

6. CPU Scheduling

Sungyoung Lee

*College of Engineering
KyungHee University*

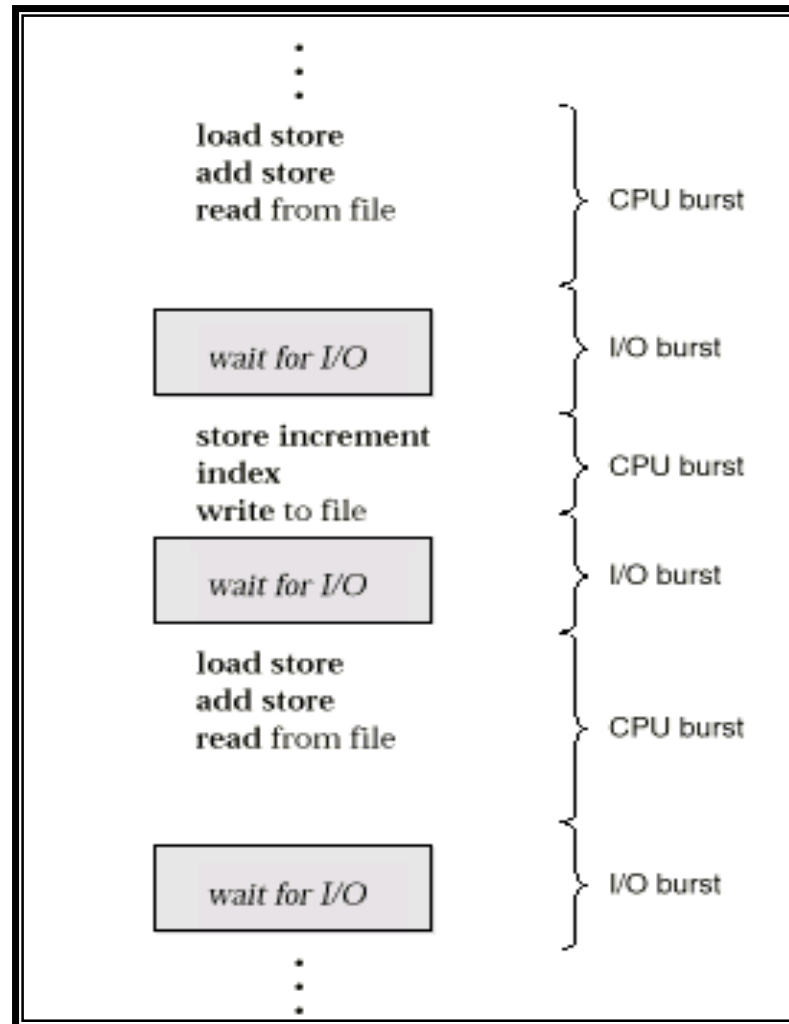
Contents

- n *Basic Concepts*
- n *Scheduling Criteria*
- n *Scheduling Algorithms*
- n *Multiple-Processor Scheduling*
- n *Real-Time Scheduling*
- n *Algorithm Evaluation*

Basic Concepts

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

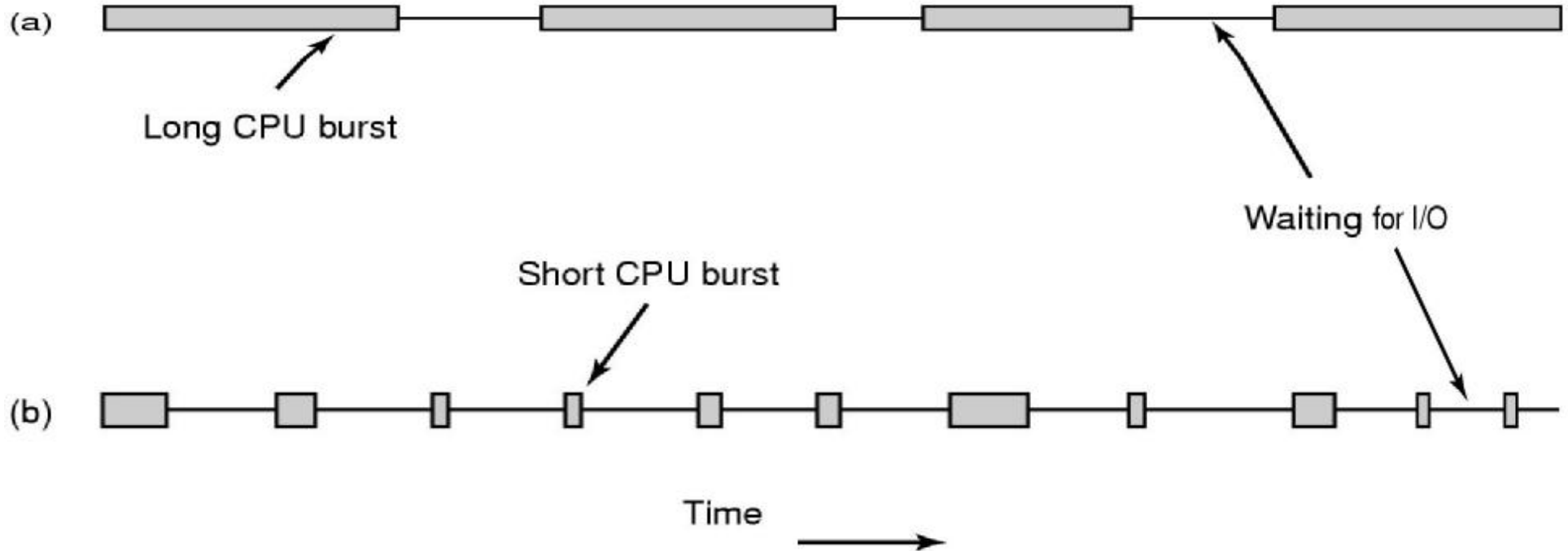
Alternating Sequence of CPU And I/O Bursts



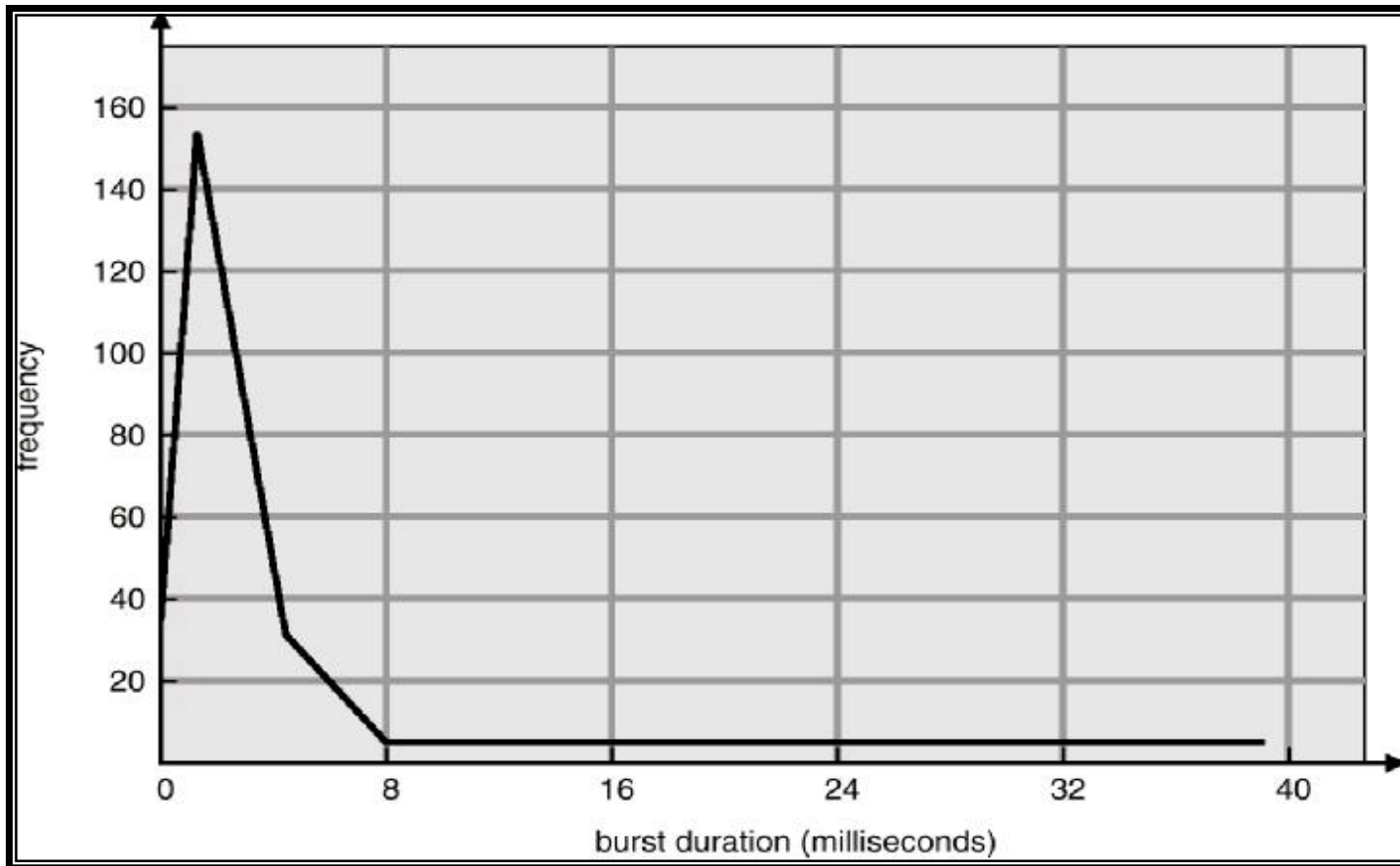
CPU burst vs. I/O burst

n (a) A CPU-bound process

n (b) An I/O-bound process

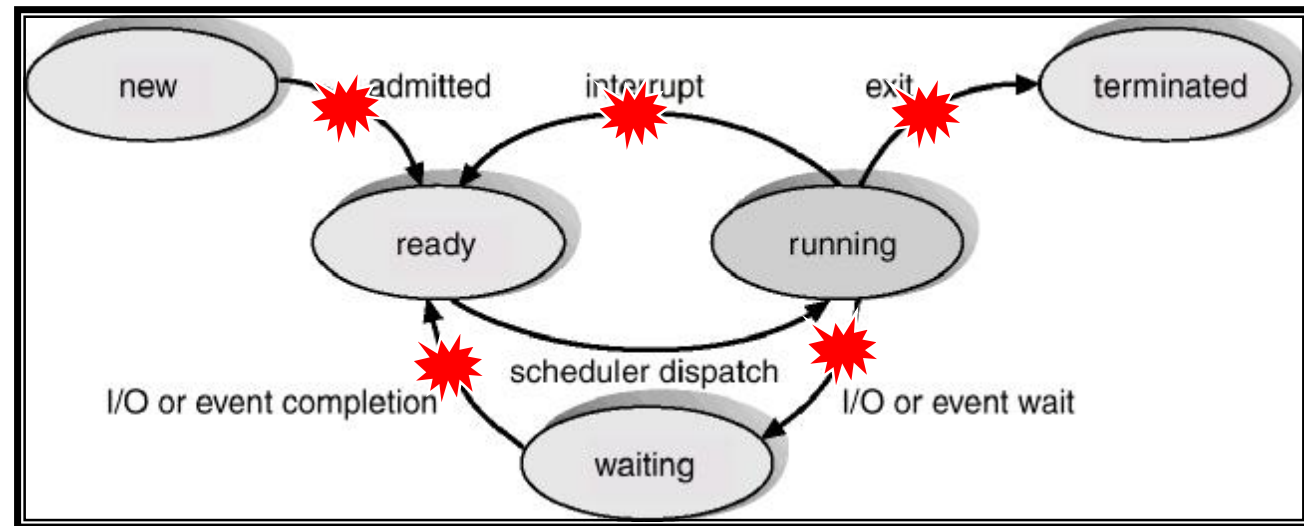


Histogram of CPU-burst Times



CPU Scheduler

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



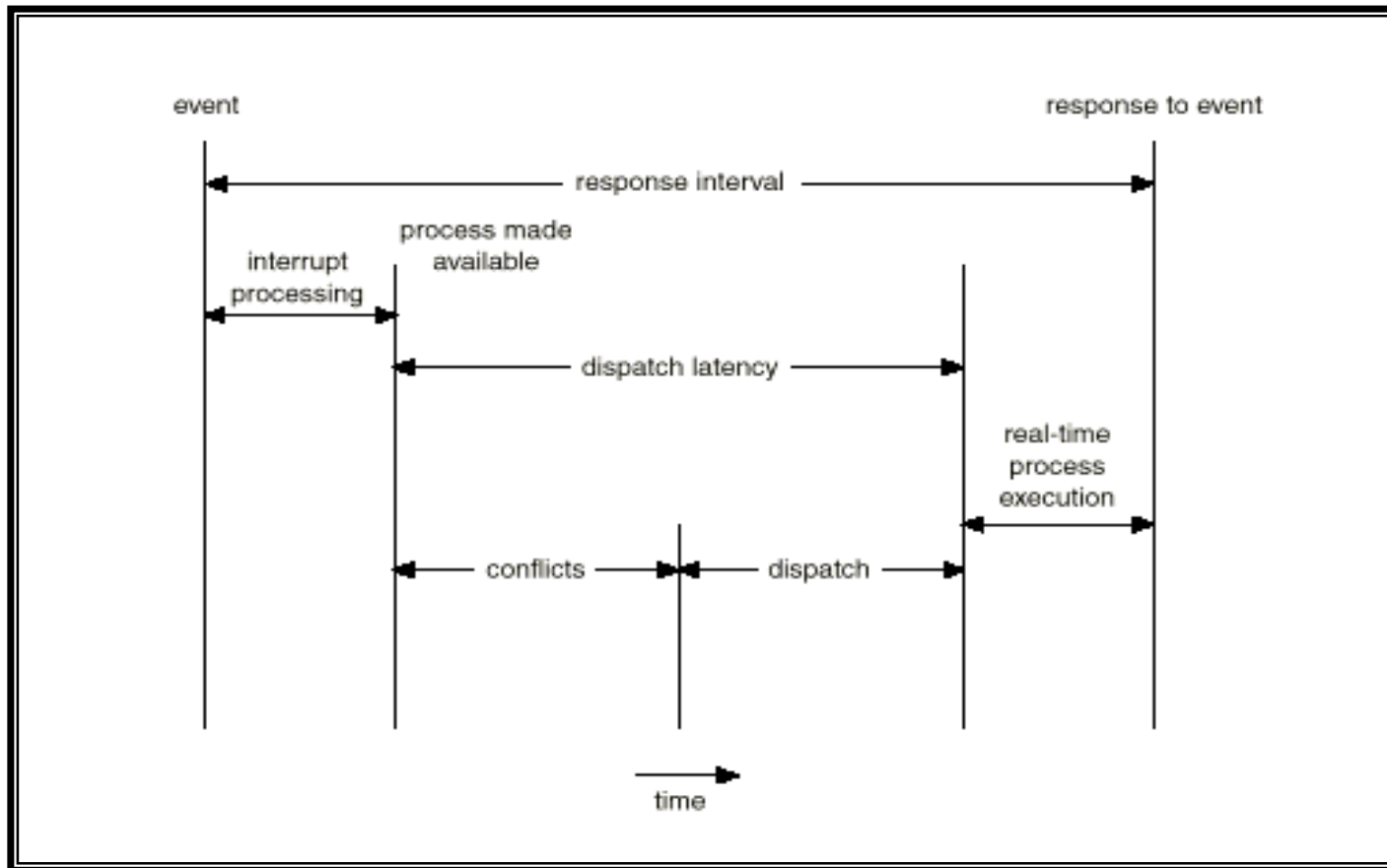
- n Scheduling under 1 and 4 is *nonpreemptive*
- n All other scheduling is *preemptive*

Dispatcher

- n Dispatcher module gives control of the CPU to the process selected by the short-term scheduler; this involves:
 - ü switching context
 - ü switching to user mode
 - ü jumping to the proper location in the user program to restart that program

- n *Dispatch latency*
 - ü time it takes for the dispatcher to stop one process and start another running

Dispatch Latency



Preemptive vs. Non-preemptive

n Non-preemptive scheduling

- ü The scheduler waits for the running job to explicitly (voluntarily) block
- ü Scheduling takes place only when
 - § A process switches from running to waiting state
 - § A process terminates

n Preemptive scheduling

- ü The scheduler can interrupt a job and force a context switch
- ü What happens
 - § If a process is preempted in the midst of updating the shared data?
 - § If a process in system call is preempted?

Scheduling Criteria

n CPU utilization

- ü keep the CPU as busy as possible

n Throughput

- ü # of processes that complete their execution per time unit

n Turnaround time

- ü amount of time to execute a particular process

n Waiting time

- ü amount of time a process has been waiting in the ready queue

n Response time

- ü 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

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

Scheduling Goals

n All systems

- ü **Fairness**: giving each process a fair share of the CPU
- ü **Balance**: keeping all parts of the system busy

n Batch systems

- ü **Throughput**: maximize jobs per hour
- ü **Turnaround time**: minimize time between submission and termination
- ü **CPU utilization**: keep the CPU busy all the time

Scheduling Goals (Cont'd)

n Interactive systems

- ü **Response time**: minimize average time spent on ready queue
- ü **Waiting time**: minimize average time spent on wait queue
- ü **Proportionality**: meet users' expectations

n Real-time systems

- ü **Meeting deadlines**: avoid losing data
- ü **Predictability**: avoid quality degradation in multimedia systems

n Starvation

- ü A situation where a process is prevented from making progress because another process has the resource it requires.
 - § Resource could be the CPU or a lock
- ü A poor scheduling policy can cause starvation
 - § If a high-priority process always prevents a low-priority process from running on the CPU
- ü Synchronization can also cause starvation
 - § One thread always beats another when acquiring a lock
 - § Constant supply of readers always blocks out writers

n First-Come, First-Served

- ü Jobs are scheduled in order that they arrive
- ü “Real-world” scheduling of people in lines
 - § e.g. supermarket, bank tellers, McDonalds, etc.
- ü Typically, non-preemptive
- ü Jobs are treated equally: no starvation

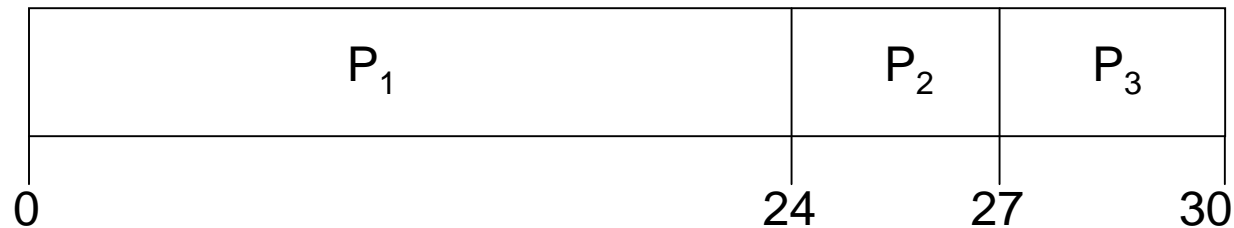
n Problems

- ü Average waiting time can be large if small jobs wait behind long ones
 - § Basket vs. cart
- ü May lead to poor overlap of I/O and CPU

First-Come, First-Served (FCFS) Scheduling

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

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



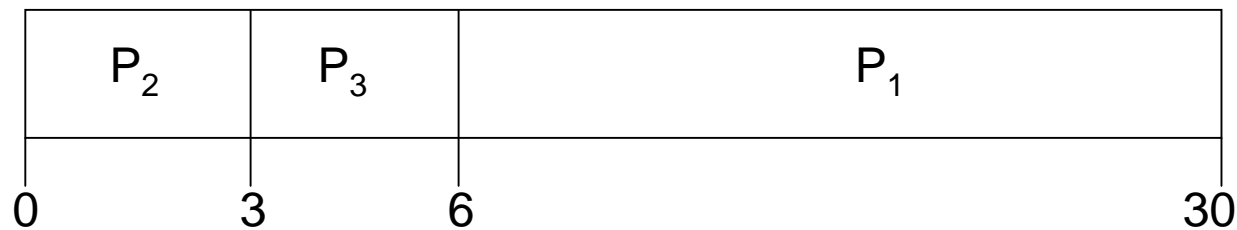
- n Waiting time for $P_1 = 0$; $P_2 = 24$; $P_3 = 27$
- n Average waiting time: $(0 + 24 + 27)/3 = 17$

FCFS Scheduling (Cont'd)

n Suppose that the processes arrive in the order

$$P_2, P_3, P_1.$$

n The Gantt chart for the schedule is:



n Waiting time for $P_1 = 6$; $P_2 = 0$; $P_3 = 3$

n Average waiting time: $(6 + 0 + 3)/3 = 3$

n Much better than previous case

n *Convoy effect*

ü short process behind long process

n Shortest Job First

- ü Choose the job with the smallest expected CPU burst
- ü Can prove that SJF has optimal min. average waiting time
 - § Only when all jobs are available simultaneously
- ü Non-preemptive

n Problems

- ü Impossible to know size of future CPU burst
- ü Can you make a reasonable guess?
- ü Can potentially starve

Shortest-Job-First (SJF) Scheduling

- n Associate with each process the length of its next CPU burst. Use these lengths to schedule the process with the shortest time

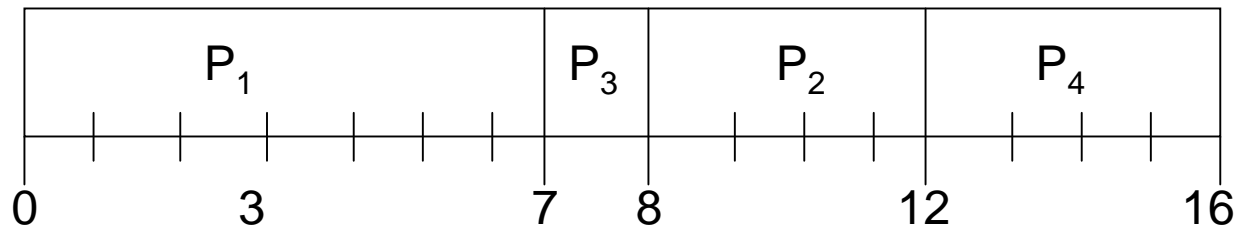
- n Two schemes:
 - ü Nonpreemptive
 - § Once CPU given to the process it cannot be preempted until completes its CPU burst
 - ü Preemptive
 - § If a new process arrives with CPU burst length less than remaining time of current executing process, preempt
 - § This scheme is known as the Shortest-Remaining-Time-First (SRTF)

- n SJF is *optimal*
 - ü 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

n SJF (non-preemptive)

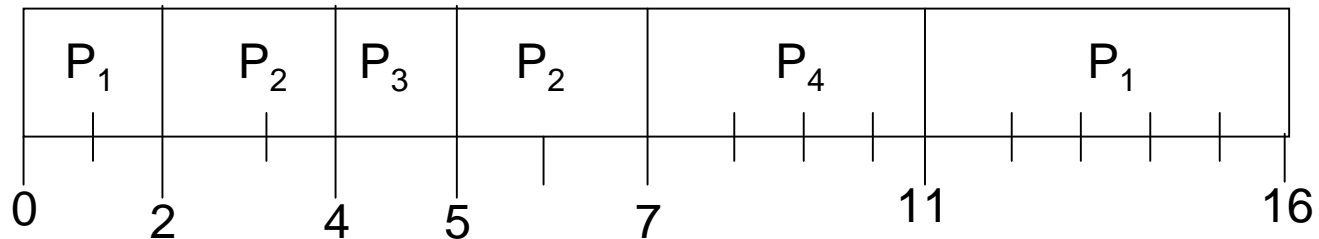


n Average waiting time = $(0 + 6 + 3 + 7)/4 = 4$

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

n SJF (preemptive) (= SRTF)



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

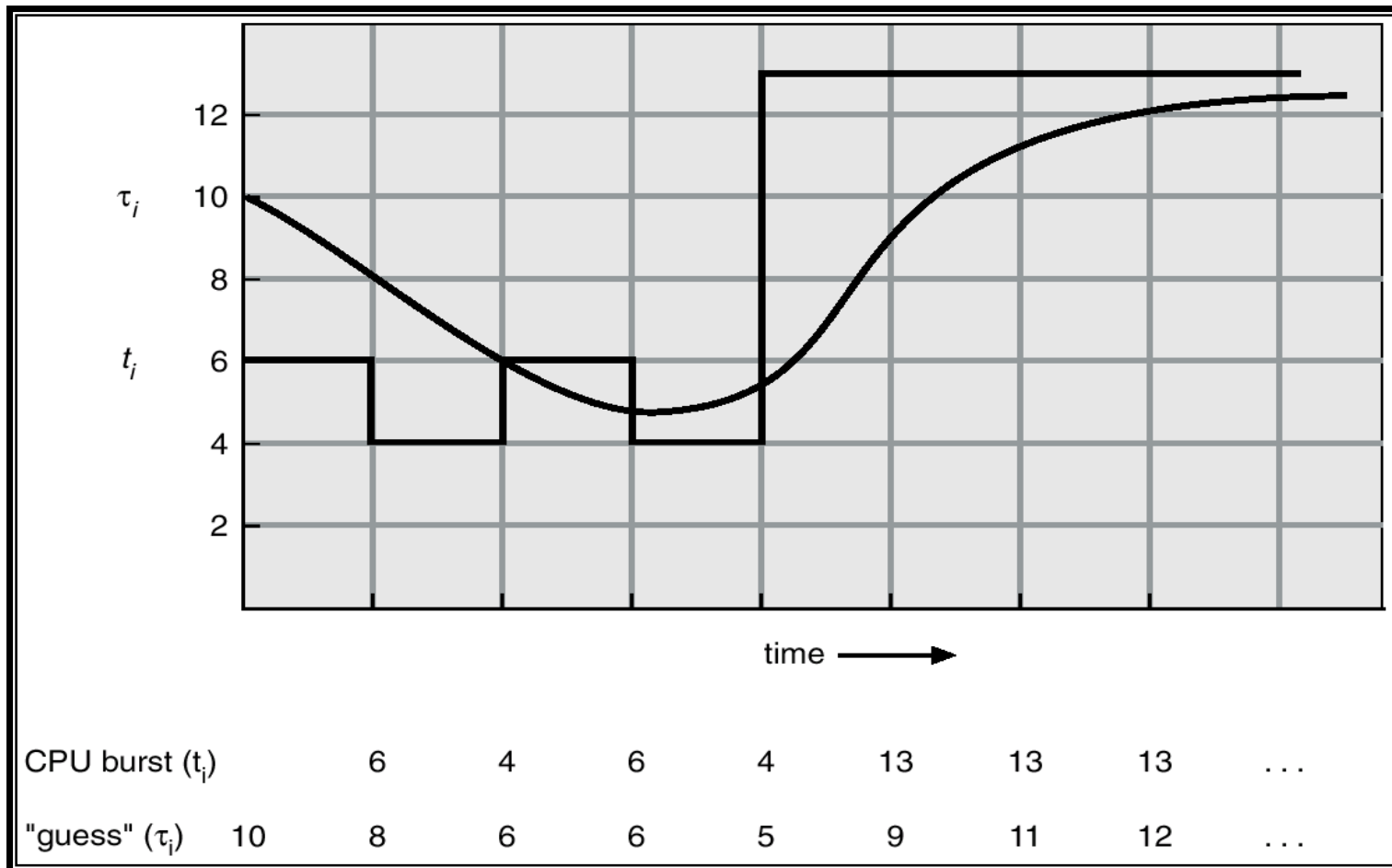
Determining Length of Next CPU Burst

- n Can only estimate the length
- n 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. t_{n+1} = predicted value for the next CPU burst
3. $a, 0 \leq a \leq 1$
4. Define :

$$t_{n+1} = a t_n + (1 - a)t_n.$$

Prediction of the Length of the Next CPU Burst



Examples of Exponential Averaging

n $\alpha = 0$

ü $\tau_{n+1} = \tau_n$

ü Recent history does not count

n $\alpha = 1$

ü $\tau_{n+1} = t_n$

ü Only the actual last CPU burst counts

n If we expand the formula, we get:

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

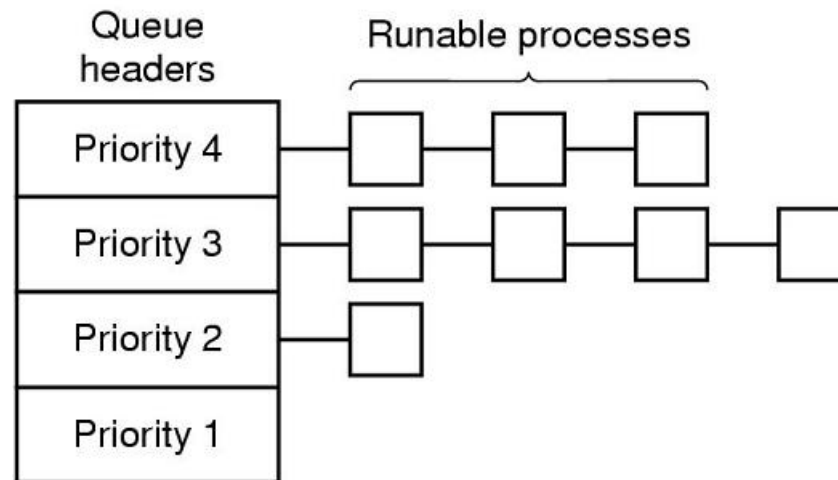
n Since both α and $(1 - \alpha)$ are less than or equal to 1, each successive term has less weight than its predecessor

Priority Scheduling

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

Priority Scheduling

- n Abstractly modeled as multiple “priority queues”
 - ü Put ready job on Q associated with its priority



Round Robin (RR)

- n Each process gets a small unit of CPU time (*time quantum*), usually 10-100 milliseconds
 - ü After this time has elapsed, the process is preempted and added to the end of the ready queue

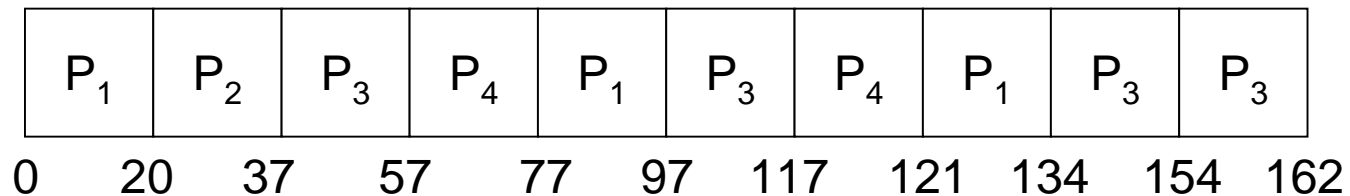
- n If there are n processes in the ready queue and the time quantum is q , then each process gets $1/n$ of the CPU time in chunks of at most q time units at once
 - ü No process waits more than $(n-1)q$ time units

- n Performance
 - ü q large \Rightarrow FIFO
 - ü q small \Rightarrow q must be large with respect to context switch, otherwise overhead is too high

Example of RR with Time Quantum = 20

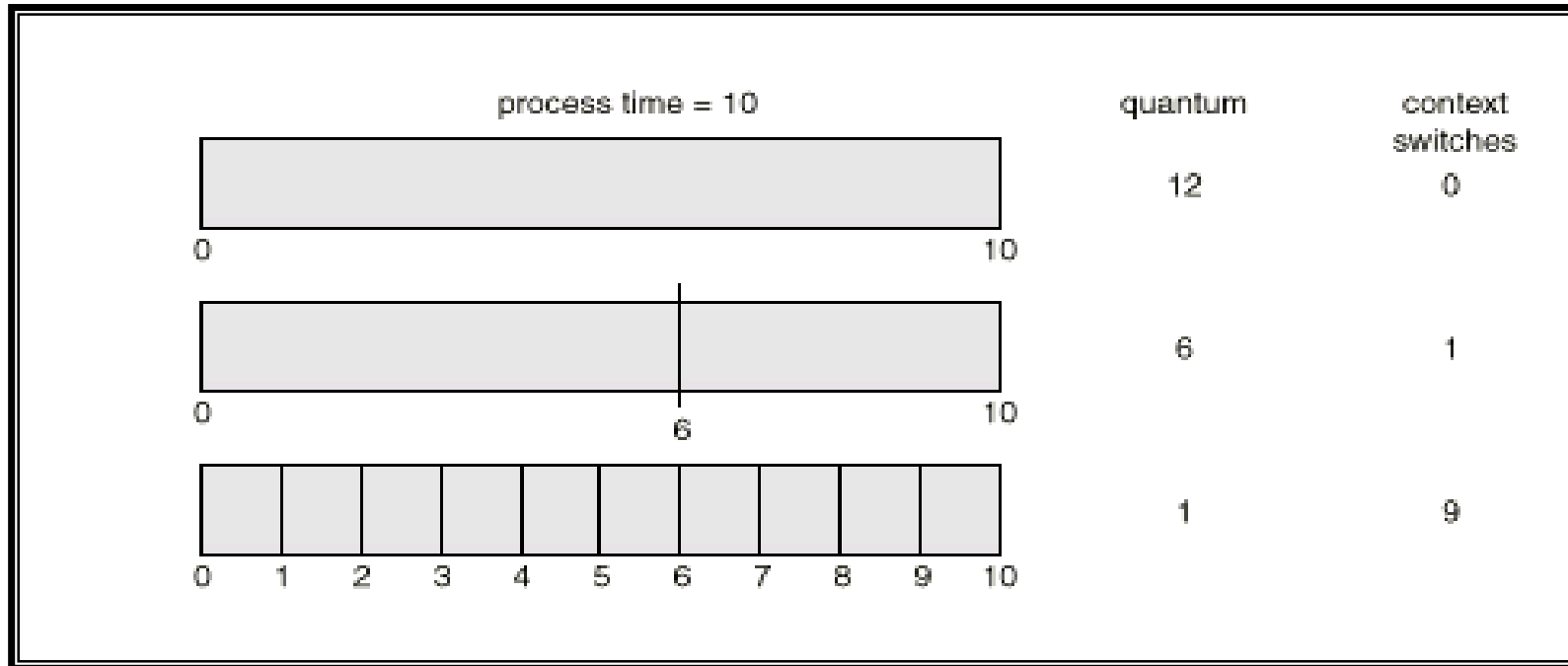
<u>Process</u>	<u>Burst Time</u>
P_1	53
P_2	17
P_3	68
P_4	24

n The Gantt chart is:

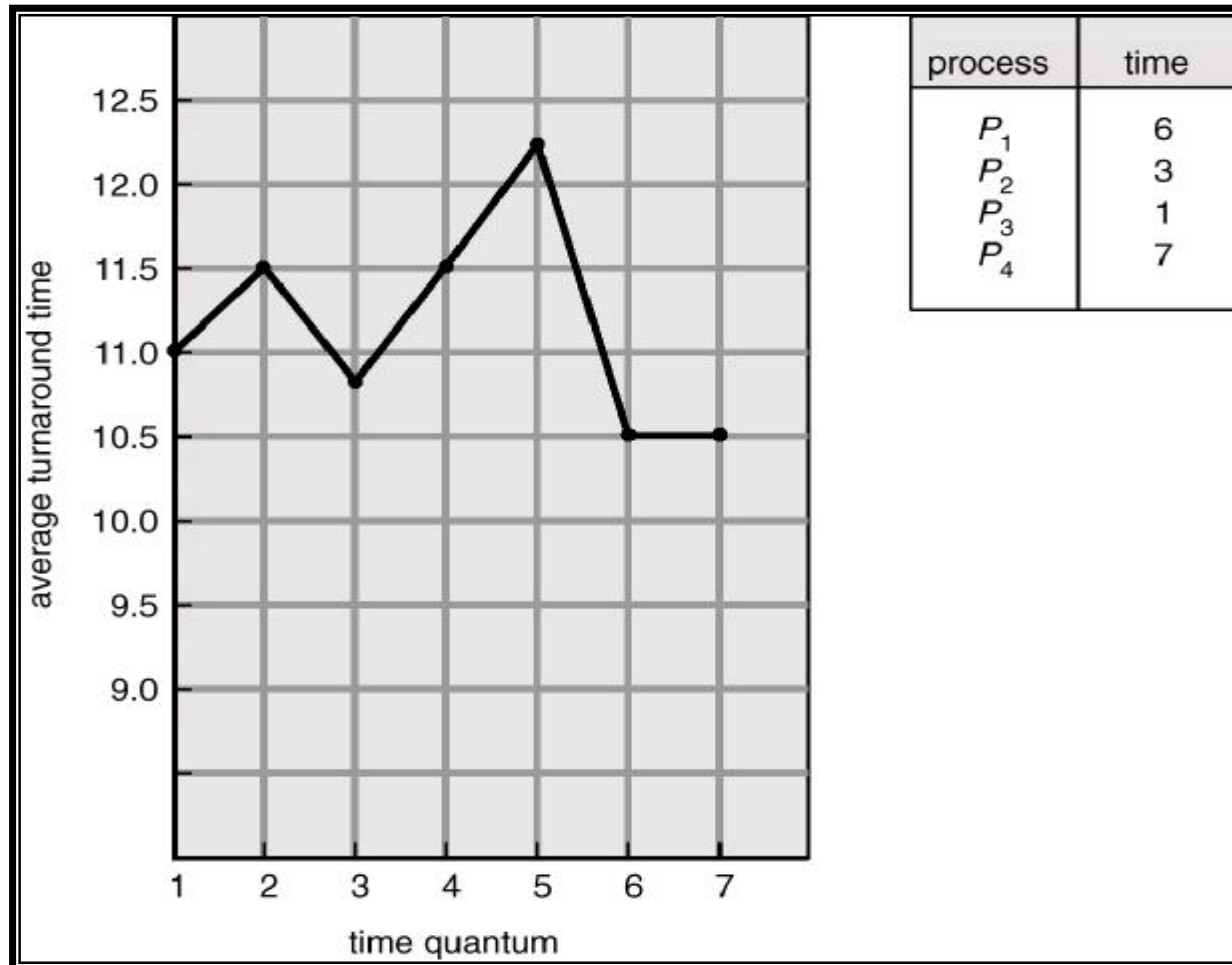


n Typically, higher average turnaround than SJF, but better *response*

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

n What do you set the quantum to be?

- ü quantum $\rightarrow \infty$: FIFO

- ü quantum $\rightarrow 0$: processor sharing

- ü If small, then context switches are frequent incurring high overhead (CPU utilization drops)

- ü If large, then response time drops

- ü A rule of thumb: 80% of the CPU bursts should be shorter than the time quantum

n Treats all jobs equally

- ü Multiple background jobs?

Combining Algorithms

- n Scheduling algorithms can be combined in practice
 - ü Have multiple queues
 - ü Pick a different algorithm for each queue
 - ü Have a mechanism to schedule among queues
 - ü And maybe, move processes between queues

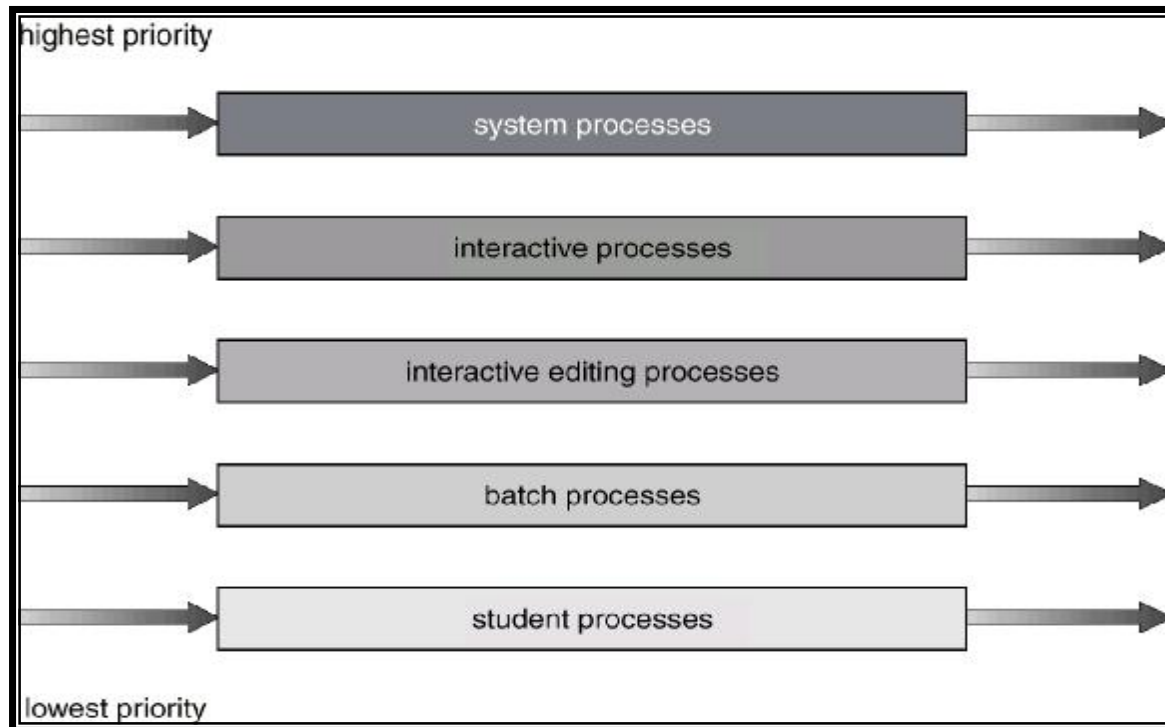
Multilevel Queue

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

- n Each queue has its own scheduling algorithm:
 - ü foreground – RR
 - ü background – FCFS

- n Scheduling must be done between the queues
 - ü Fixed priority scheduling
 - § (i.e., serve all from foreground then from background) Possibility of starvation
 - ü 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

- n A process can move between the various queues
 - ü aging can be implemented this way

- n Multilevel-feedback-queue scheduler defined by the following parameters:
 - ü number of queues
 - ü scheduling algorithms for each queue
 - ü method used to determine when to upgrade a process
 - ü method used to determine when to demote a process
 - ü method used to determine which queue a process will enter when that process needs service

Example of Multilevel Feedback Queue

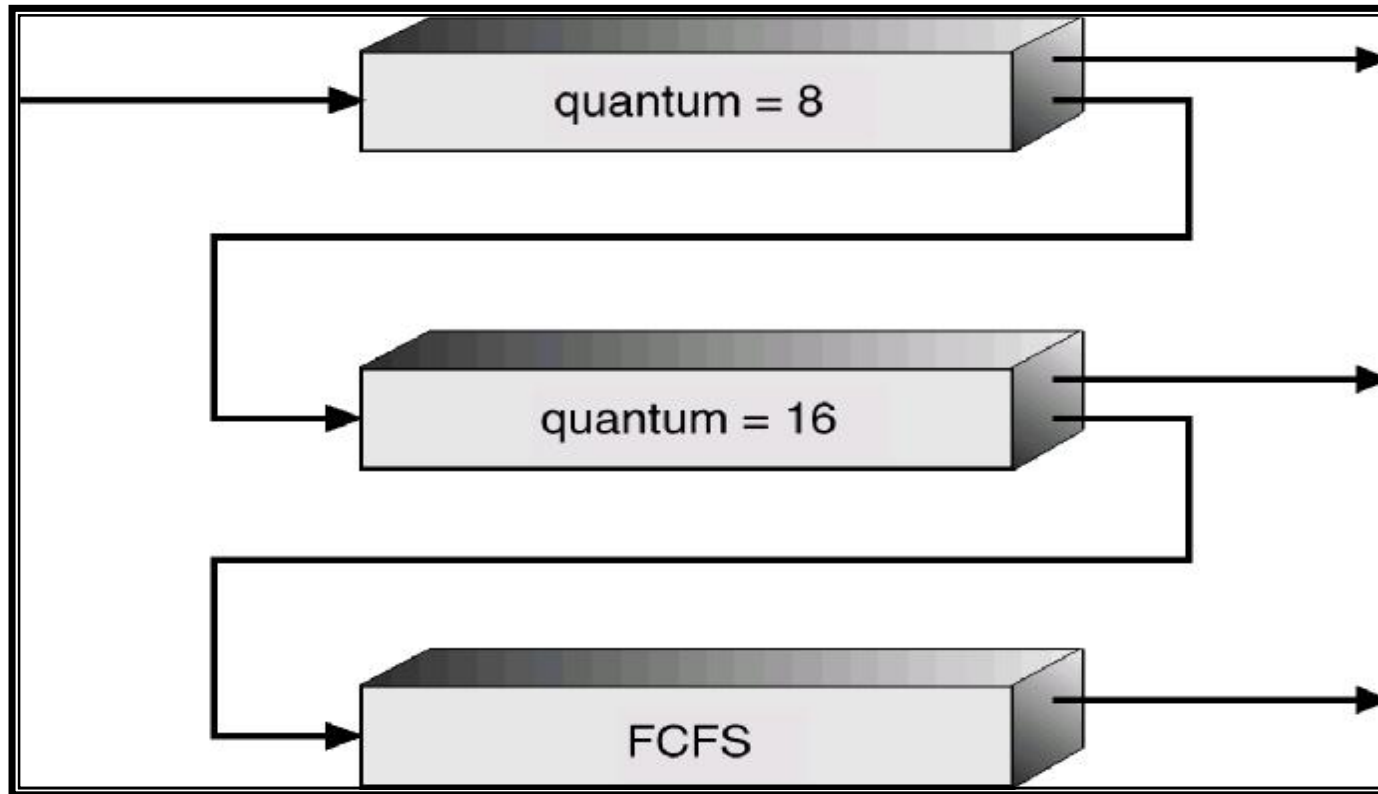
n Three queues:

- ü Q_0 – time quantum 8 milliseconds
- ü Q_1 – time quantum 16 milliseconds
- ü Q_2 – FCFS

n Scheduling

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

Multilevel Feedback Queues



- n The canonical UNIX scheduler uses a MLFQ
 - ü 3 – 4 classes spanning ~170 priority levels
 - § Timeshare, System, Real-time, Interrupt (Solaris 2)
 - ü Priority scheduling across queues, RR within a queue
 - § The process with the highest priority always runs
 - § Processes with the same priority are scheduled RR
 - ü Processes dynamically change priority
 - § Increases over time if process blocks before end of quantum
 - § Decreases over time if process uses entire quantum

UNIX Scheduler (Cont'd)

n Motivation

- ü The idea behind the UNIX scheduler is to reward interactive processes over CPU hogs
- ü Interactive processes typically run using short CPU bursts
 - § They do not finish quantum before waiting for more input
- ü Want to minimize response time
 - § Time from keystroke (putting process on ready queue) to executing the handler (process running)
 - § Don't want editor to wait until CPU hog finishes quantum
- ü This policy delays execution of CPU-bound jobs

Multiple-Processor Scheduling

- n CPU scheduling more complex when multiple CPUs are available
- n *Homogeneous processors* within a multiprocessor
 - ü UMA (Uniform Memory Access)
- n *Load sharing*
- n *Asymmetric multiprocessing*
 - ü Only one processor accesses the system data structures, alleviating the need for data sharing
 - ü Not efficient

Real-Time Scheduling

n *Hard real-time systems*

- ü required to complete a critical task within a guaranteed amount of time

n *Soft real-time computing*

- ü requires that critical processes receive priority over less fortunate ones

n *Static vs. Dynamic priority scheduling*

- ü Static: Rate-Monotonic algorithm

- ü Dynamic: EDF (Earliest Deadline First) algorithm

n Hard real-time

- ü Must complete a critical task within a guaranteed amount of time
- ü Resource reservation
 - § A process is submitted along with its resource requirements
- ü Requires worst-case timing analysis
 - § Minimize unavoidable and unforeseeable variation in the amount of time to execute a particular process
 - § Very difficult in a system with secondary storage or virtual memory
- ü Typically composed of special-purpose software running on dedicated hardware with limited functionality

Real-Time Scheduling (Cont'd)

n Soft real-time

ü Less restrictive

§ Multimedia, high-speed interactive graphics, etc.

ü May cause an unfair allocation of resources and may result in longer delays, or even starvations, for some processes

ü Requirements

§ The system must have priority scheduling, and real-time processes must have the highest priority

(The priority of real-time processes must not degrade over time)

§ Dispatch latency must be small

Real-Time Scheduling (Cont'd)

n Problem

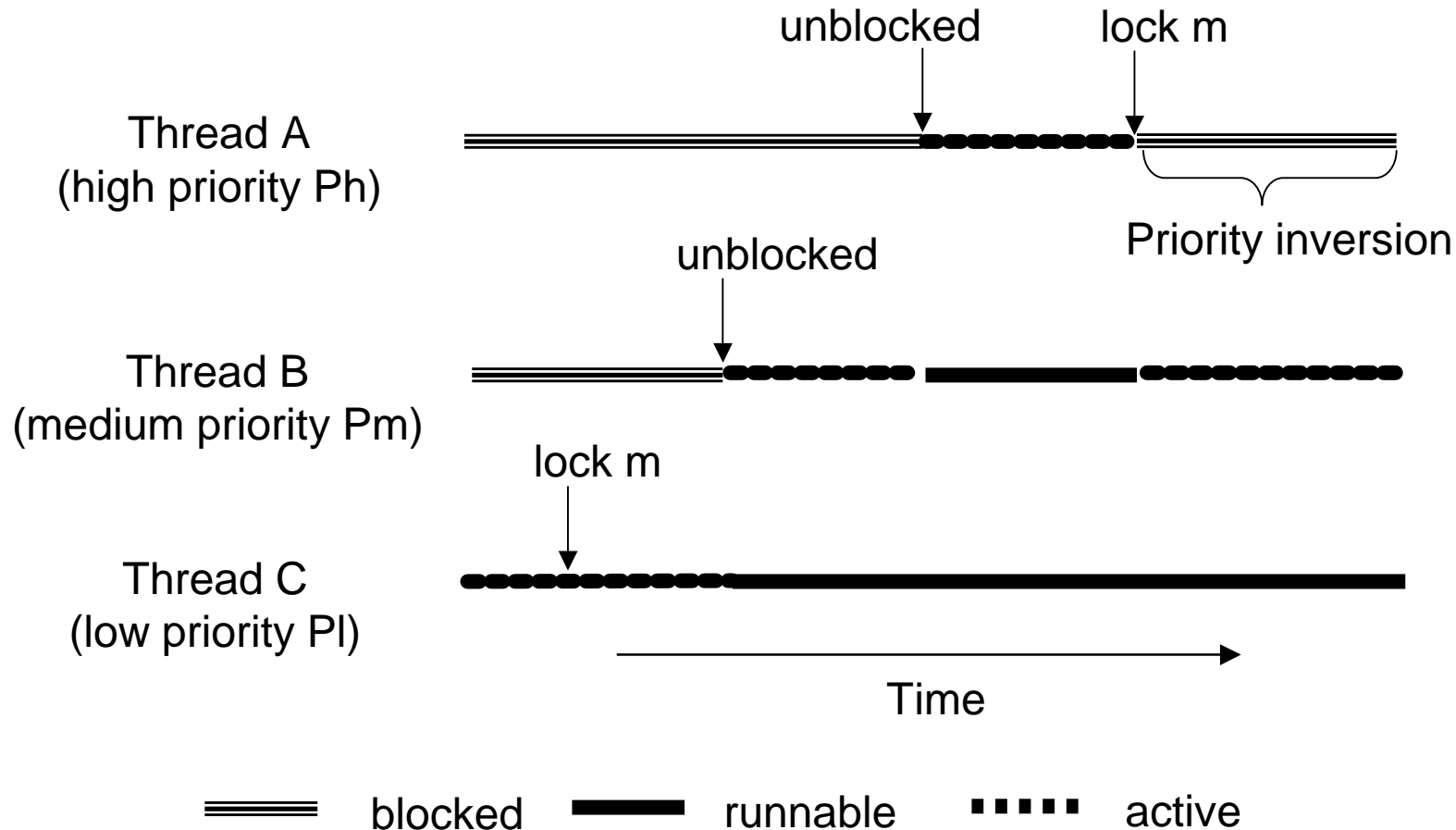
- ü Most versions of UNIX are forced to wait either for a system call to complete or for an I/O block to take place before doing a context switch

n Preempting system calls

- ü Insert preemption points
 - § Still dispatch latency can be large
- ü Make the entire kernel preemptible.
 - § All kernel data structures must be protected
 - § Solaris 2

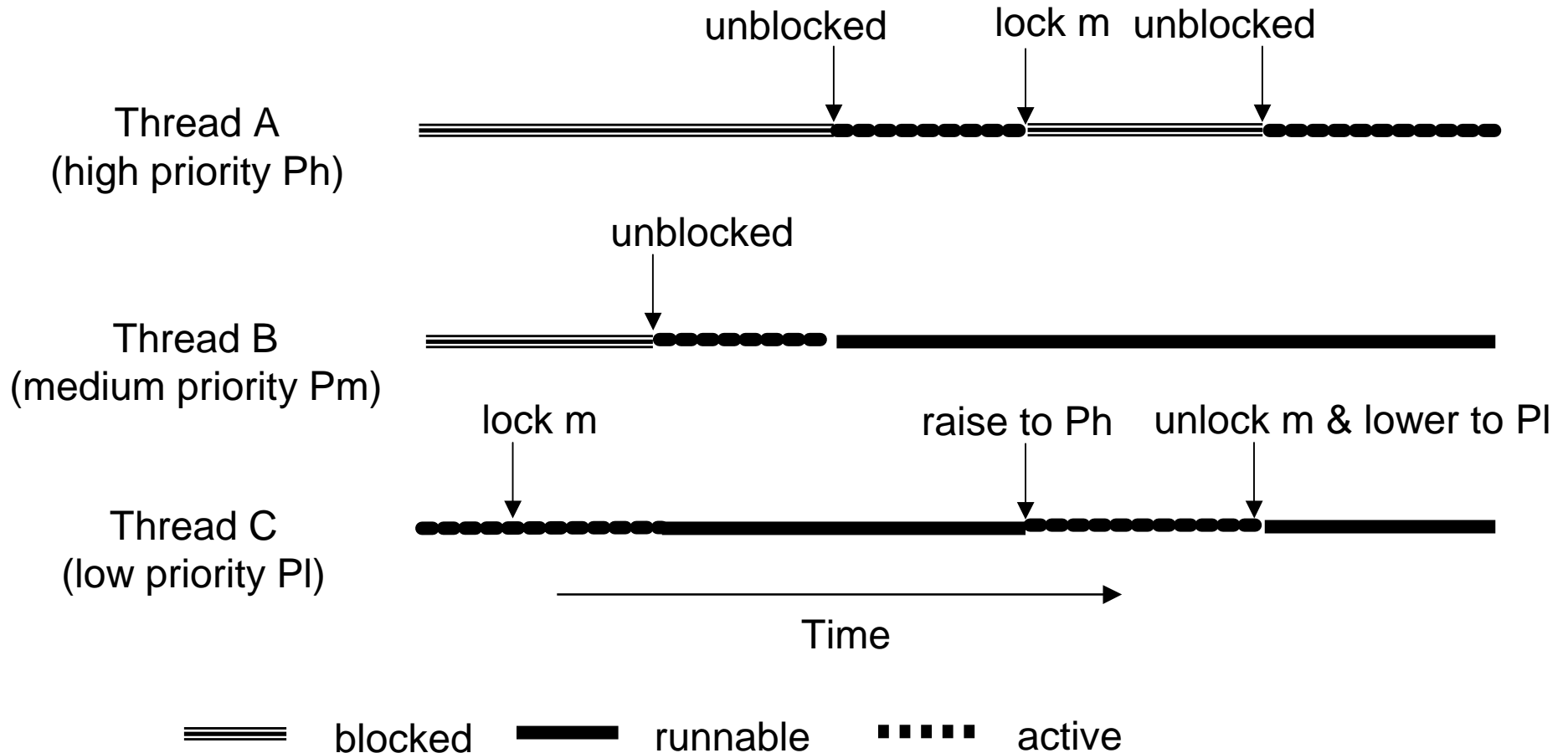
Real-Time Scheduling (Cont'd)

n Priority inversion problem



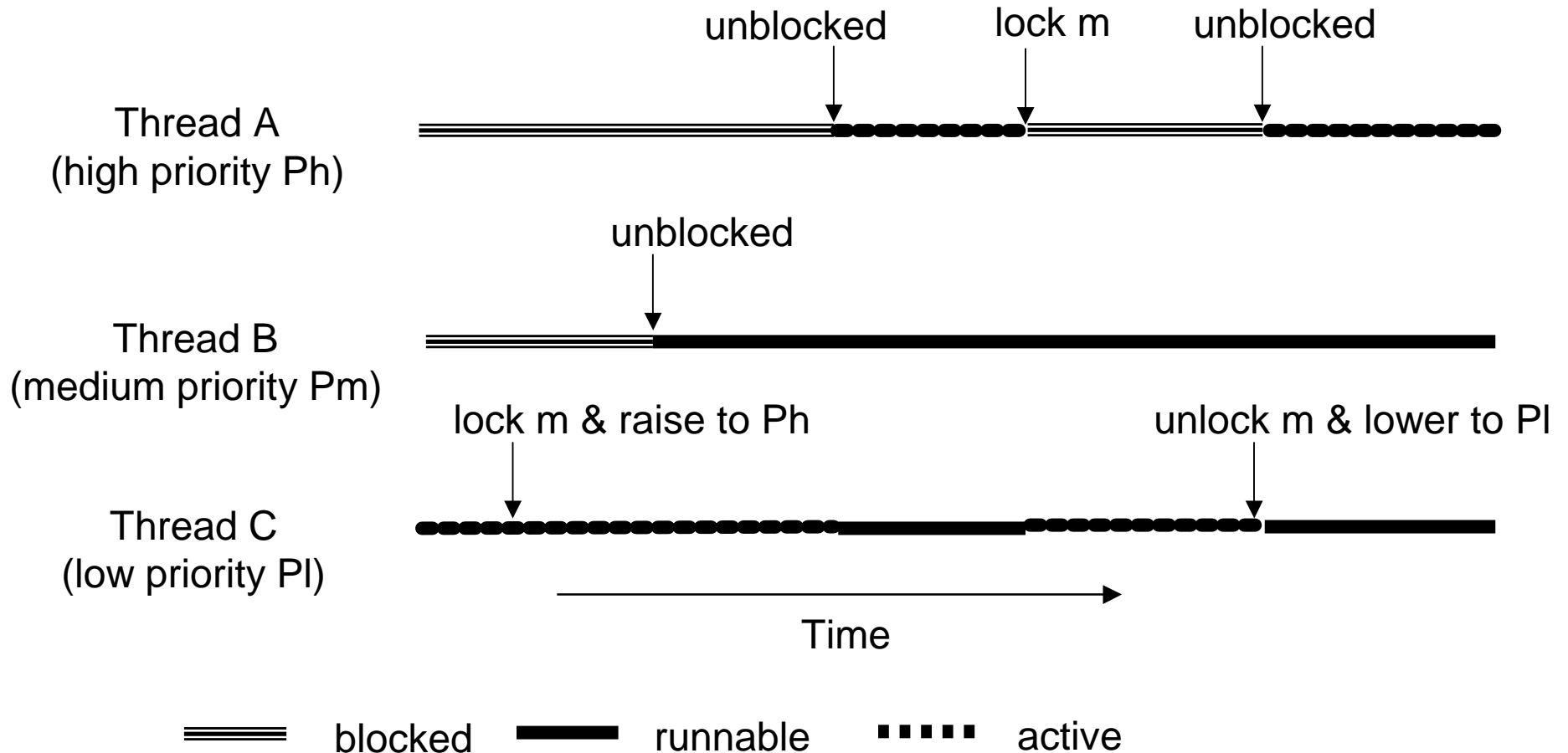
Real-Time Scheduling (Cont'd)

n Priority inheritance protocol



Real-Time Scheduling (Cont'd)

n Priority ceiling protocol



Algorithm Evaluation

n Deterministic modeling

- ü Takes a particular predetermined workload and defines the performance of each algorithm for that workload

n Queueing models

- ü Mathematical models used to compute expected system parameters

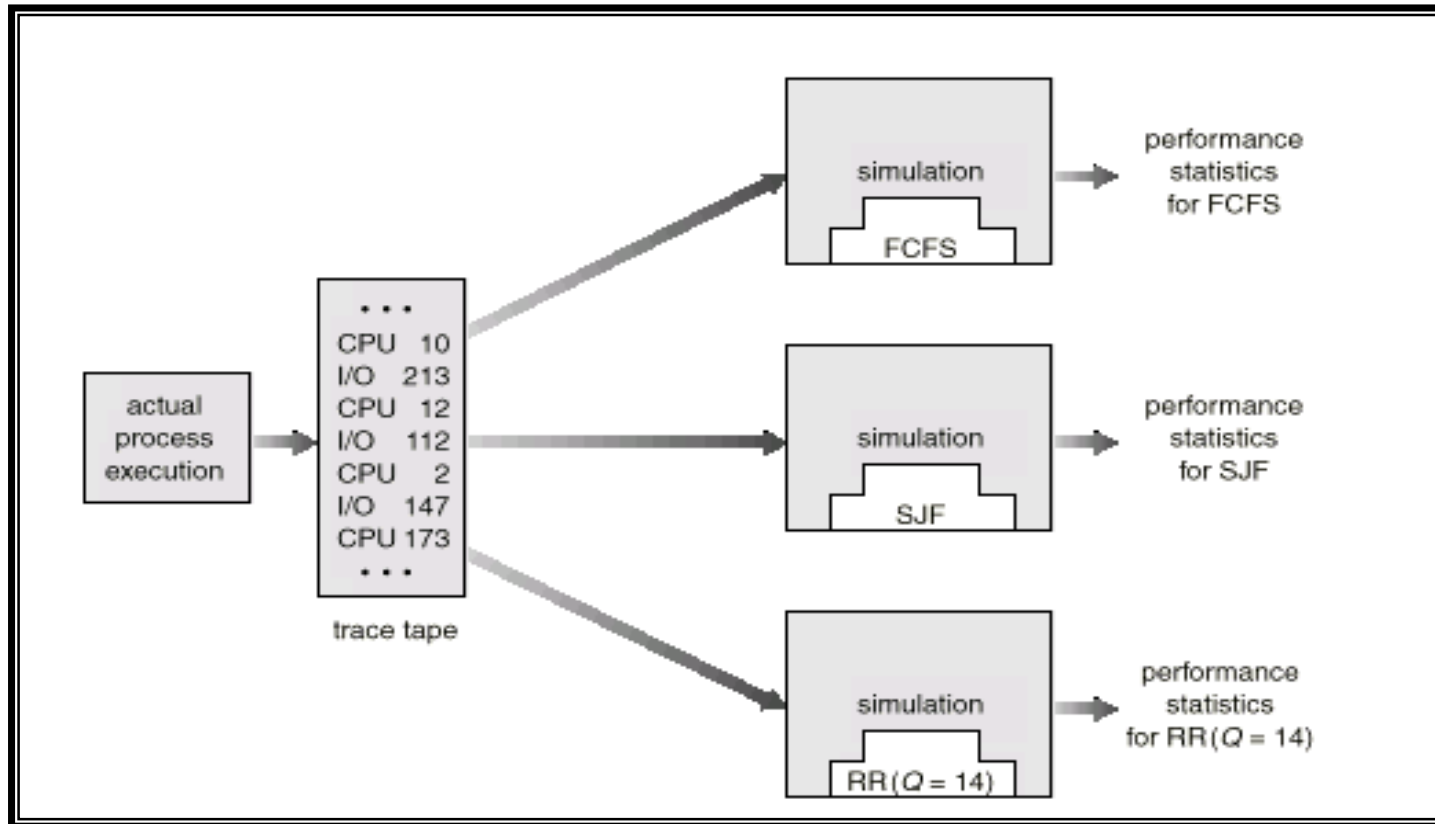
n Simulation

- ü Algorithmic models which simulate a simplified version of a system using statistical input
- ü Trace tape (or trace data)
- ü Cf) *Emulation*

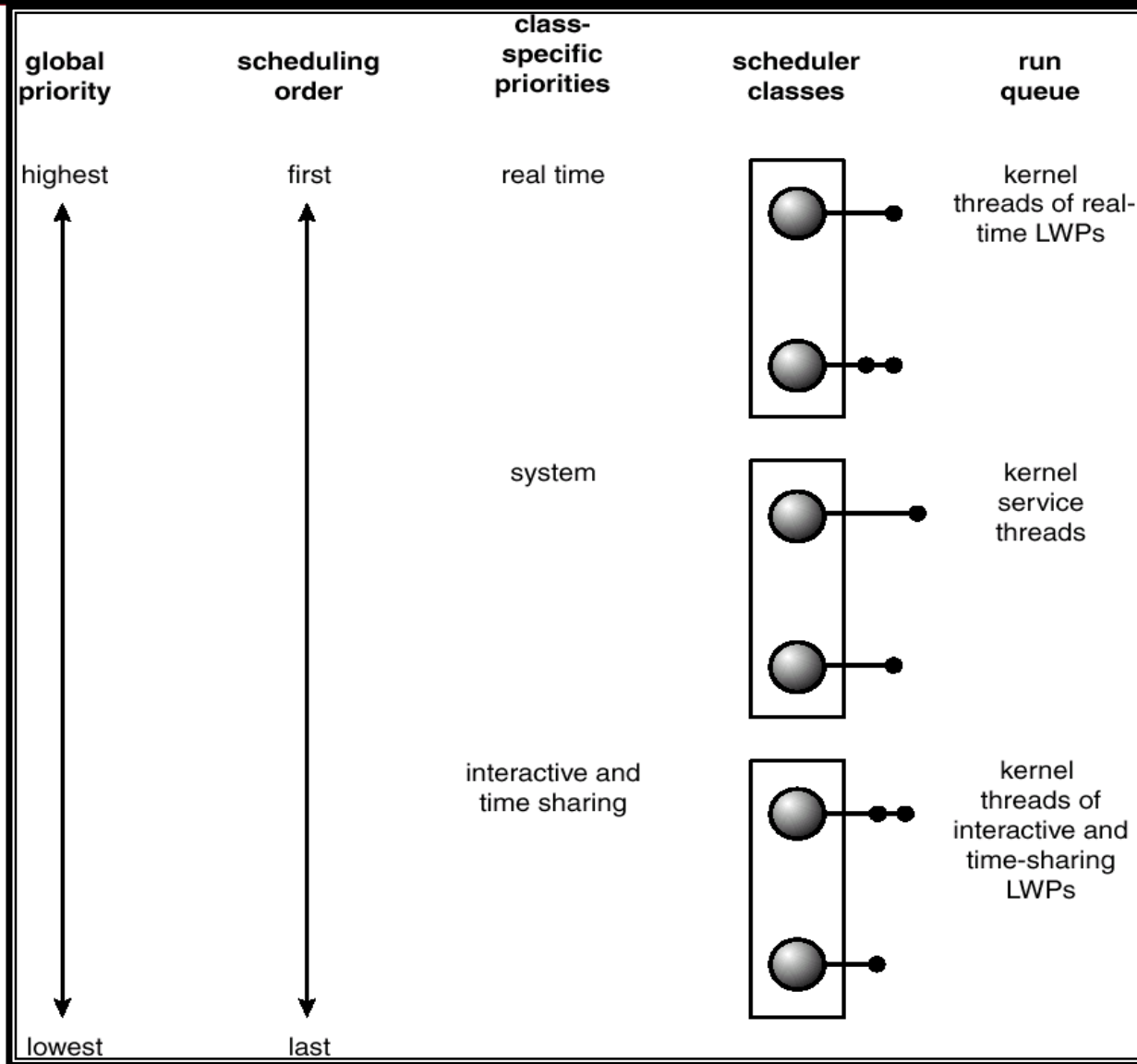
n Implementation

- ü Direct implementation of the system under test, with appropriate benchmarks

Evaluation of CPU Schedulers by Simulation



Solaris 2 Scheduling



Windows 2000 Priorities

	real-time	high	above normal	normal	below normal	idle priority
time-critical	31	15	15	15	15	15
highest	26	15	12	10	8	6
above normal	25	14	11	9	7	5
normal	24	13	10	8	6	4
below normal	23	12	9	7	5	3
lowest	22	11	8	6	4	2
idle	16	1	1	1	1	1