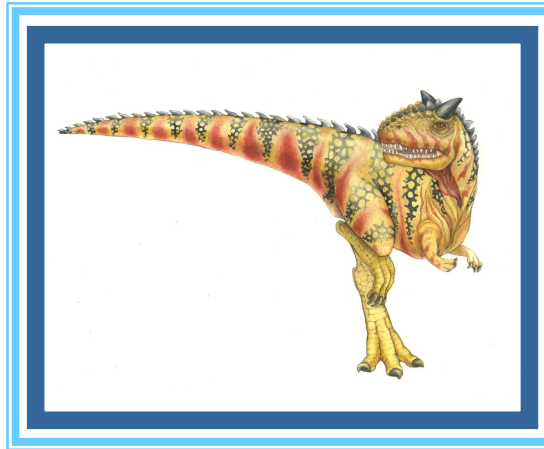


# Chapter 6: CPU Scheduling

---





# Chapter 6: CPU Scheduling

---

- Basic Concepts
- Scheduling Criteria
- Scheduling Algorithms
- Thread Scheduling





# Objectives

---

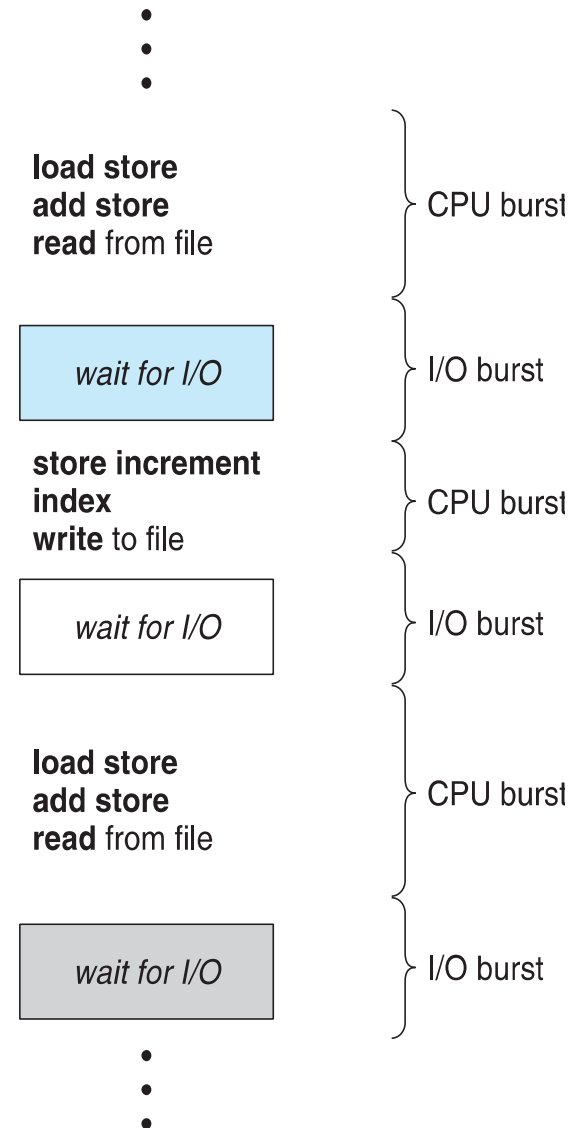
- To introduce CPU scheduling, which is the basis for multiprogrammed operating systems
- To describe various CPU-scheduling algorithms
- To discuss evaluation criteria for selecting a CPU-scheduling algorithm for a particular system

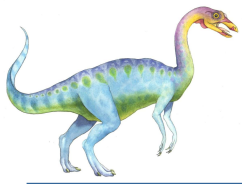




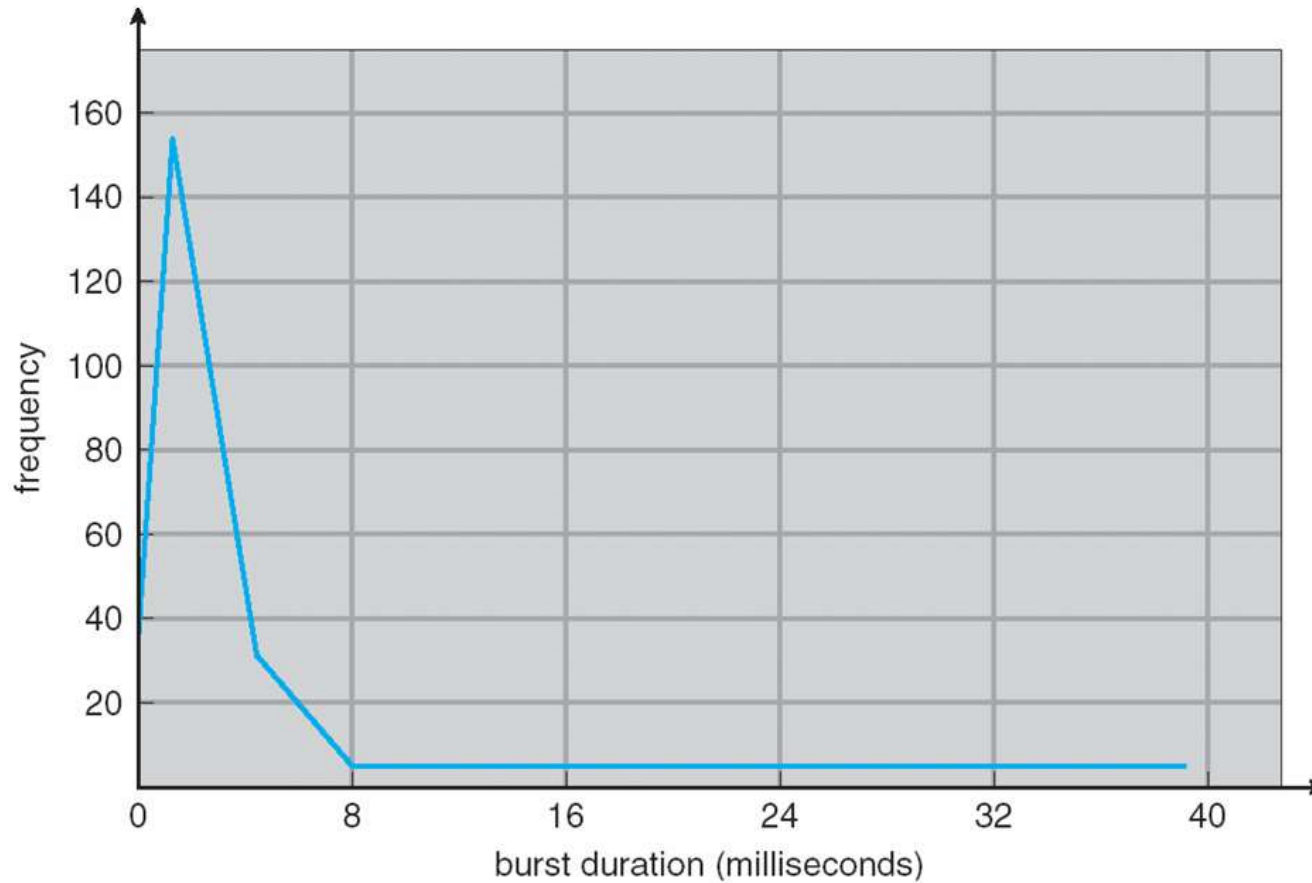
# Basic Concepts

- ❑ 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
- ❑ The following figure shows that we usually have a large number of short CPU bursts and a small number of long CPU bursts.





# 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
  - ❑ A ready queue can be implemented as a FIFO queue, a priority queue, a tree, or simply an unordered linked list
  - ❑ The records in the queues are generally process control blocks (PCBs) of the processes
- ❑ CPU scheduling decisions may take place when a process:
  1. Switches from running to waiting state (no choice in terms of scheduling)
  2. Switches from running to ready state (there is a choice)
  3. Switches from waiting to ready (there is a choice)
  4. Terminates (no choice in terms of scheduling)





- ❑ Scheduling under 1 and 4 is **nonpreemptive (or cooperative)**
- ❑ All other scheduling is **preemptive**
- ❑ **Nonpreemptive scheduling: Once the 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 waiting state.**
- ❑ **Unfortunately, preemptive scheduling can result in race conditions when data are shared among several processes.**
- ❑ **Preemptive Scheduling needs to**
  - ❑ Consider access to shared data (the race condition)
  - ❑ Consider preemption while in kernel mode (might also occur while accessing shared kernel data)
  - ❑ Consider interrupts occurring during crucial OS activities (the sections of code affected by interrupts must be guarded from simultaneous use; otherwise input might be lost or output overwritten; these sections of code are not accessed concurrently by several processes, they disable interrupts at entry and re-enable interrupts at exit.)





# 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
- ❑ The dispatcher should be as fast as possible, since it is invoked during every process switch







# 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 (The interval from the time of submission of a process to the time of completion is the turnaround time.)
- ❑ **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, This is the time it takes to start responding, not the time it takes to output the response (for time-sharing environment). In an interactive system this is more important than the turn around time which is limited by the speed of the output device.





# Scheduling Algorithm Optimization Criteria

---

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





# 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:  $(0 + 24 + 27)/3 = 17$



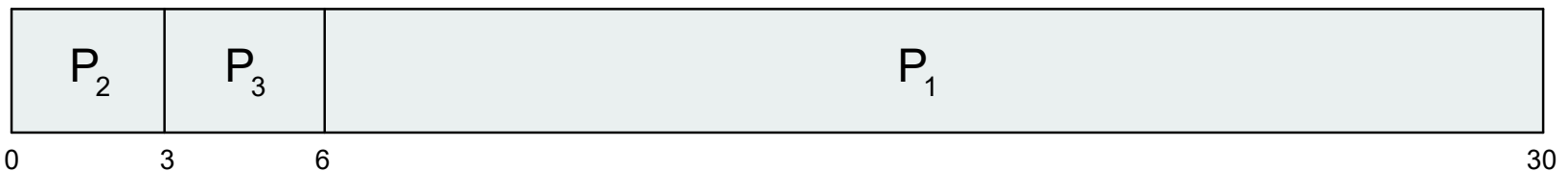


# FCFS Scheduling (Cont.)

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:  $(6 + 0 + 3)/3 = 3$  (a substantial reduction)
- Much better than previous case
- Thus, the average waiting time under an FCFS policy is generally not minimal and may vary substantially if the processes' CPU burst times vary greatly
- **Convoy effect** - short process behind long process
  - Consider one CPU-bound and many I/O-bound processes
  - This effect results in lower CPU and device utilization than might be possible if the shorter processes were allowed to go first. (this is bad of course)





# 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

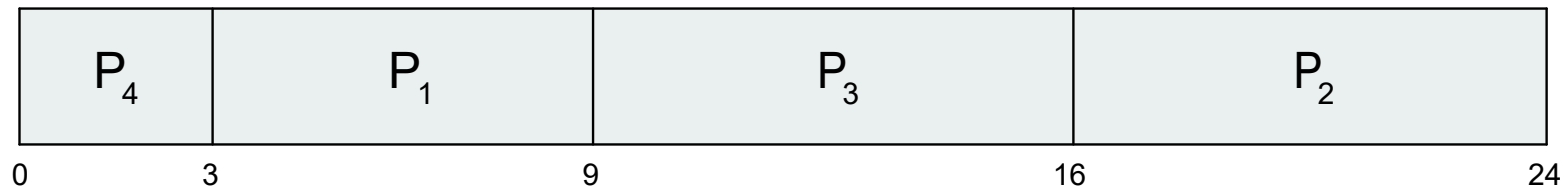




# Example of SJF

<u>Process</u>	<u>Burst Time</u>
$P_1$	6
$P_2$	8
$P_3$	7
$P_4$	3

## □ SJF scheduling chart



- Average waiting time =  $(3 + 16 + 9 + 0) / 4 = 7$
- By comparison, if we were using the FCFS scheduling scheme, the average waiting time would be 10.25 milliseconds.
- Notice that the FCFS scheduling algorithm is non-pre-emptive





- The SJF scheduling algorithm **is provably optimal**, in that it gives the minimum average waiting time for a given set of processes.
- Moving a short process before along one decreases the waiting time of the short process more than it increases the waiting time of the long process.
- Consequently, the average waiting time decreases.
  
- Although the SJF algorithm is optimal, it cannot be implemented at the level of short-term CPU scheduling. With short-term scheduling, **there is no way to know the length of the next CPU burst.**





# Determining Length of Next CPU Burst

- With short-term scheduling, there is no way to know the length of the next CPU burst
- Can only estimate the length – should be similar to the previous ones
  - Then pick process with shortest predicted next CPU burst
- Can be done by using the length of previous CPU bursts, **using exponential averaging**

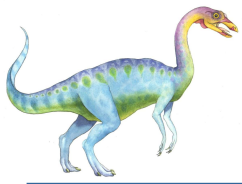
1.  $t_n$  = actual length of  $n^{th}$  CPU burst
2.  $\tau_{n+1}$  = predicted value for the next CPU burst
3.  $\alpha, 0 \leq \alpha \leq 1$
4. Define:

$$\tau_{n+1} = \alpha t_n + (1 - \alpha)\tau_n.$$

- Commonly,  $\alpha$  set to  $\frac{1}{2}$  **(to give recent and past history equal weights)**
- Preemptive version called **shortest-remaining-time-first**







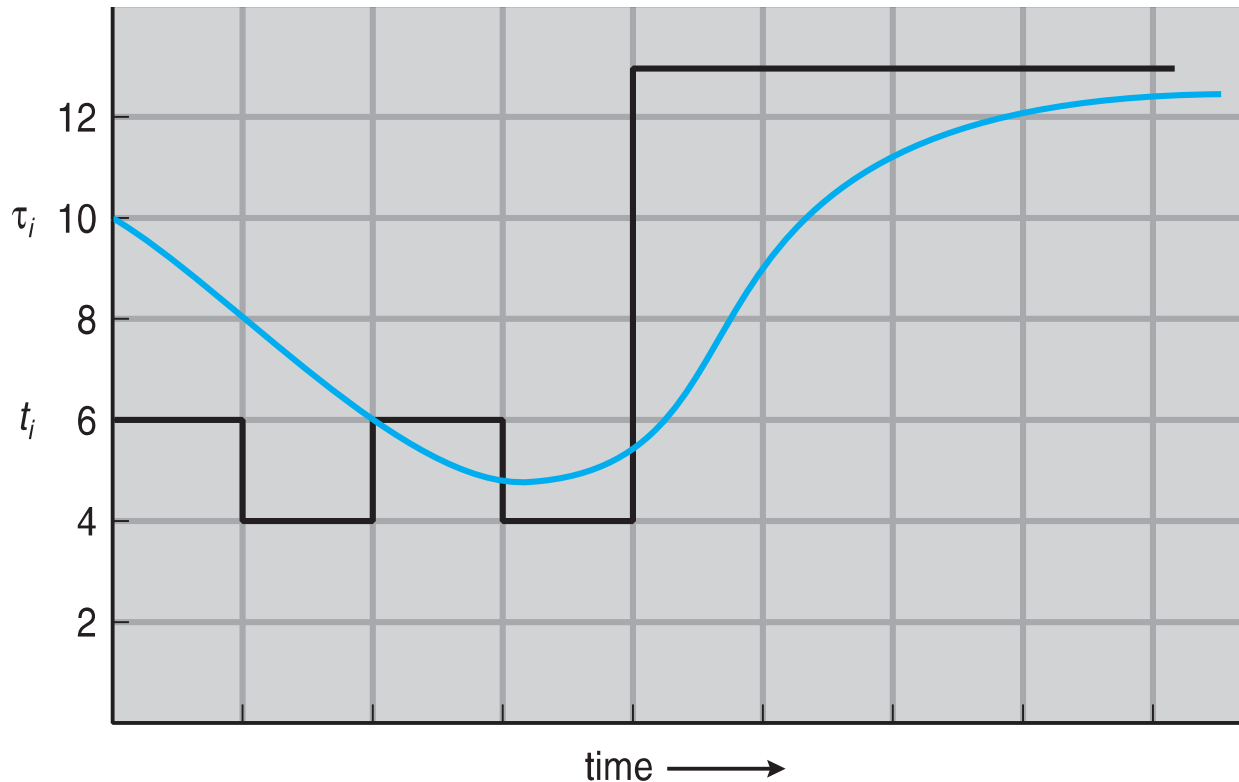
- Assuming  $\alpha$  is  $\frac{1}{2}$  and  $T_0=10$





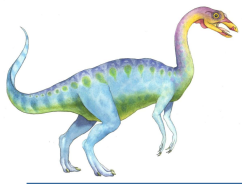
# Prediction of the Length of the Next CPU Burst

- Assuming  $\alpha$  is  $\frac{1}{2}$  and  $T_0=10$



CPU burst ( $t_i$ )	6	4	6	4	13	13	13	...	
"guess" ( $\tau_i$ )	10	8	6	6	5	9	11	12	...





# Examples of Exponential Averaging

- $\alpha = 0$ 
  - $\tau_{n+1} = \tau_n$
  - Recent history does not count
- $\alpha = 1$ 
  - $\tau_{n+1} = \alpha t_n$
  - Only the actual last CPU burst counts
- If we expand the formula by substituting for  $T_n$ , we get:

$$\begin{aligned}\tau_{n+1} = & \alpha t_n + (1 - \alpha)\alpha t_{n-1} + \dots \\ & + (1 - \alpha)^j \alpha t_{n-j} + \dots \\ & + (1 - \alpha)^{n+1} \tau_0\end{aligned}$$

- Since both  $\alpha$  and  $(1 - \alpha)$  are less than or equal to 1, each successive term has less weight than its predecessor



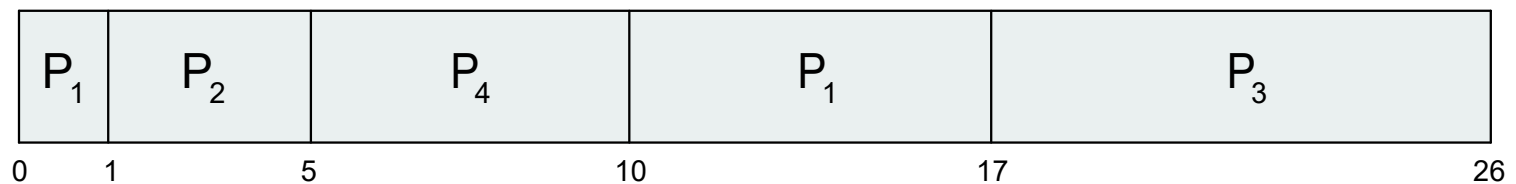


# Example of Shortest-remaining-time-first

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

<u>Process</u>	<u>Arrival Time</u>	<u>Burst Time</u>
$P_1$	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





- The next CPU burst of a newly arrived process may be shorter than what is left of the currently executing process.
- A preemptive SJF algorithm will preempt the currently executing process,
- whereas a nonpreemptive SJF algorithm will allow the currently running process to finish its CPU burst.
- Preemptive SJF scheduling is sometimes called shortest-remaining-time-first scheduling.





# Priority Scheduling

---

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





# Example of Priority Scheduling

<u>Process</u>	<u>Burst Time</u>	<u>Priority</u>
$P_1$	10	3
$P_2$	1	1
$P_3$	2	4
$P_4$	1	5
$P_5$	5	2

□ Priority scheduling Gantt Chart



□ Average waiting time = 8.2 msec





# 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 (**because the process may release the CPU before its time slice is over**).
- ❑ No process waits more than  $(n-1)q$  time units **until its next time quantum**
- ❑ For example, with five processes and a time quantum of 20 milliseconds, each process will get up to 20 milliseconds every 100 milliseconds







- The performance of the RR algorithm depends heavily on the size of the time quantum.
  - At one extreme, if the time quantum is **extremely large**, the RR policy **s the same as the FCFS policy**.
  - In contrast, if the time quantum is **extremely small** (say, 1 millisecond), the RR approach can result in **a large number of context switches**.
- 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**
  - If the context-switch time is approximately 10 percent of the time quantum, then about 10 percent of the CPU time will be spent in context switching.

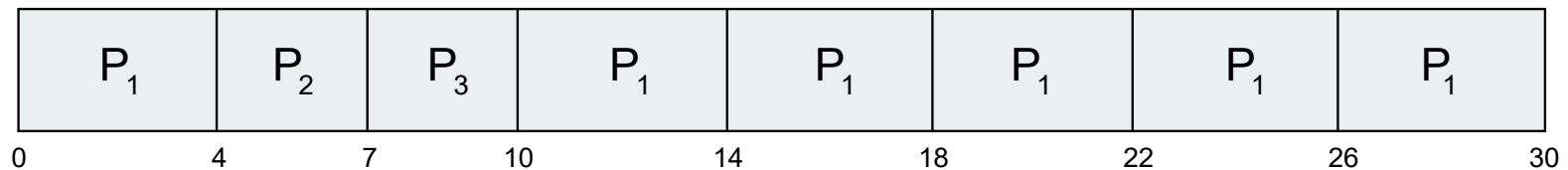




# Example of RR with Time Quantum = 4

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

□ The Gantt chart is:



□ **the average waiting time**

- $P_1$  waits for 6 milliseconds(10-4),
  - $P_2$  waits for 4 milliseconds,
  - and  $P_3$  waits for 7 milliseconds.
- Thus, the average waiting time is  $17/3 = 5.66$  milliseconds.



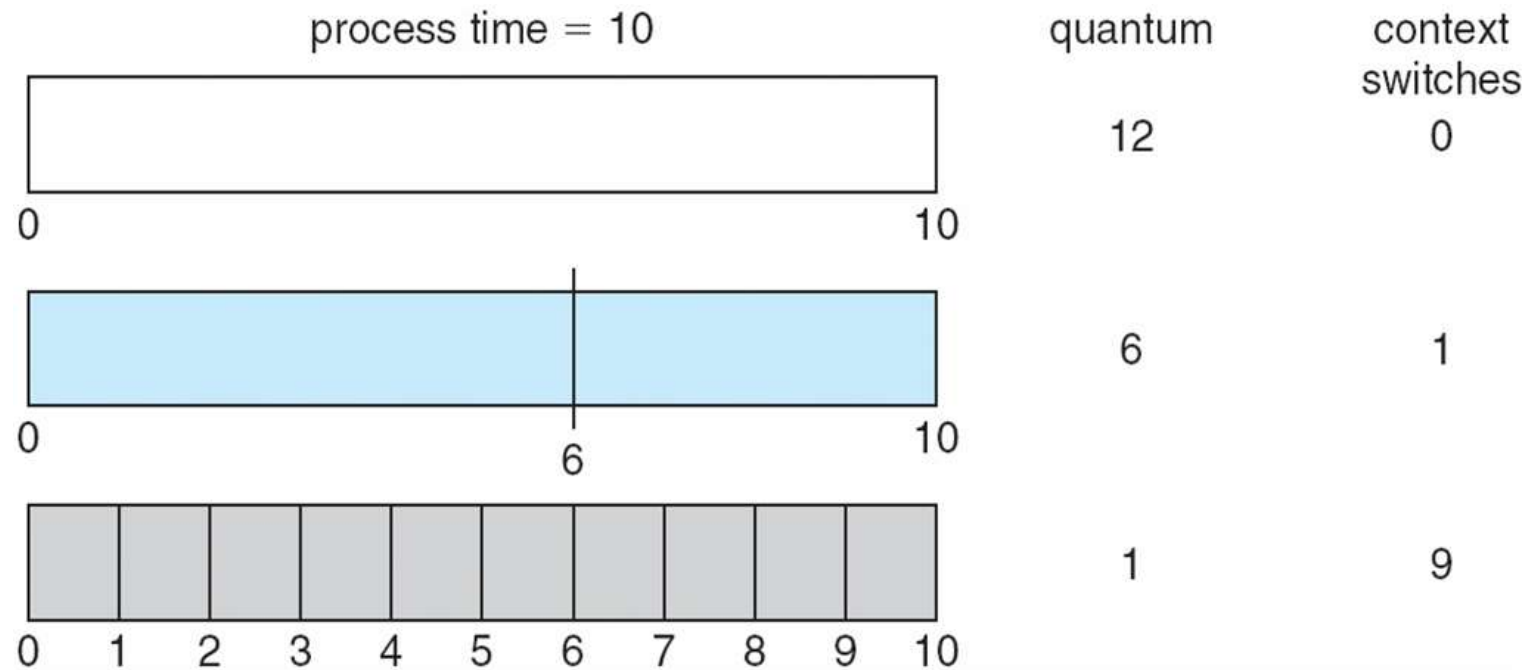


- Typically, **higher average turnaround** than SJF, but better ***response***
- $q$  should be large compared to context switch time
- In modern machines,  $q$  is usually 10ms to 100ms, context switch  $< 10$  ms



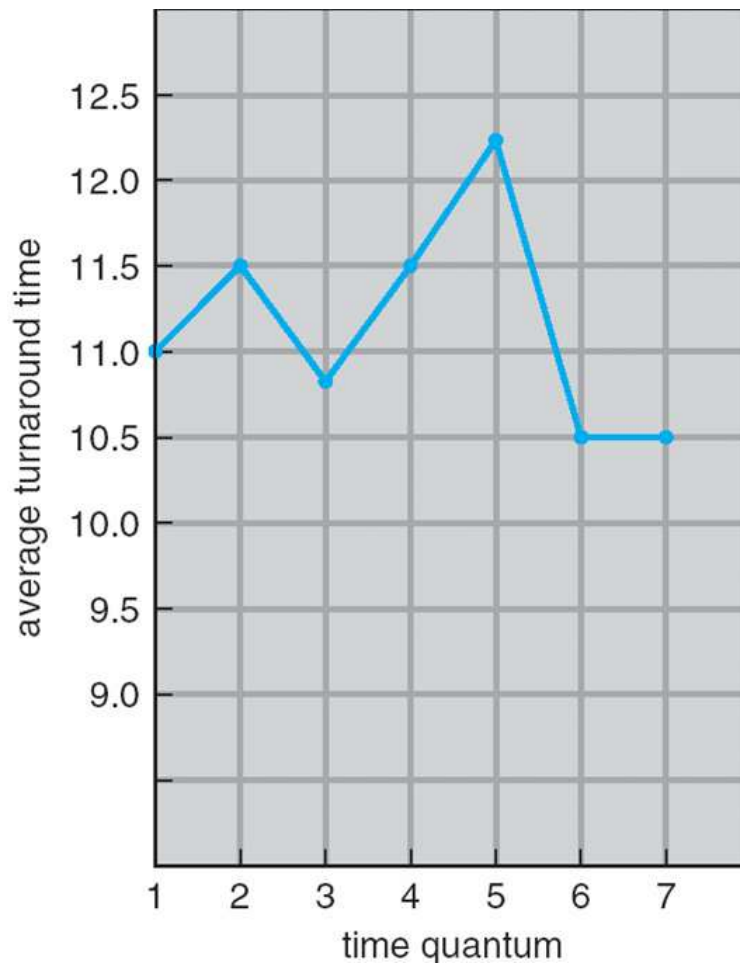


# Time Quantum and Context Switch Time





# Turnaround Time Varies With The Time Quantum



process	time
$P_1$	6
$P_2$	3
$P_3$	1
$P_4$	7

80% of CPU bursts  
should be shorter than  $q$





# Multilevel Queue

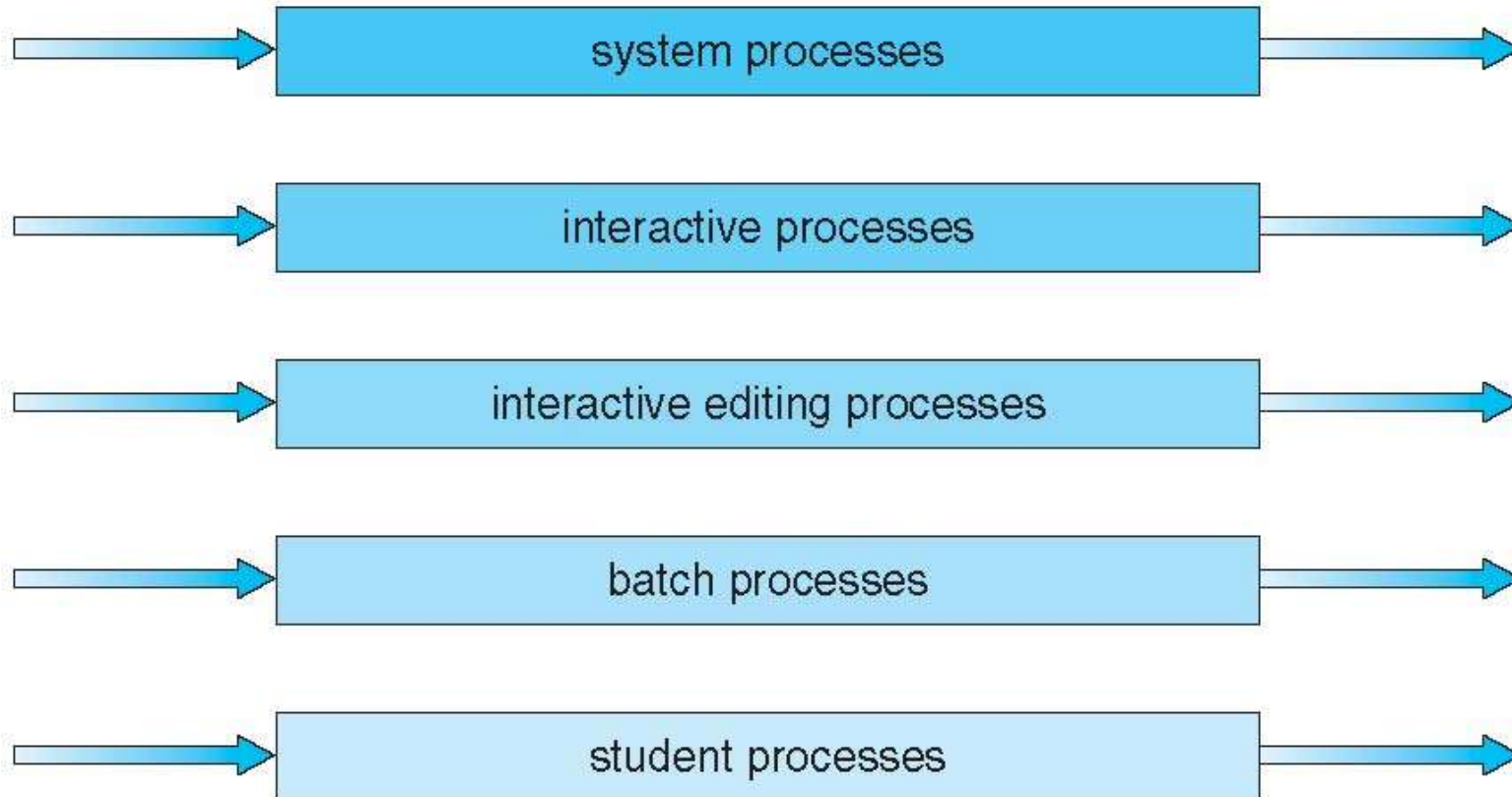
- Ready queue is partitioned into separate queues, eg:
  - **foreground** (interactive)
  - **background** (batch)
- Process permanently in a given queue
- Each queue has its own scheduling algorithm:
  - foreground – RR
  - background – FCFS
- Scheduling must be done between the queues:
  - Fixed priority scheduling; (i.e., serve all from foreground then from background). Possibility of starvation.
  - 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
  - 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:
  - number of queues
  - scheduling algorithms for each queue
  - method used to determine when to upgrade a process
  - method used to determine when to demote a process
  - method used to determine which queue a process will enter when that process needs service





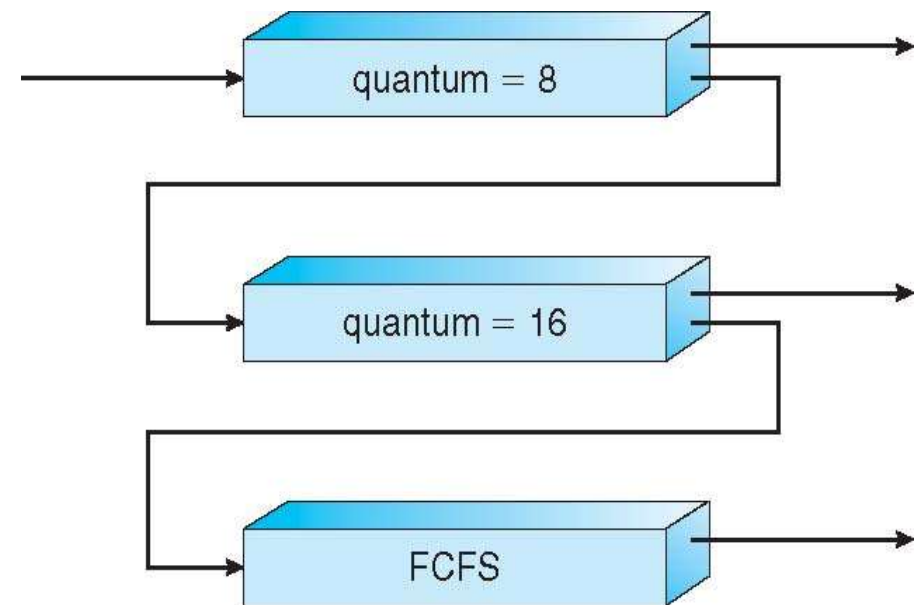


# Example of Multilevel Feedback Queue

- Three queues:
  - $Q_0$  – RR with time quantum 8 milliseconds
  - $Q_1$  – RR time quantum 16 milliseconds
  - $Q_2$  – FCFS

- Scheduling

- A new job enters queue  $Q_0$  which is served FCFS
  - ▶ When it gains CPU, job receives 8 milliseconds
  - ▶ If it does not finish in 8 milliseconds, job is moved to queue  $Q_1$
- At  $Q_1$  job is again served FCFS and receives 16 additional milliseconds
  - ▶ If it still does not complete, it is preempted and moved to queue  $Q_2$





# Thread Scheduling

---

- Distinction between user-level and kernel-level threads
- When threads supported, threads scheduled, not processes
- Many-to-one and many-to-many models, thread library schedules user-level threads to run on LWP
  - Known as **process-contention scope (PCS)** since scheduling competition is within the process
  - Typically done via priority set by programmer
- Kernel thread scheduled onto available CPU is **system-contention scope (SCS)** – competition among all threads in system





# Pthread Scheduling

---

- API allows specifying either PCS or SCS during thread creation
  - PTHREAD\_SCOPE\_PROCESS schedules threads using PCS scheduling
  - PTHREAD\_SCOPE\_SYSTEM schedules threads using SCS scheduling
- Can be limited by OS – Linux and Mac OS X only allow PTHREAD\_SCOPE\_SYSTEM





# Pthread Scheduling API

---

```
#include <pthread.h>
#include <stdio.h>
#define NUM_THREADS 5
int main(int argc, char *argv[]) {
    int i, scope;
    pthread_t tid[NUM_THREADS];
    pthread_attr_t attr;
    /* get the default attributes */
    pthread_attr_init(&attr);
    /* first inquire on the current scope */
    if (pthread_attr_getscope(&attr, &scope) != 0)
        fprintf(stderr, "Unable to get scheduling scope\n");
    else {
        if (scope == PTHREAD_SCOPE_PROCESS)
            printf("PTHREAD_SCOPE_PROCESS");
        else if (scope == PTHREAD_SCOPE_SYSTEM)
            printf("PTHREAD_SCOPE_SYSTEM");
        else
            fprintf(stderr, "Illegal scope value.\n");
    }
}
```





# Pthread Scheduling API

```
/* set the scheduling algorithm to PCS or SCS */
pthread_attr_setscope(&attr, PTHREAD_SCOPE_SYSTEM);
/* create the threads */
for (i = 0; i < NUM_THREADS; i++)
    pthread_create(&tid[i], &attr, runner, NULL);
/* now join on each thread */
for (i = 0; i < NUM_THREADS; i++)
    pthread_join(tid[i], NULL);
}

/* Each thread will begin control in this function */
void *runner(void *param)
{
    /* do some work ... */
    pthread_exit(0);
}
```



# End of Chapter 6

---

