



IF2230 CPU Scheduling

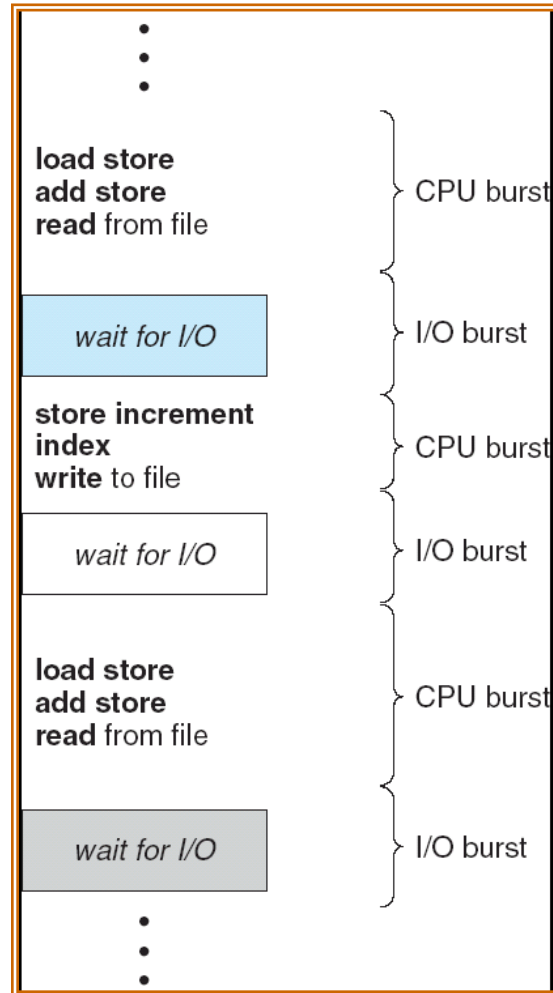


Basic Concepts

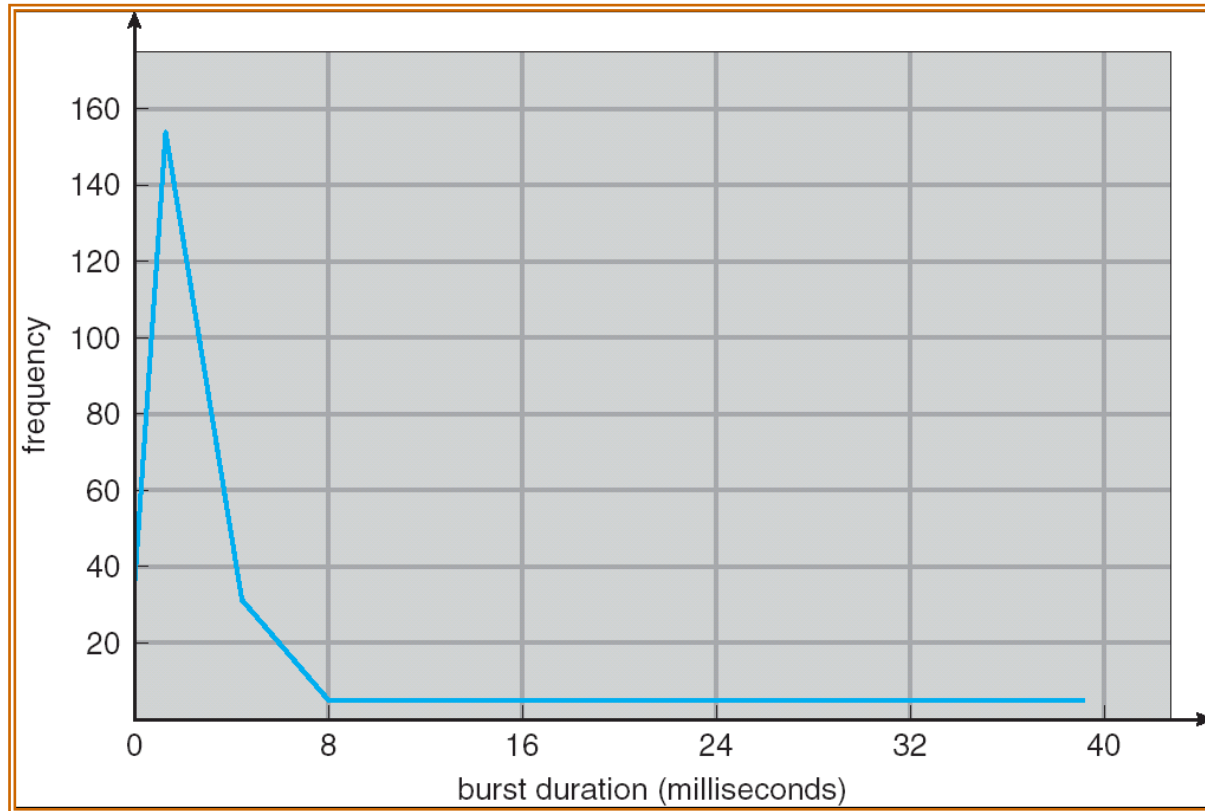
- ▶ Maximum CPU utilization obtained with multiprogramming
- ▶ CPU–I/O Burst Cycle – Process execution consists of a *cycle* of CPU execution and I/O wait
- ▶ CPU burst distribution



Alternating Sequence of CPU And I/O Bursts



Histogram of CPU-burst Times



CPU Scheduler

- ▶ Selects from among the processes in memory that are ready to execute, and allocates the CPU to one of them
- ▶ 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*



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)



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

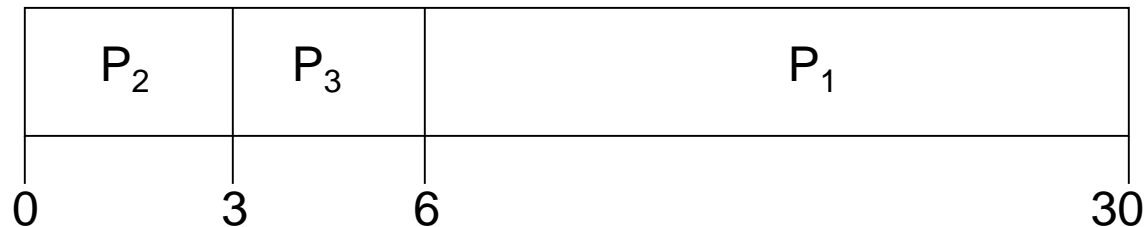


FCFS Scheduling (Cont.)

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: $(6 + 0 + 3)/3 = 3$
- ▶ Much better than previous case
- ▶ *Convoy effect* short process behind long process



Shortest-Job-First (SJF) Scheduling

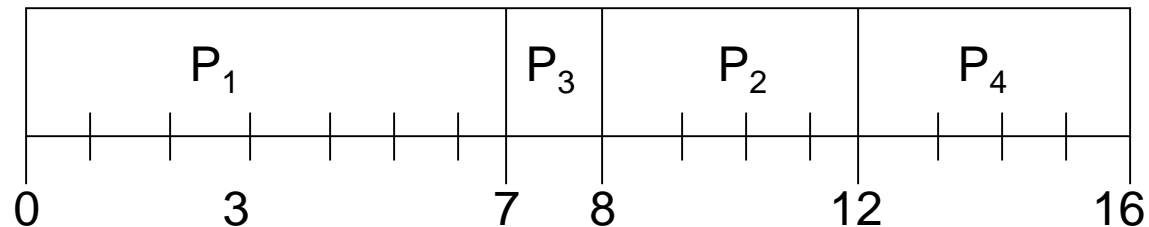
- ▶ Associate with each process the length of its next CPU burst. Use these lengths to schedule the process with the shortest time
- ▶ Two schemes:
 - ▶ nonpreemptive – once CPU given to the process it cannot be preempted until completes its CPU burst
 - ▶ preemptive – if a new process arrives with CPU burst length less than remaining time of current executing process, preempt. This scheme is known as the Shortest-Remaining-Time-First (SRTF)
- ▶ SJF is optimal – gives minimum average waiting time for a given set of processes



Example of Non-Preemptive SJF

<u>Process</u>	<u>Arrival Time</u>	<u>Burst Time</u>
P_1	0.0	7
P_2	2.0	4
P_3	4.0	1
P_4	5.0	4

► SJF (non-preemptive)



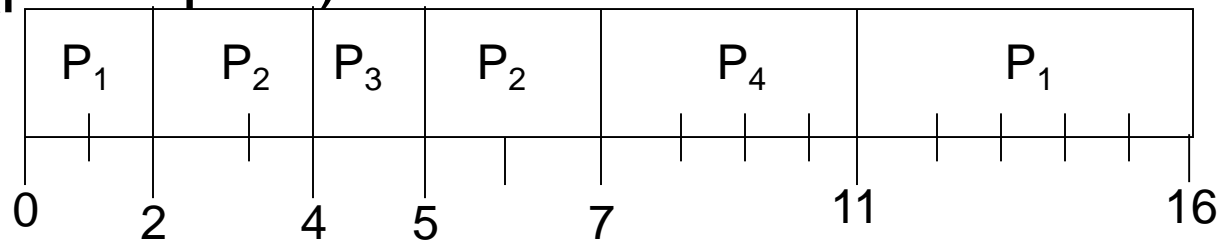
► Average waiting time = $(0 + 6 + 3 + 7)/4 = 4$



Example of Preemptive SJF

<u>Process</u>	<u>Arrival Time</u>	<u>Burst Time</u>
P_1	0.0	7
P_2	2.0	4
P_3	4.0	1
P_4	5.0	4

► SJF (preemptive)



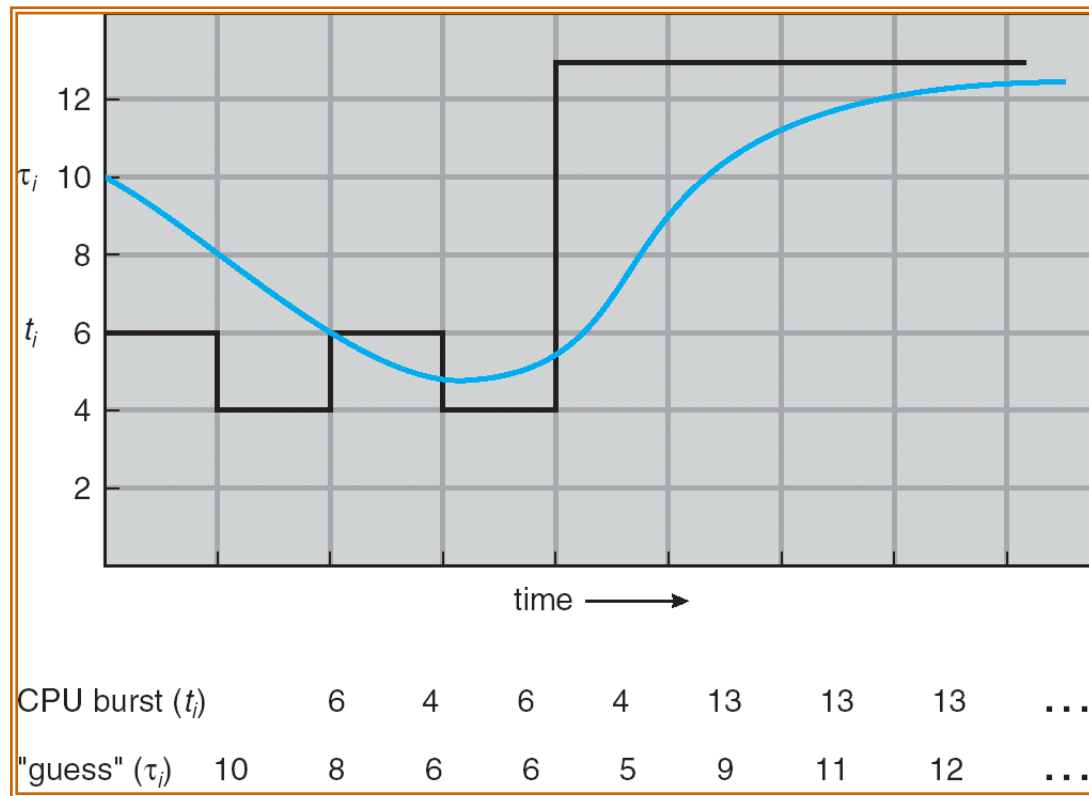
► Average waiting time = $(9 + 1 + 0 + 2)/4 = 3$

Determining Length of Next CPU Burst

- ▶ Can only estimate the length
- ▶ 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. τ_{n+1} = predicted value for the next CPU burst
 3. $\alpha, 0 \leq \alpha \leq 1$
 4. Define : $\tau_{n+1} = \alpha t_n + (1 - \alpha)\tau_n$.



Prediction of the Length of the Next CPU Burst



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



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 a priority scheduling where priority is the 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



Round Robin (RR)

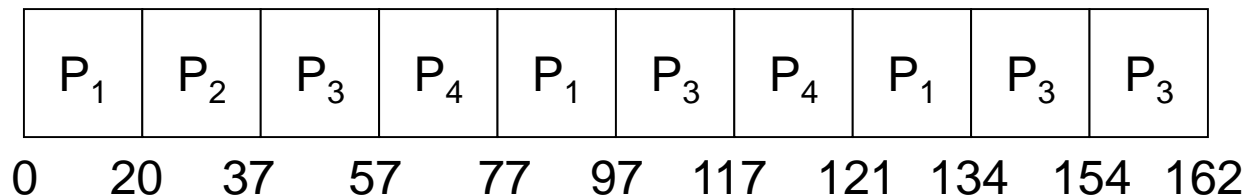
- ▶ Each process gets a small unit of CPU time (*time quantum*), usually 10-100 milliseconds. After this time has elapsed, the process is preempted and added to the end of the ready queue.
- ▶ If there are n processes in the ready queue and the time quantum is q , then each process gets $1/n$ of the CPU time in chunks of at most q time units at once. No process waits more than $(n-1)q$ time units.
- ▶ Performance
 - ▶ q large \Rightarrow FIFO
 - ▶ q small $\Rightarrow q$ must be large with respect to context switch, otherwise overhead is too high



Example of RR with Time Quantum = 20

<u>Process</u>	<u>Burst Time</u>
P_1	53
P_2	17
P_3	68
P_4	24

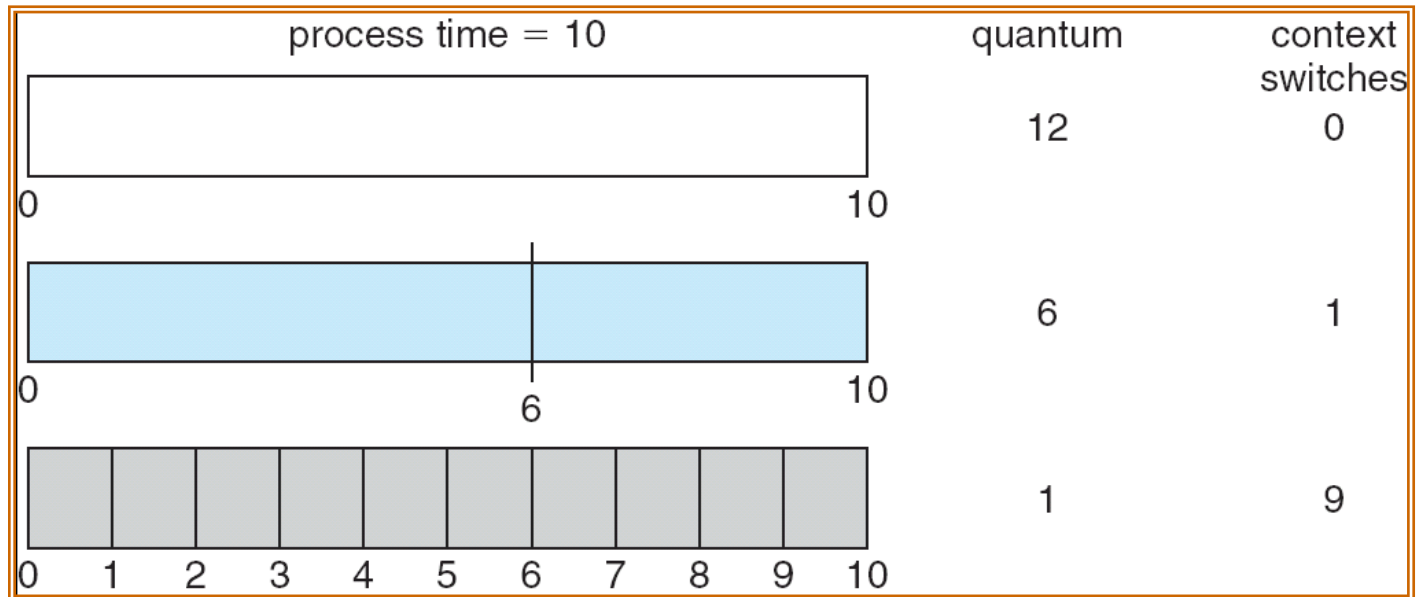
- ▶ The Gantt chart is:



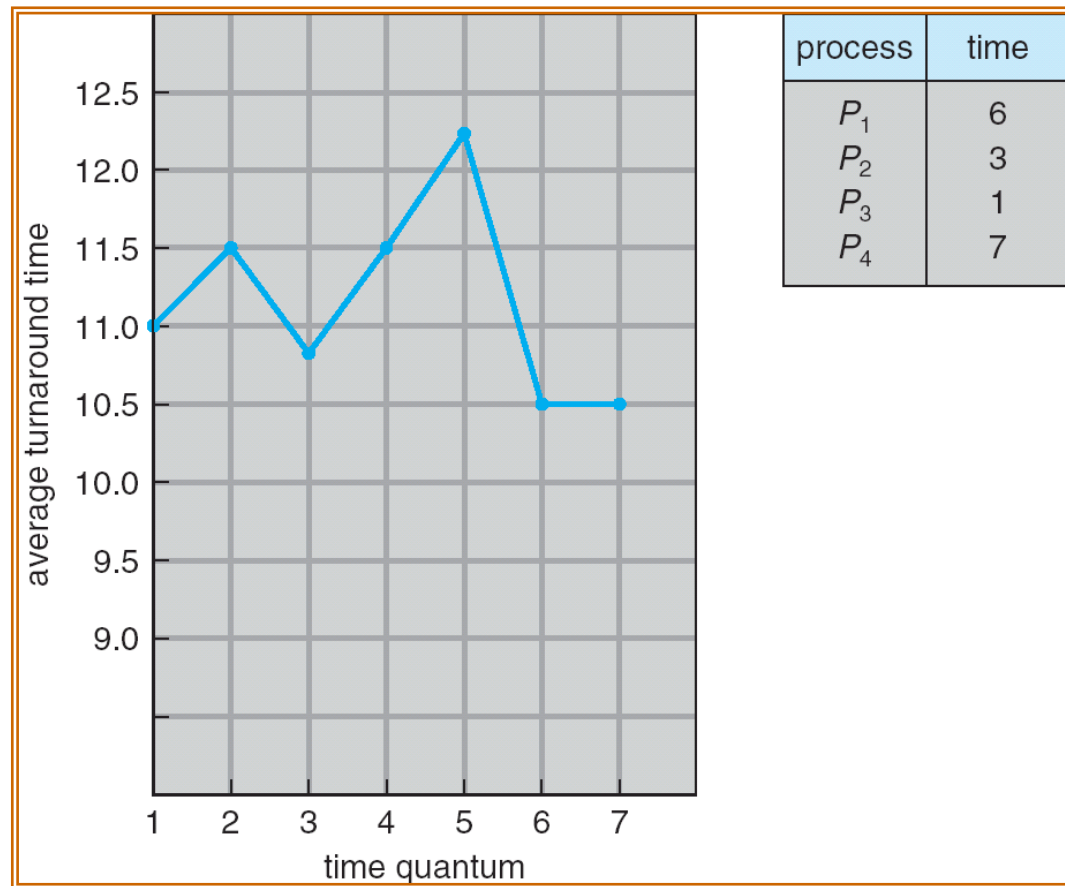
- ▶ Typically, higher average turnaround than SJF, but better *response*



Time Quantum and Context Switch Time



Turnaround Time Varies With The Time Quantum

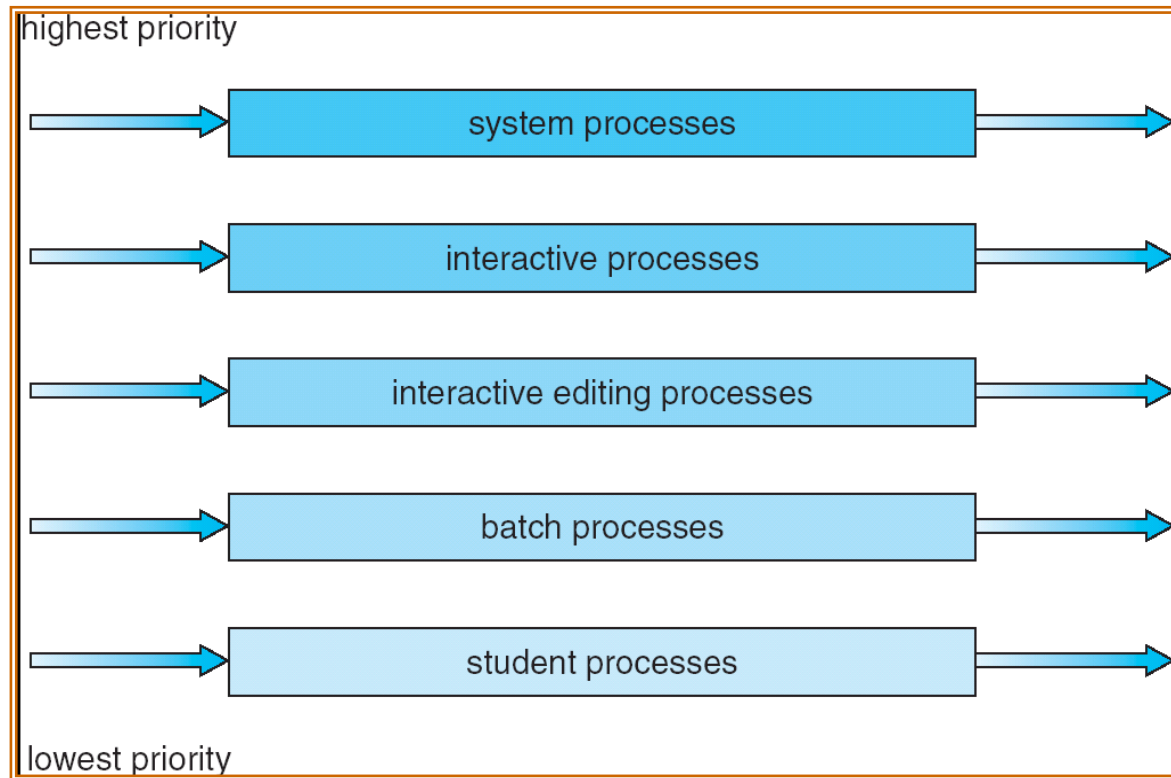


Multilevel Queue

- ▶ Ready queue is partitioned into separate queues:
foreground (interactive)
background (batch)
- ▶ 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



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:**

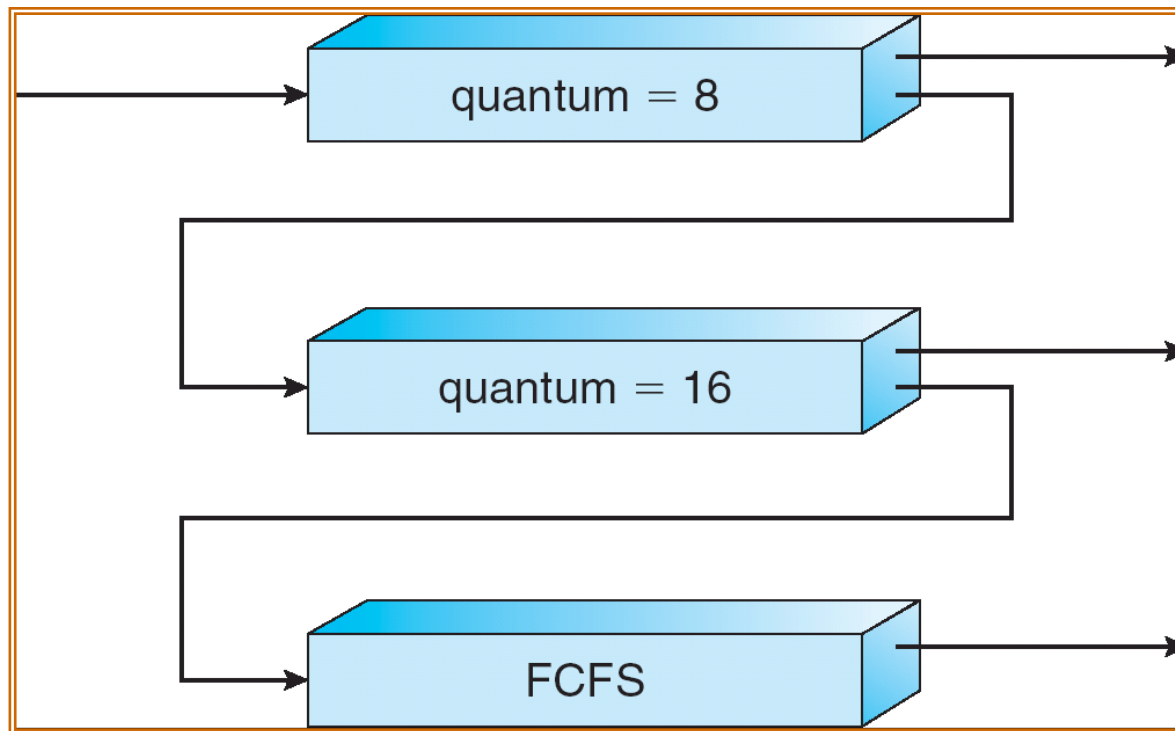
- ▶ Q_0 – RR with time quantum 8 milliseconds
- ▶ Q_1 – RR time quantum 16 milliseconds
- ▶ Q_2 – FCFS

- ▶ **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 .



Multilevel Feedback Queues



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);
}
```

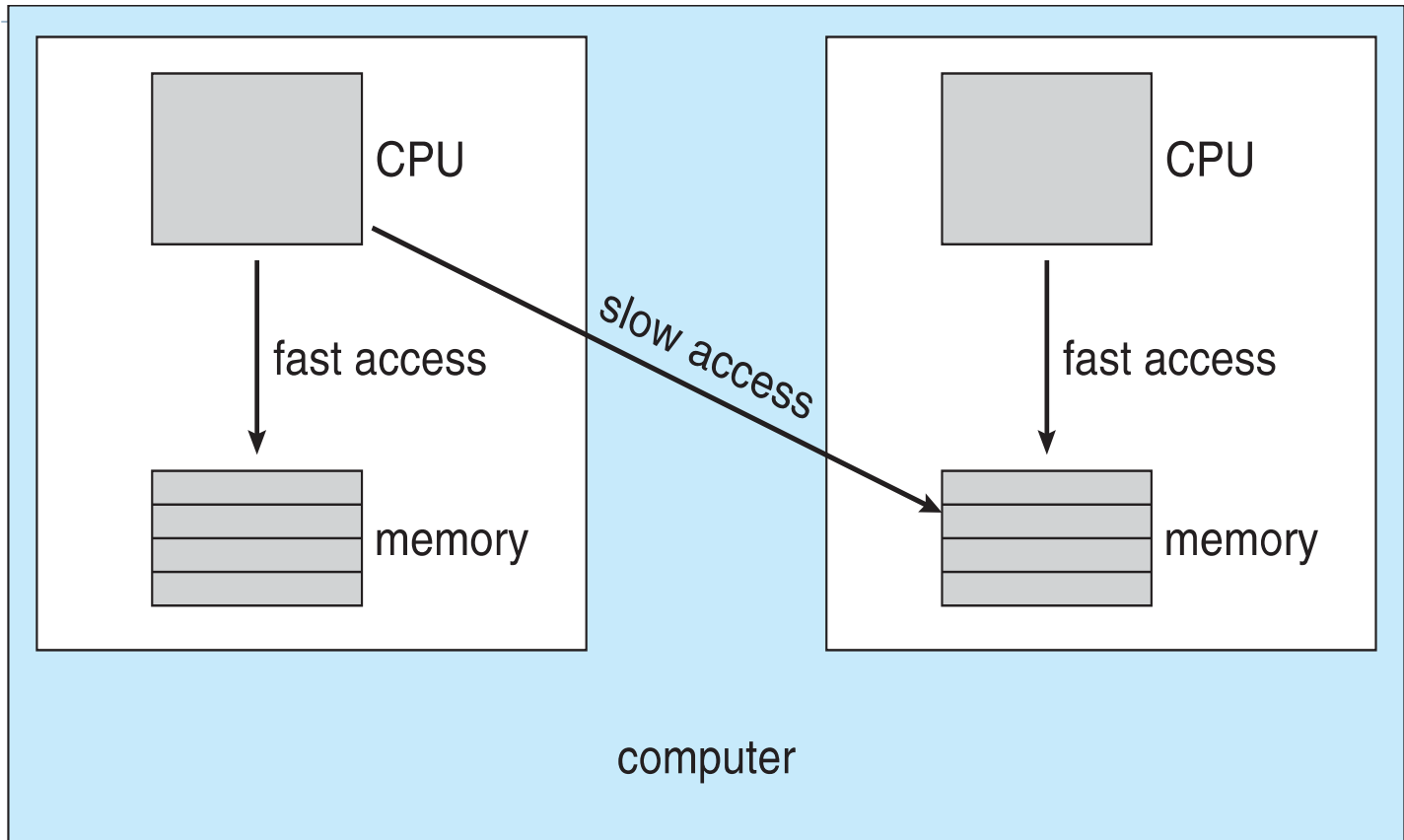


Multiple-Processor Scheduling

- ▶ CPU scheduling more complex when multiple CPUs are available
- ▶ *Homogeneous processors* within a multiprocessor
- ▶ *Load sharing*
- ▶ *Asymmetric multiprocessing* – only one processor accesses the system data structures, alleviating the need for data sharing
- ▶ **Symmetric multiprocessing (SMP)** – each processor is self-scheduling, all processes in common ready queue, or each has its own private queue of ready processes
 - ▶ Currently, most common
- ▶ **Processor affinity** – process has affinity for processor on which it is currently running
 - ▶ **soft affinity**
 - ▶ **hard affinity**
 - ▶ Variations including **processor sets**



NUMA and CPU Scheduling



Note that memory-placement algorithms can also consider affinity



Multiple-Processor Scheduling – Load Balancing

- ▶ If SMP, need to keep all CPUs loaded for efficiency
- ▶ **Load balancing** attempts to keep workload evenly distributed
- ▶ **Push migration** – periodic task checks load on each processor, and if found pushes task from overloaded CPU to other CPUs
- ▶ **Pull migration** – idle processors pulls waiting task from busy processor

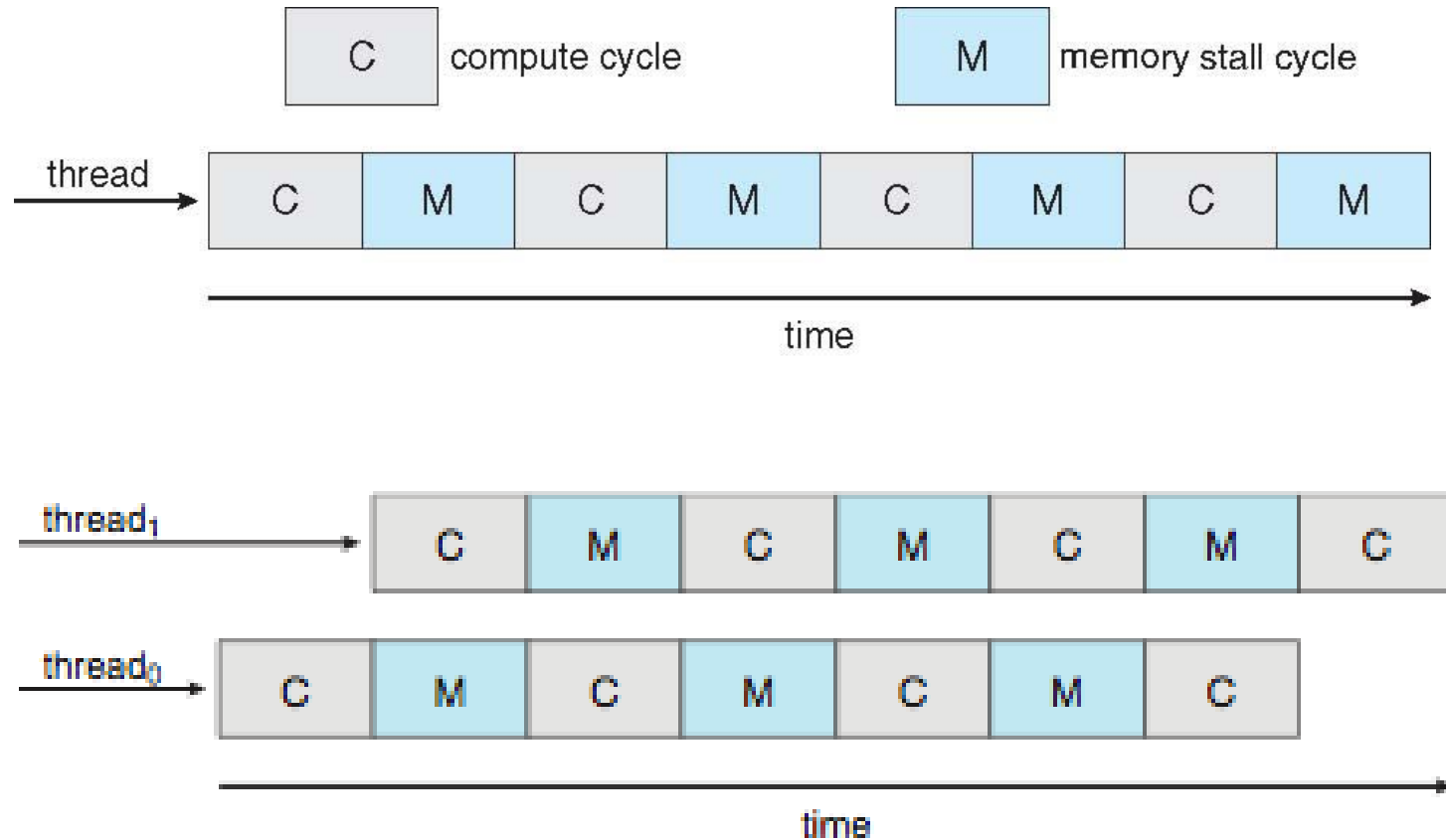


Multicore Processors

- ▶ Recent trend to place multiple processor cores on same physical chip
- ▶ Faster and consumes less power
- ▶ Multiple threads per core also growing
 - ▶ Takes advantage of memory stall to make progress on another thread while memory retrieve happens

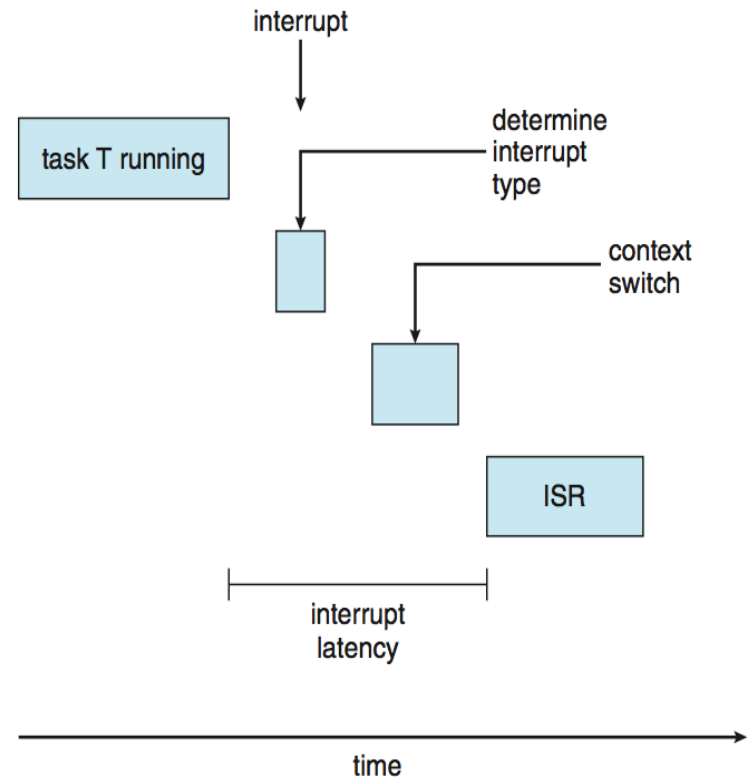


Multithreaded Multicore System



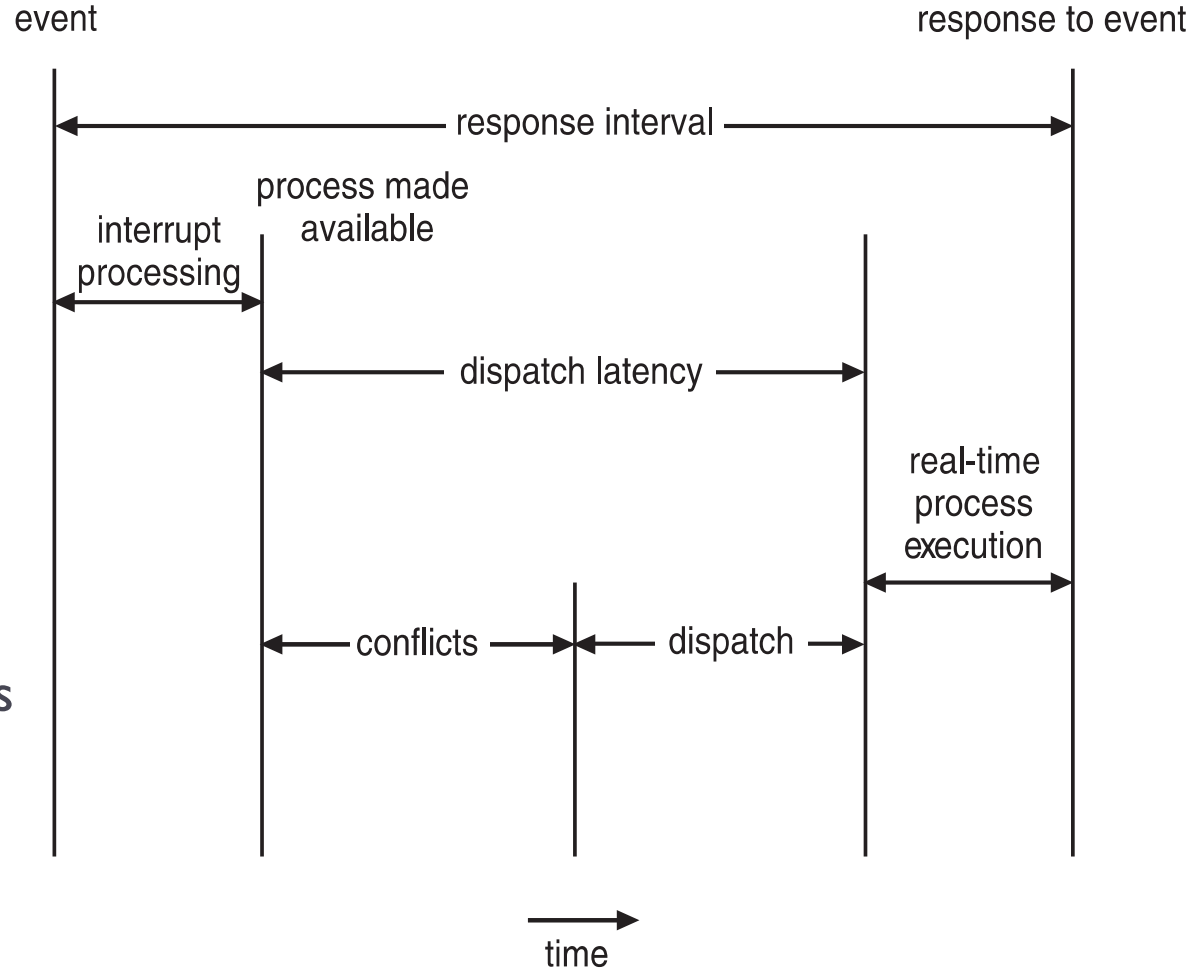
Real-Time CPU Scheduling

- ▶ Can present obvious challenges
- ▶ **Soft real-time systems** – no guarantee as to when critical real-time process will be scheduled
- ▶ **Hard real-time systems** – task must be serviced by its deadline
- ▶ Two types of latencies affect performance
 1. Interrupt latency – time from arrival of interrupt to start of routine that services interrupt
 2. Dispatch latency – time for schedule to take current process off CPU and switch to another



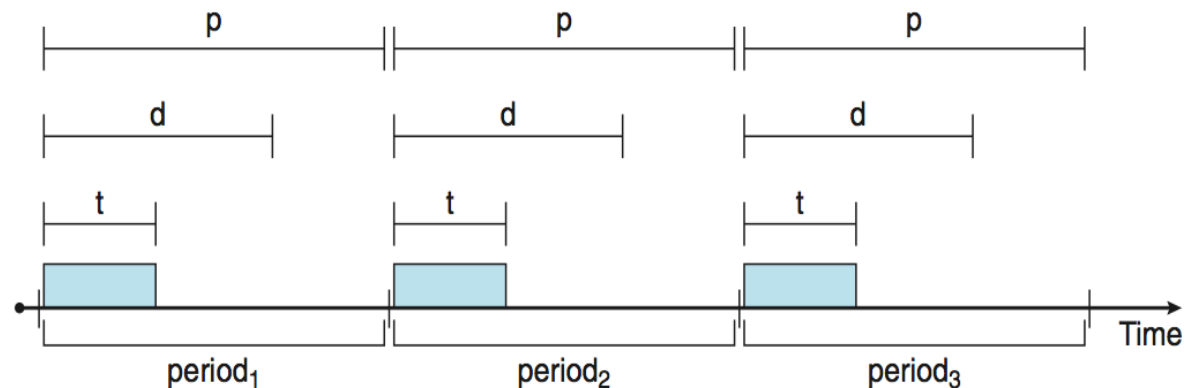
Real-Time CPU Scheduling (Cont.)

- ▶ **Conflict phase of dispatch latency:**
 1. Preemption of any process running in kernel mode
 2. Release by low-priority process of resources needed by high-priority processes



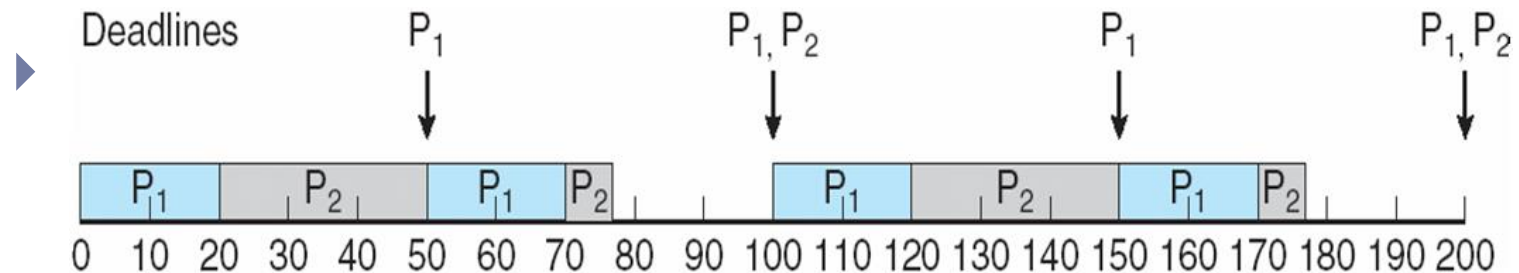
Priority-based Scheduling

- ▶ For real-time scheduling, scheduler must support preemptive, priority-based scheduling
 - ▶ But only guarantees soft real-time
- ▶ For hard real-time must also provide ability to meet deadlines
- ▶ Processes have new characteristics: **periodic** ones require CPU at constant intervals
 - ▶ Has processing time t , deadline d , period p
 - ▶ $0 \leq t \leq d \leq p$
 - ▶ **Rate** of periodic task is $1/p$

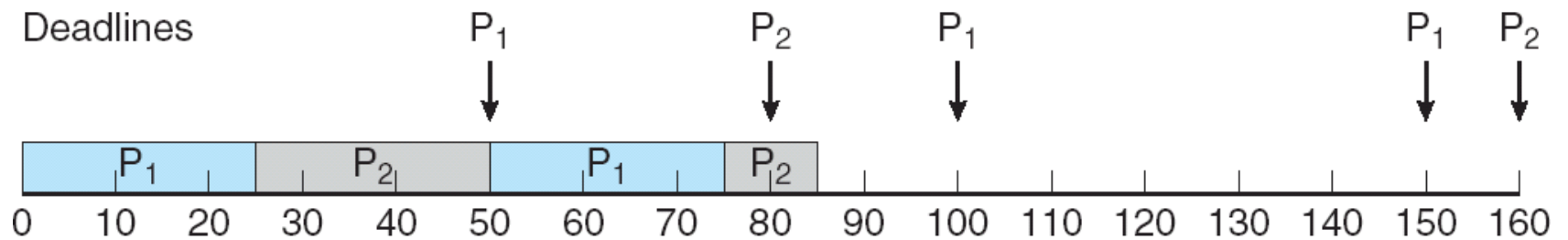


Rate Monotonic Scheduling

- ▶ A priority is assigned based on the inverse of its period
- ▶ Shorter periods = higher priority;
- ▶ Longer periods = lower priority



Missed Deadlines with Rate Monotonic Scheduling

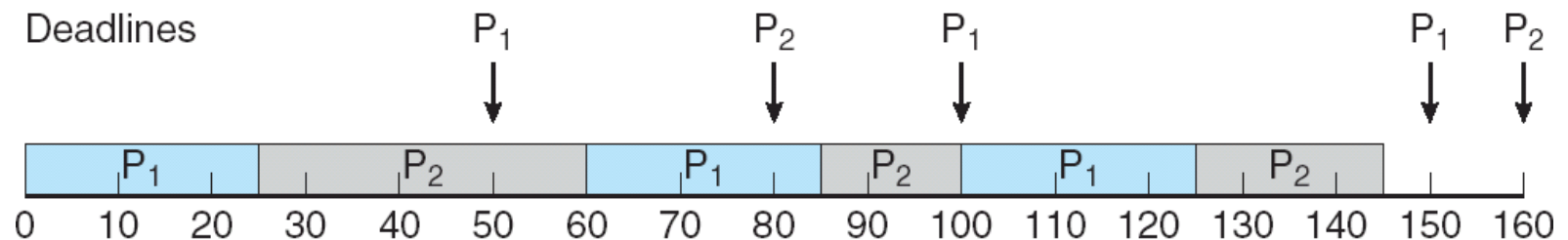


Earliest Deadline First Scheduling (EDF)

- Priorities are assigned according to deadlines:

the earlier the deadline, the higher the priority;

the later the deadline, the lower the priority



POSIX Real-Time Scheduling

- The POSIX.1b standard
- API provides functions for managing real-time threads
- Defines two scheduling classes for real-time threads:
 1. SCHED_FIFO - threads are scheduled using a FCFS strategy with a FIFO queue. There is no time-slicing for threads of equal priority
 2. SCHED_RR - similar to SCHED_FIFO except time-slicing occurs for threads of equal priority
- Defines two functions for getting and setting scheduling policy:
 1. `pthread_attr_getsched_policy(pthread_attr_t *attr, int *policy)`
 2. `pthread_attr_setsched_policy(pthread_attr_t *attr, int policy)`



POSIX Real-Time Scheduling API

```
#include <pthread.h>
#include <stdio.h>
#define NUM_THREADS 5
int main(int argc, char *argv[])
{
    int i, policy;
    pthread_t tid[NUM_THREADS];
    pthread_attr_t attr;
    /* get the default attributes */
    pthread_attr_init(&attr);
    /* get the current scheduling policy */
    if (pthread_attr_getschedpolicy(&attr, &policy) != 0)
        fprintf(stderr, "Unable to get policy.\n");
    else {
        if (policy == SCHED_OTHER) printf("SCHED_OTHER\n");
        else if (policy == SCHED_RR) printf("SCHED_RR\n");
        else if (policy == SCHED_FIFO) printf("SCHED_FIFO\n");
    }
}
```



POSIX Real-Time Scheduling API (Cont.)

```
/* set the scheduling policy - FIFO, RR, or OTHER */
if (pthread_attr_setschedpolicy(&attr, SCHED_FIFO) != 0)
    fprintf(stderr, "Unable to set policy.\n");
/* 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);
}
```



Virtualization and Scheduling

- ▶ Virtualization software schedules multiple guests onto CPU(s)
- ▶ Each guest doing its own scheduling
 - ▶ Not knowing it doesn't own the CPUs
 - ▶ Can result in poor response time
 - ▶ Can effect time-of-day clocks in guests
- ▶ Can undo good scheduling algorithm efforts of guests



Algorithm Evaluation

- ▶ Define the criteria for evaluation
- ▶ Deterministic modeling – takes a particular predetermined workload and defines the performance of each algorithm for that workload
- ▶ Queueing models
- ▶ Simulation
- ▶ Implementation





End of Chapter 5



Linux Scheduling Through Version 2.5

- ▶ Prior to kernel version 2.5, ran variation of standard UNIX scheduling algorithm
- ▶ Version 2.5 moved to constant order $O(1)$ scheduling time
 - ▶ Preemptive, priority based
 - ▶ Two priority ranges: time-sharing and real-time
 - ▶ **Real-time** range from 0 to 99 and **nice** value from 100 to 140
 - ▶ Map into global priority with numerically lower values indicating higher priority
 - ▶ Higher priority gets larger q
 - ▶ Task run-able as long as time left in time slice (**active**)
 - ▶ If no time left (**expired**), not run-able until all other tasks use their slices
 - ▶ All run-able tasks tracked in per-CPU **runqueue** data structure
 - ▶ Two priority arrays (active, expired)
 - ▶ Tasks indexed by priority
 - ▶ When no more active, arrays are exchanged
 - ▶ Worked well, but poor response times for interactive processes

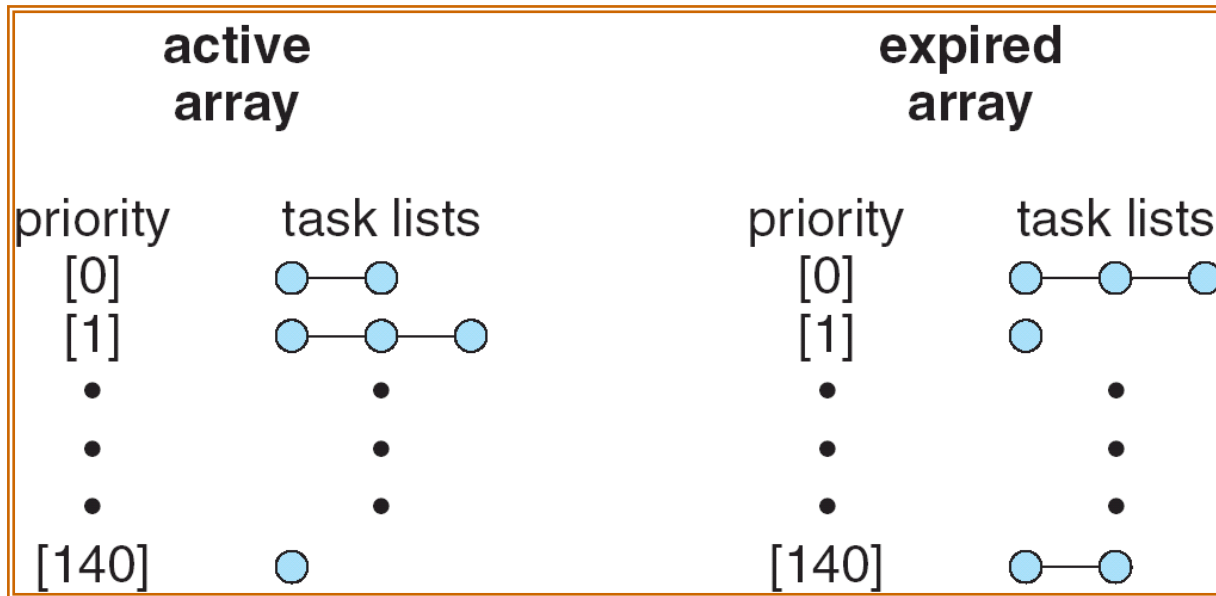


The Relationship Between Priorities and Time-slice length

numeric priority	relative priority		time quantum
0	highest	real-time tasks	200 ms
•			
•			
•			
99			
100		other tasks	10 ms
•			
•			
•			
140			
	lowest		



List of Tasks Indexed According to Priorities

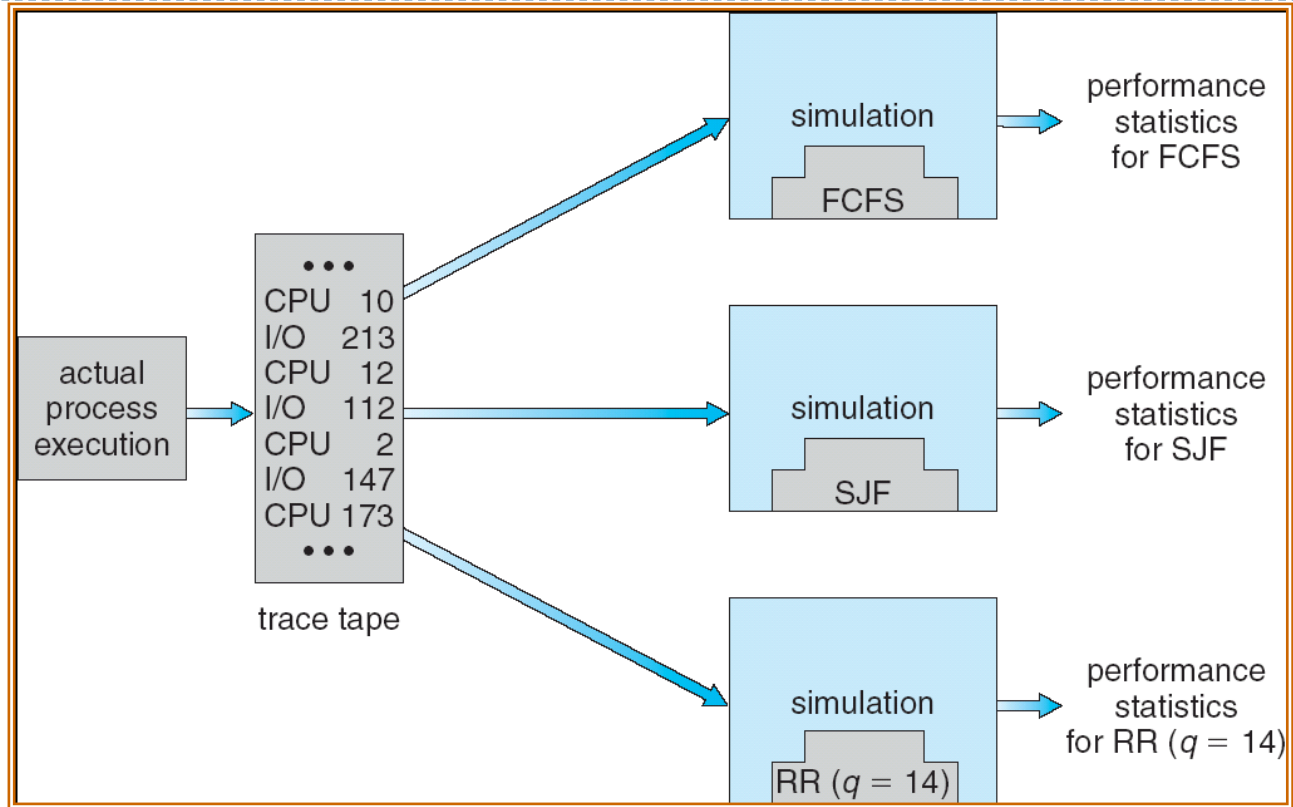


Linux Scheduling in Version 2.6.23 +

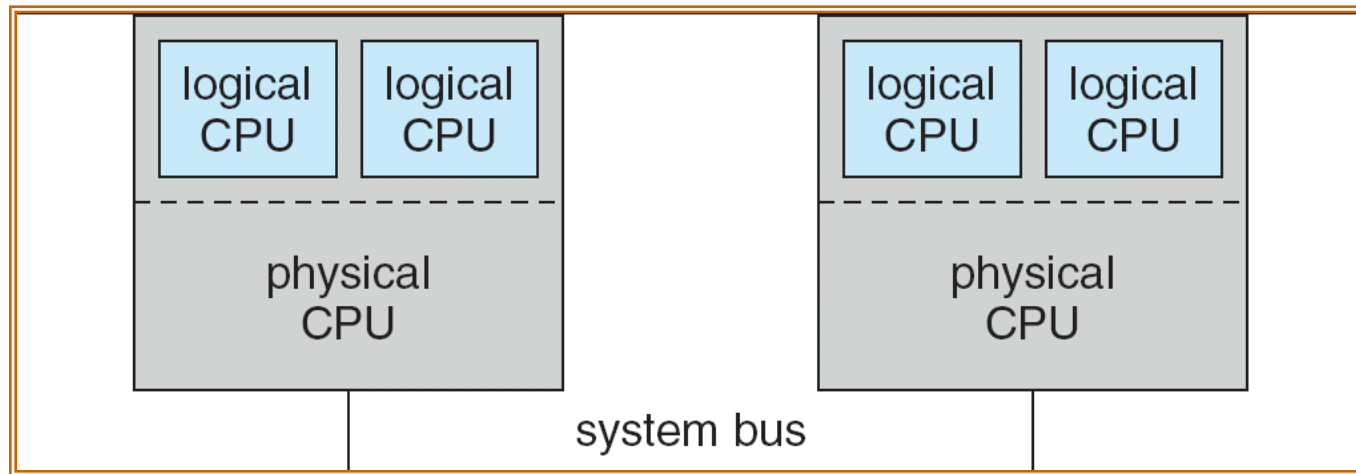
- ▶ **Completely Fair Scheduler (CFS)**
 - ▶ **Scheduling classes**
 - ▶ Each has specific priority
 - ▶ Scheduler picks highest priority task in highest scheduling class
 - ▶ Rather than quantum based on fixed time allotments, based on proportion of CPU time
 - ▶ 2 scheduling classes included, others can be added
 1. default
 2. real-time
 - ▶ Quantum calculated based on **nice value** from -20 to +19
 - ▶ Lower value is higher priority
 - ▶ Calculates **target latency** – interval of time during which task should run at least once
 - ▶ Target latency can increase if say number of active tasks increases
 - ▶ CFS scheduler maintains per task **virtual run time** in variable **vruntime**
 - ▶ Associated with decay factor based on priority of task – lower priority is higher decay rate
 - ▶ Normal default priority yields virtual run time = actual run time
 - ▶ To decide next task to run, scheduler picks task with lowest virtual run time
-



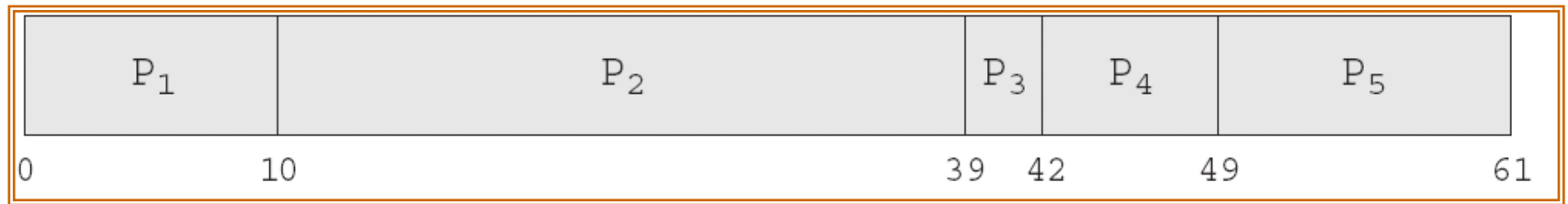
5.15



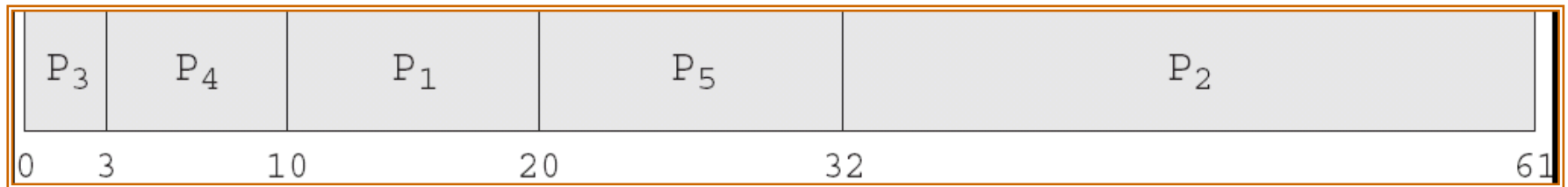
5.08



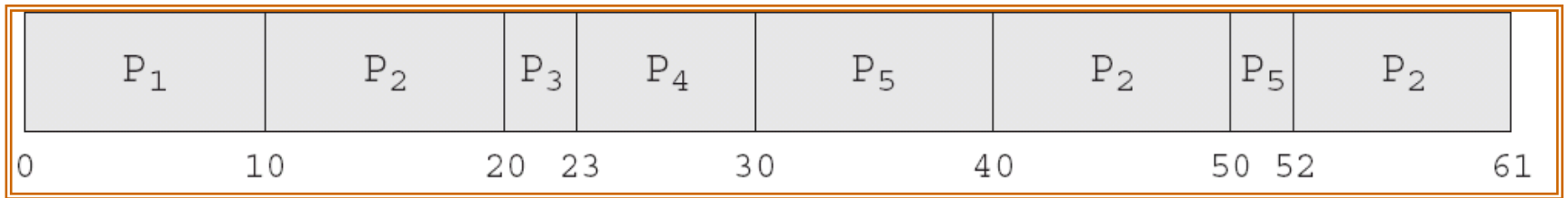
In-5.7



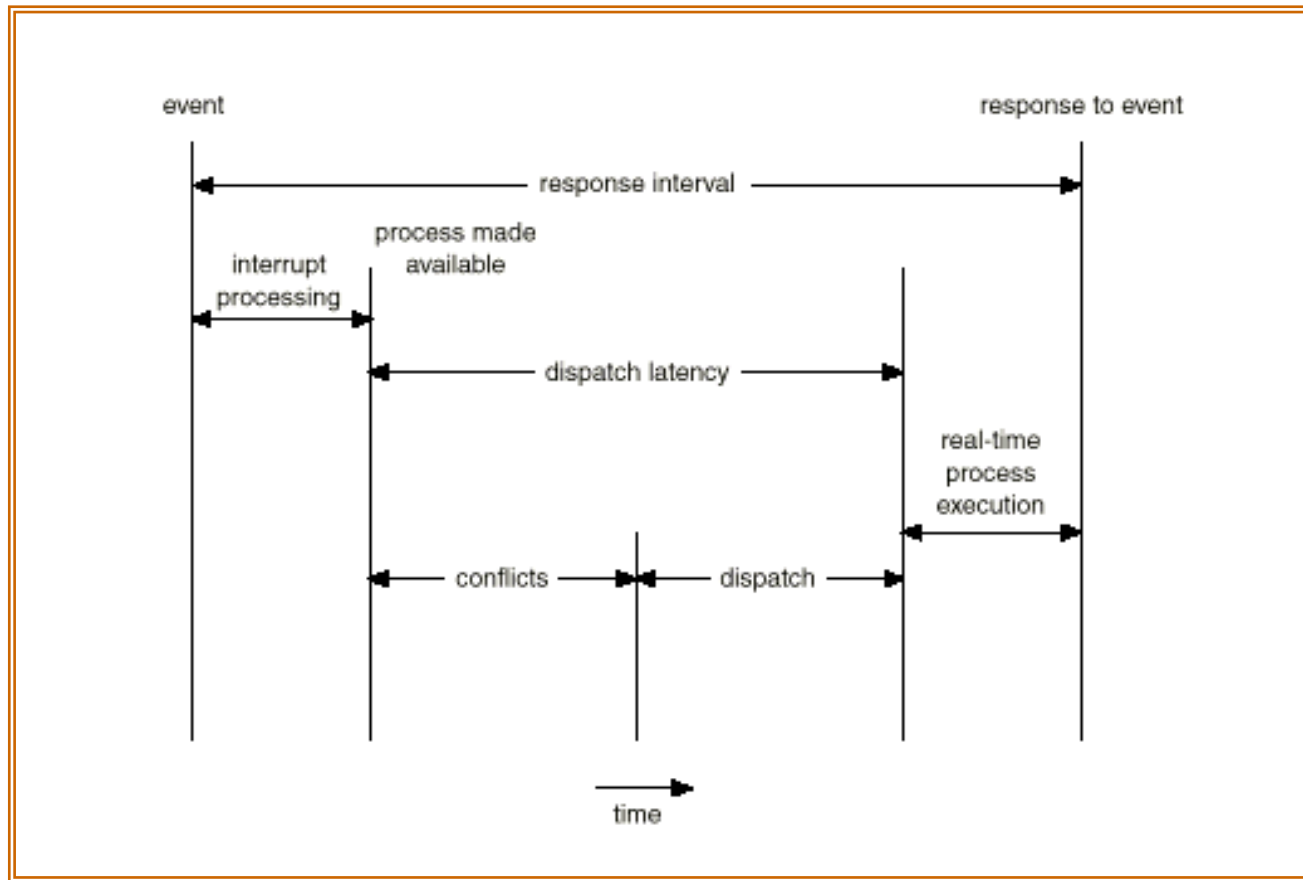
In-5.8



In-5.9



Dispatch Latency



Java Thread Scheduling

- ▶ JVM Uses a Preemptive, Priority-Based Scheduling Algorithm
- ▶ FIFO Queue is Used if There Are Multiple Threads With the Same Priority



Java Thread Scheduling (cont)

JVM Schedules a Thread to Run When:

1. The Currently Running Thread Exits the Runnable State
2. A Higher Priority Thread Enters the Runnable State

* Note – the JVM Does Not Specify Whether Threads are Time-Sliced or Not



Time-Slicing

Since the JVM Doesn't Ensure Time-Slicing, the `yield()` Method

May Be Used:

```
while (true) {  
    // perform CPU-intensive task  
    ...  
    Thread.yield();  
}
```

This Yields Control to Another Thread of Equal Priority



Thread Priorities

<u>Priority</u>	<u>Comment</u>
Thread.MIN_PRIORITY Thread Priority	Minimum
Thread.MAX_PRIORITY Priority	Maximum Thread
Thread.NORM_PRIORITY	Default Thread Priority

Priorities May Be Set Using setPriority() method:
setPriority(Thread.NORM_PRIORITY + 2);

