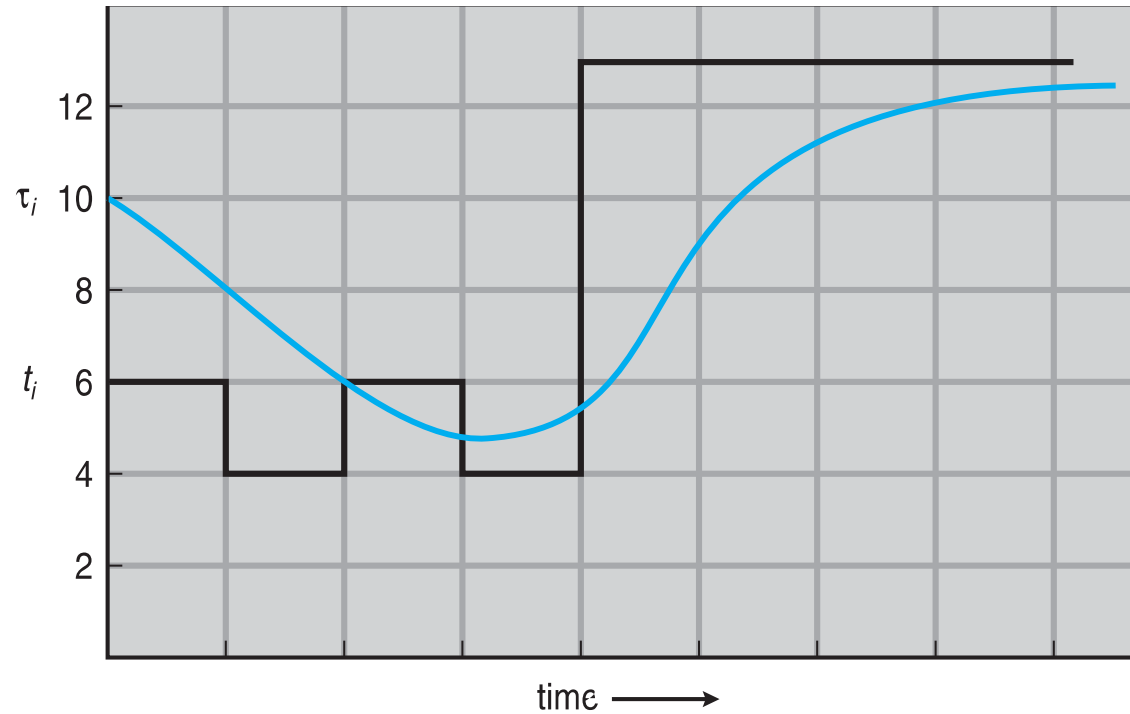


- Average waiting time = $(3 + 16 + 9 + 0) / 4 = 7$

Length of Next CPU Burst

- Can only estimate the length – should be similar to the previous one
 - | Then pick process with shortest predicted next CPU burst
- Can be done by using the length of previous CPU bursts, using exponential averaging
 - | t_n : actual length of n^{th} CPU burst time
 - | T_n : predicted value for the next CPU burst
 - | A : is between 0 and 1
 - | $T_{n+1} = A t_n + (1-A) T_n$
- Commonly, A set to $\frac{1}{2}$
- Preemptive version called **shortest-remaining-time-first**

Prediction of the Next CPU Burst



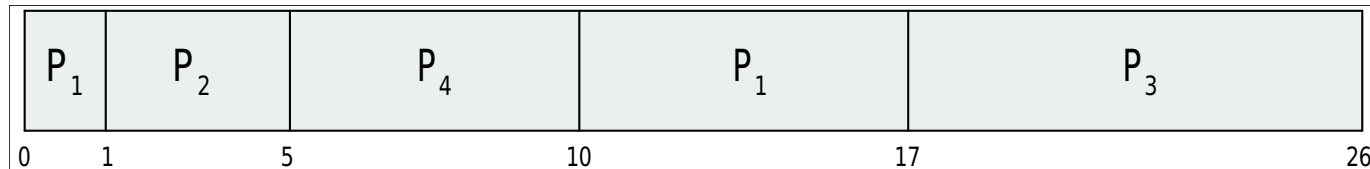
| | | | | | | | | |
|----------------------|----|---|---|---|----|----|----|-----|
| CPU burst (t_i) | 6 | 4 | 6 | 4 | 13 | 13 | 13 | ... |
| "guess" (τ_i) | 10 | 8 | 6 | 6 | 9 | 11 | 12 | ... |

Shortest-remaining-time-first

- Now we add the concepts of varying arrival times and preemption to the analysis

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

- Preemptive* SJF Gantt Chart



- Average waiting time = $[(10-1)+(1-1)+(17-2)+5-3)]/4 = 26/4 = 6.5$ msec

Priority Scheduling

- A priority number (integer) is associated with each process
- The CPU is allocated to the process with the highest priority (smallest integer highest priority)
 - l Preemptive
 - l Nonpreemptive
- SJF is priority scheduling where priority is the inverse of predicted next CPU burst time
- Problem **Starvation** – low priority processes may never execute
- Solution **Aging** – as time progresses increase the priority of the process

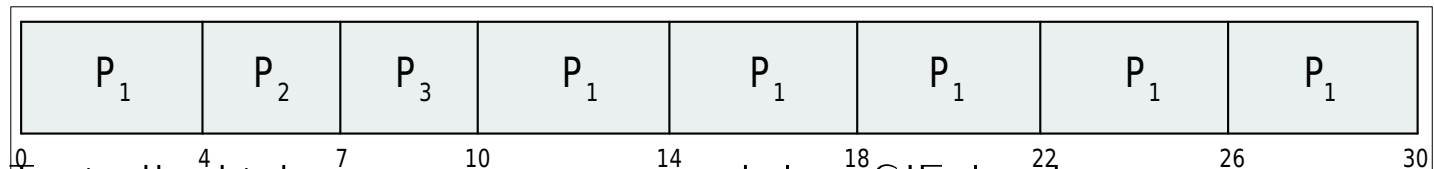
Round Robin (RR)

- 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.
- If there are n processes in the ready queue and the time quantum is q , then each process gets $1/n$ of the CPU time in chunks of at most q time units at once. No process waits more than $(n-1)q$ time units.
- Timer interrupts every quantum to schedule next process
- Performance
 - l q large : FIFO
 - l q small : q must be large with respect to context switch, otherwise overhead is too high

RR with Time Quantum = 4

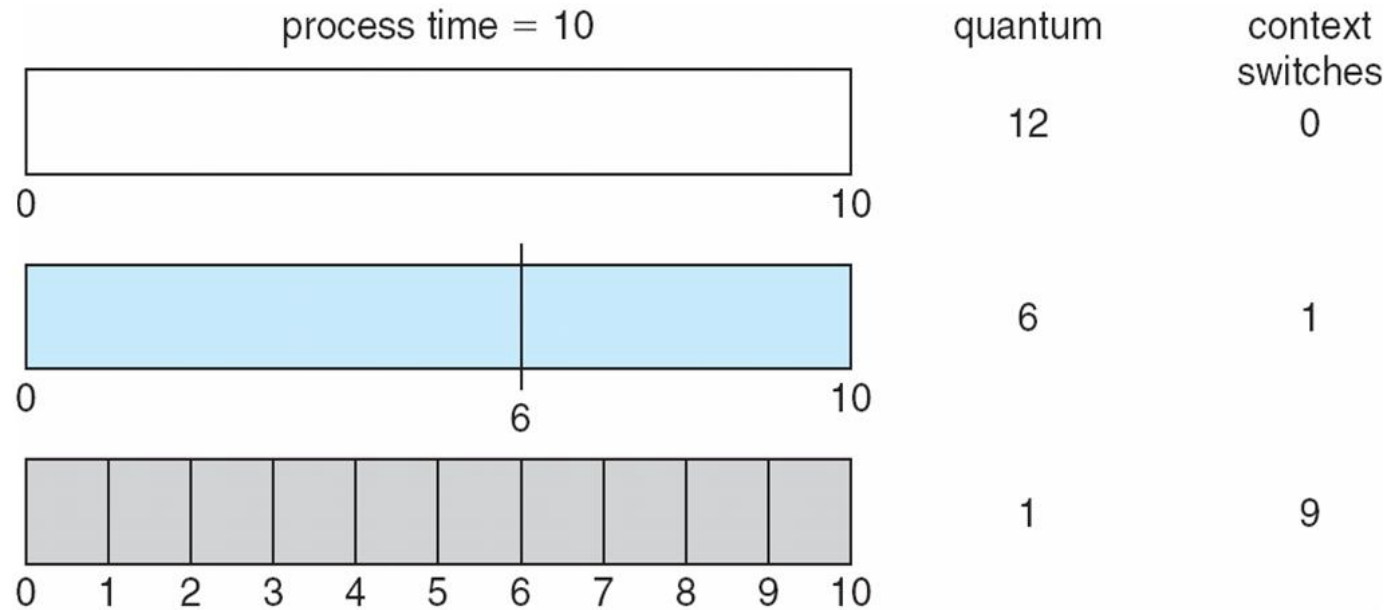
| Process | Burst Time |
|---------|------------|
| P_1 | 24 |
| P_2 | 3 |
| P_3 | 3 |

- The Gantt chart is:



- Typically, higher average turnaround than SJF, but better *response*
- q should be large compared to context switch time
- q usually 10ms to 100ms, context switch < 10 usec

Time Quantum & Context Switch Time

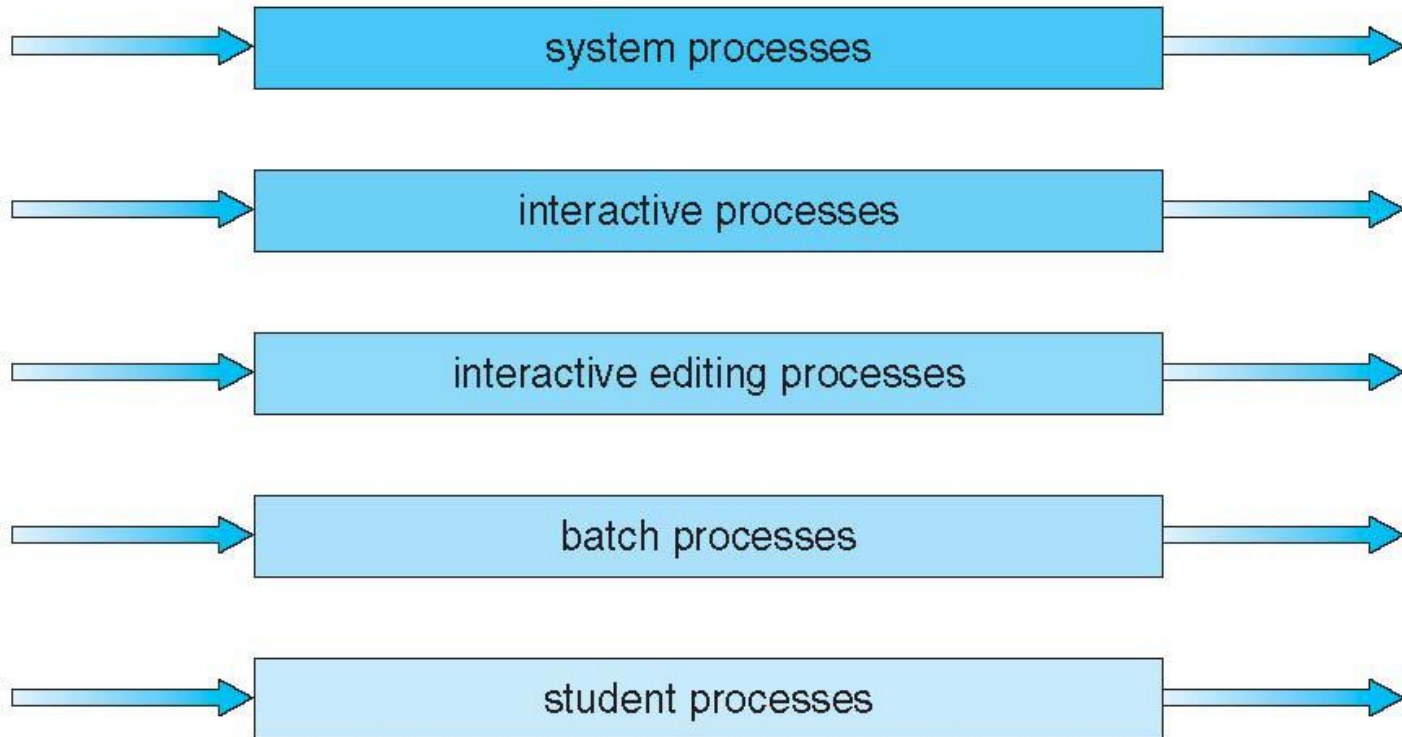


Multilevel Queue

- Ready queue is partitioned into separate queues, eg:
 - l **foreground** (interactive)
 - l **background** (batch)
- Process permanently in a given queue
- Each queue has its own scheduling algorithm:
 - l foreground – RR
 - l background – FCFS
- Scheduling must be done between the queues:
 - l Fixed priority scheduling; (i.e., serve all from foreground then from background). Possibility of starvation.
 - l Time slice – each queue gets a certain amount of CPU time which it can schedule amongst its processes; i.e., 80% to foreground in RR
 - l 20% to background in FCFS

Multilevel Queue Scheduling

highest priority



lowest priority

Multilevel Feedback Queue

- A process can move between the various queues; aging can be implemented this way
- Multilevel-feedback-queue scheduler defined by the following parameters:
 - l number of queues
 - l scheduling algorithms for each queue
 - l method used to determine when to upgrade a process
 - l method used to determine when to demote a process
 - l method used to determine which queue a process will enter when that process needs service

Multilevel Feedback Queue

- Three queues:
 - l Q_0 – time quantum 8 milliseconds
 - l Q_1 – time quantum 16 milliseconds
 - l Q_2 – FCFS
- Scheduling
 - l A new job enters queue Q_0 which is served FCFS
 - 4 When it gains CPU, job receives 8 milliseconds
 - 4 If it does not finish in 8 milliseconds, job is moved to queue Q_1
 - l At Q_1 job is again served FCFS and receives 16 additional milliseconds
 - 4 If it still does not complete, it is preempted and moved to queue Q_2

