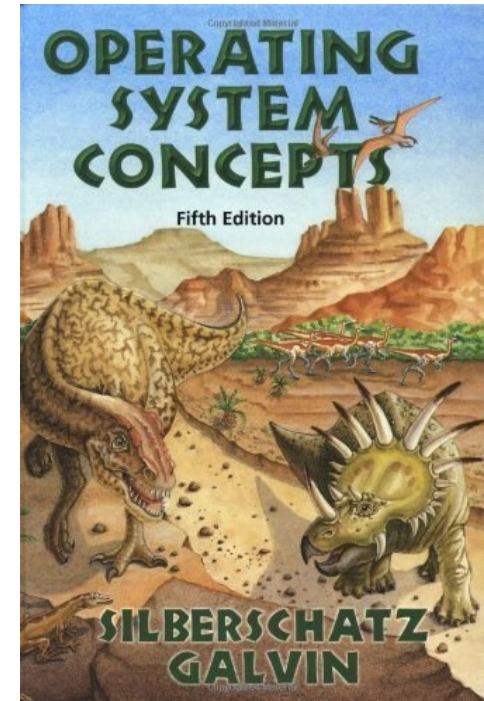


Chapter 5: CPU Scheduling

School of Computing, Gachon Univ.
Joon Yoo



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

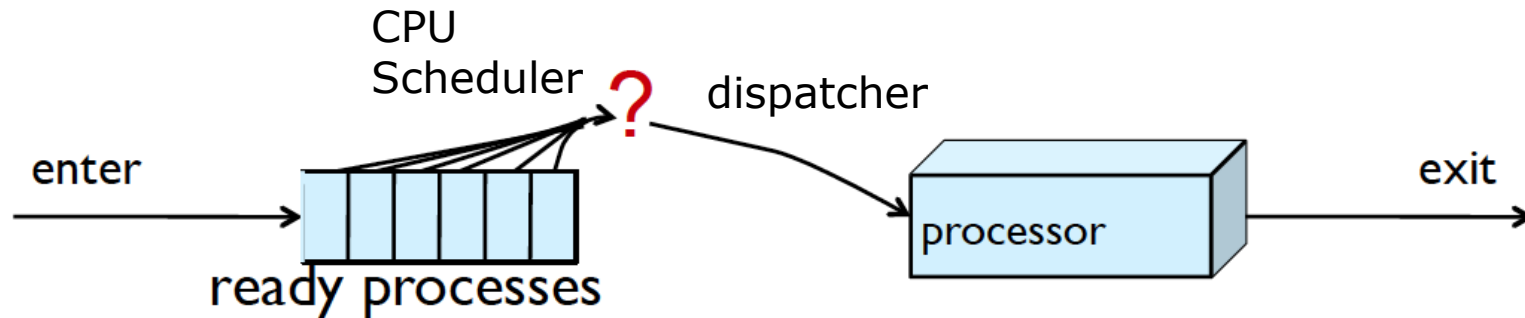
Chapter 5: CPU Scheduling

- Basic Concepts
- Scheduling Criteria
- Scheduling Algorithms
- Advanced Scheduling
 - Multiple-Processor Scheduling
 - Real-Time CPU Scheduling

Introduction

- Multitasking: multiple processes (threads) in the system with one (or more) processor cores
- The role of OS
 - Increases CPU utilization: by organizing processes (threads) so that the CPU always has one to execute
 - Timesharing: by frequent context switching
 - CPU scheduling
 - ▶ Allocate CPU time slots to processes (threads) in Ready queue
 - **Note:** Terms “**Process/Thread/CPU/Job/Task** scheduling” will be used inter-changeably in this chapter.

CPU Scheduler

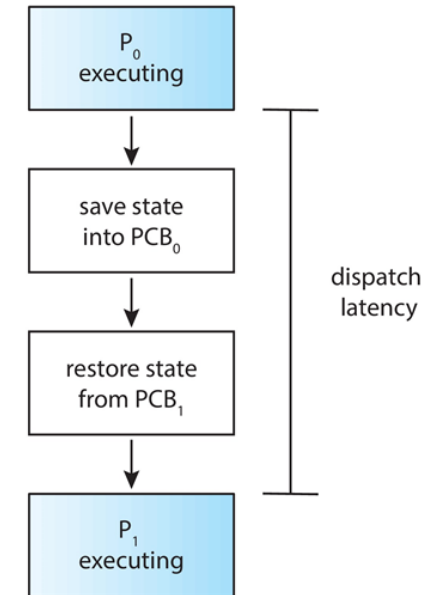


■ CPU Scheduler

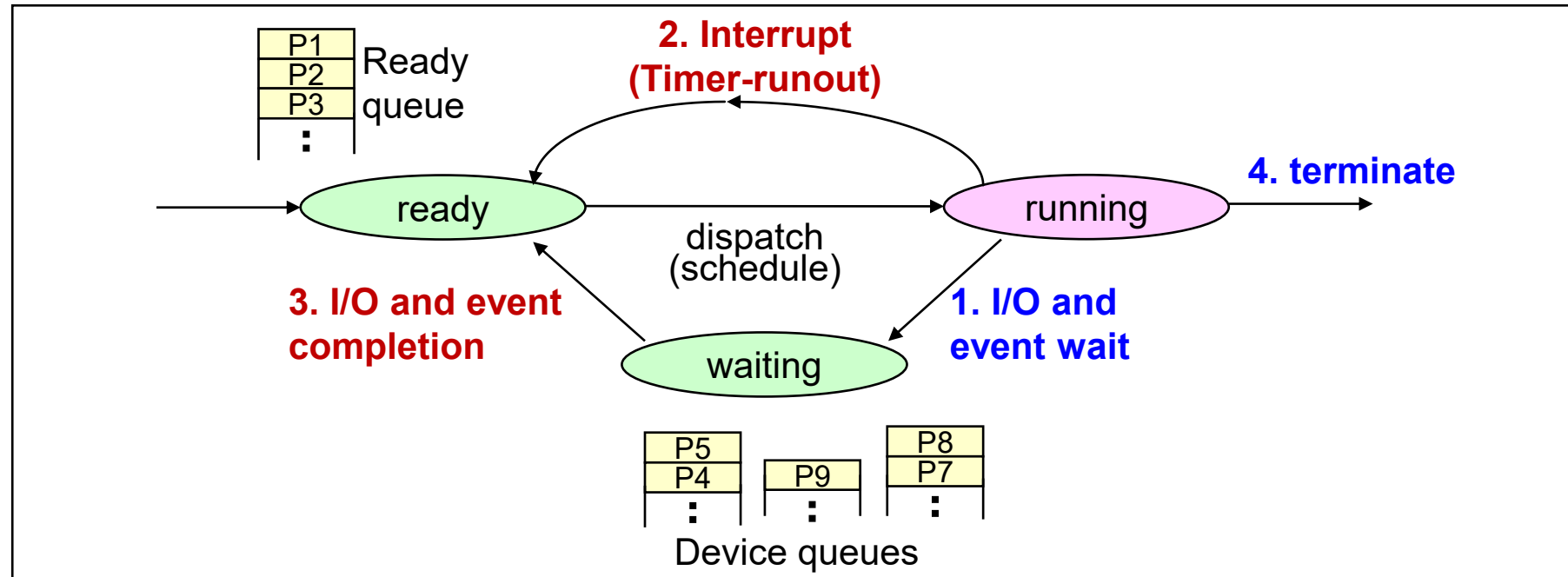
- An OS component that performs the **CPU Scheduling**
 - ▶ Which process (or thread) should I run next?
- Allocates CPU time slots to processes (or threads)

Dispatcher

- **Dispatcher**
 - **CPU Scheduler module** that gives control of the CPU to the process selected by the scheduler
 - ▶ **context switching**
 - ▶ switching to **user mode**
 - ▶ jumping to the proper location (**PC**) in the user program to restart that program

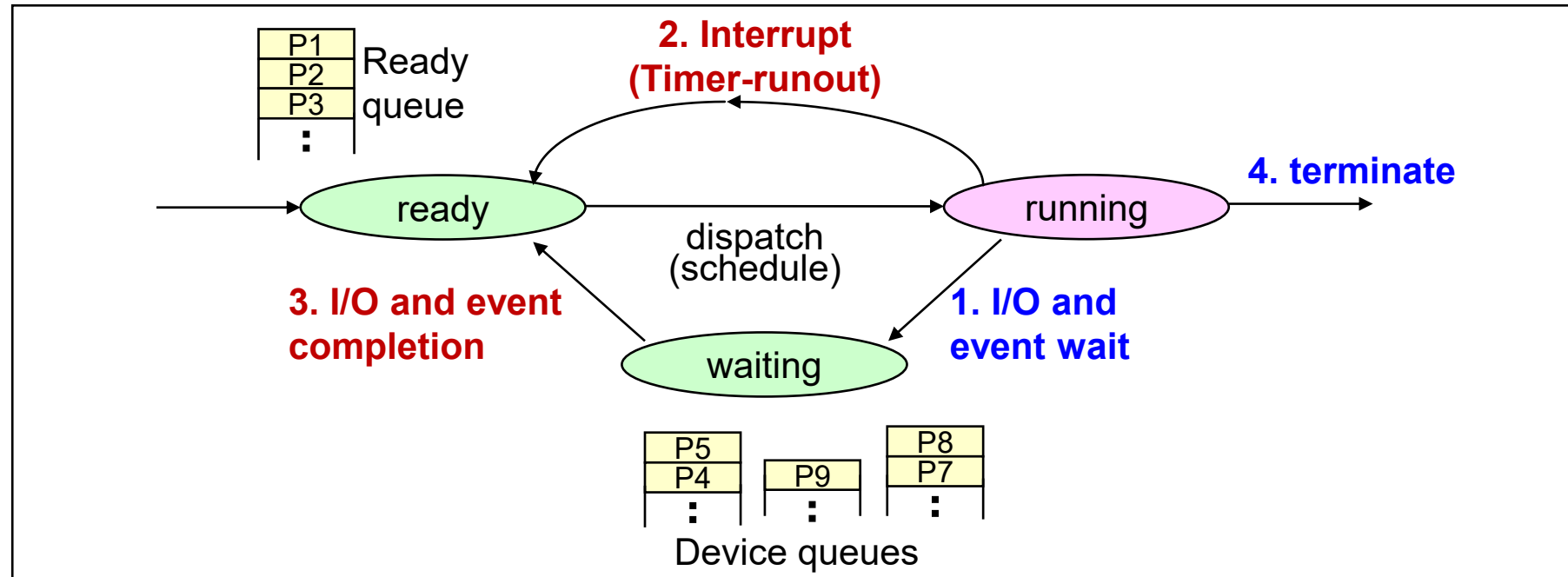


Re: Process and State Transition



- CPU scheduling decisions may take place under above four circumstances

Re: Process and State Transition



- Three states : ready, running, and waiting
- CPU scheduling decisions may take place when a process:

Must dispatch a new process

1.
2.
3.
4.

Switches from running to waiting state

Switches from running to ready state

Switches from waiting to ready

Terminates

Can make a decision!
- should scheduler dispatch a new process or continue with old one?

Preemptive (vs. Non-preemptive)

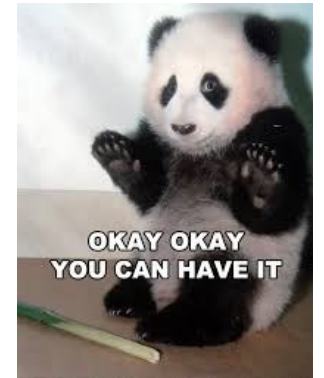
■ Non-preemptive

- Once CPU is allocated to a process, it holds it till it
 - ▶ It gives up by itself (스스로): Terminates (case 4) or blocks itself to wait for I/O – (case 1)
 - ▶ It does not give up for case 2 & 3
- Process uses the CPU until it voluntarily releases it - **No preemption**



■ Preemptive

- Currently running process may be **interrupted** and moved to the ready state by the OS
- CPU may be **preempted** to another process independent of the intention of the running process
 - ▶ Most modern OSes: Windows, Mac OS, Linux, ...



Preemptive vs. Non-preemptive

■ Non-preemptive scheduling

- Pros: Low context switch overhead
- Cons : May result in longer mean response time (not good for for time-sharing systems and real-time systems)

■ Preemptive scheduling

- Cons: race conditions (need synchronization – Ch. 6-7)

Preemptive vs. Non-preemptive

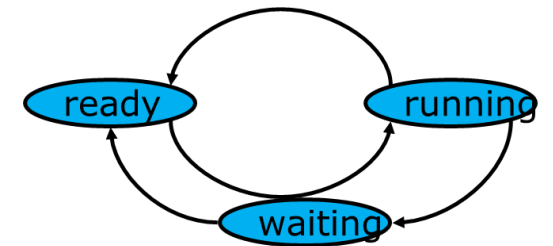
- Which scheduling policy
 - can immediately take care of Interrupts?
 - has less context-switching overhead?
 - has less response time?
 - can use time-sharing?
 - do you think is used more generally?

Chapter 5: CPU Scheduling

- Basic Concepts
- **Scheduling Criteria**
- Scheduling Algorithms
- Advanced Scheduling
 - Multiple-Processor Scheduling
 - Real-Time CPU Scheduling

Scheduling Criteria

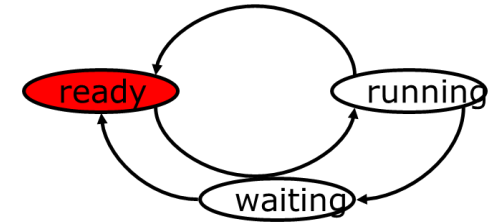
- **CPU utilization** (%)
 - Higher is better (0~100%)
 - In a real system, 40% (lighted loaded) to 90% (heavily loaded)
- **Throughput** (# of processes completed / time)
 - Higher is better
 - Long process (1 process/several seconds), short process (tens of process/second)
- **Turnaround time** (s): time for each process to complete
 - Lower is better
 - Time spent in the ready queue, executing CPU, and waiting (for I/O, ...)



Scheduling Criteria

- **Waiting time** (s): sum of time spent waiting in the ready queue

- Lower is better
- CPU-scheduling algorithm affects only the time spent on the ready queue (compare with turnaround time)



- **Response time** (s): time from request to first response
 - Lower is better
 - Used for Interactive and real-time systems
 - e.g., HTTP request to Web page load response

Chapter 5: CPU Scheduling

- Basic Concepts
- Scheduling Criteria
- **Scheduling Algorithms**
- Advanced Scheduling
 - Multiple-Processor Scheduling
 - Real-Time CPU Scheduling

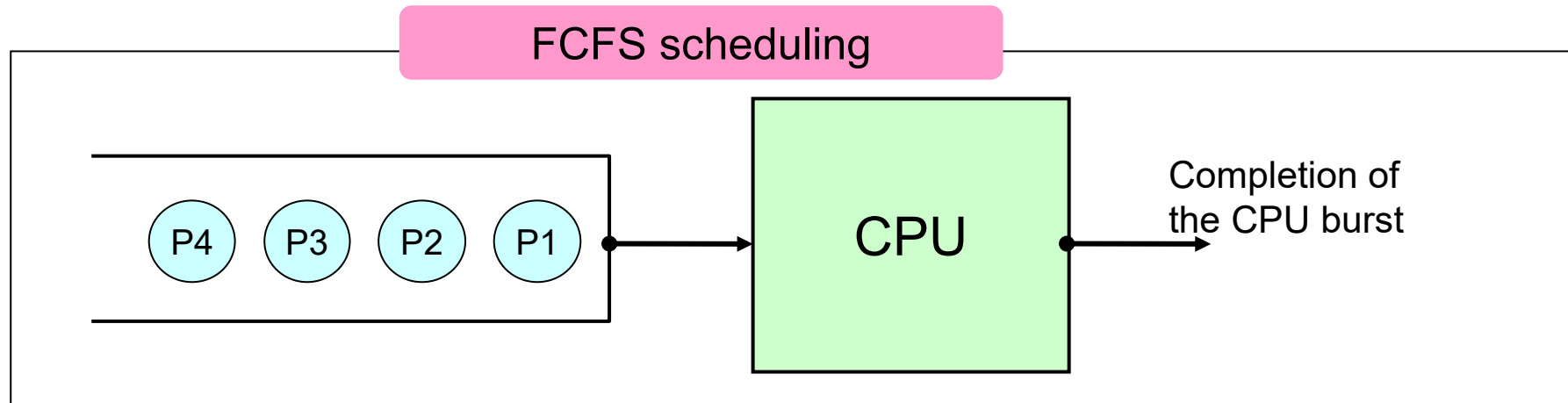
Scheduling Algorithms

- CPU Scheduling Algorithm
 - Algorithm that decides which of the processes in the ready queue is to be allocated the CPU

- Various algorithms
 - First-Come, First-Served (FCFS) Scheduling
 - Shortest-Job-First (SJF) Scheduling
 - Round-Robin Scheduling
 - Priority Scheduling
 - Earliest Deadline First Scheduling

First-Come, First-Served (FCFS)

- FCFS(First-Come-First-Served) scheduling
 - Simplest scheme using FIFO queue



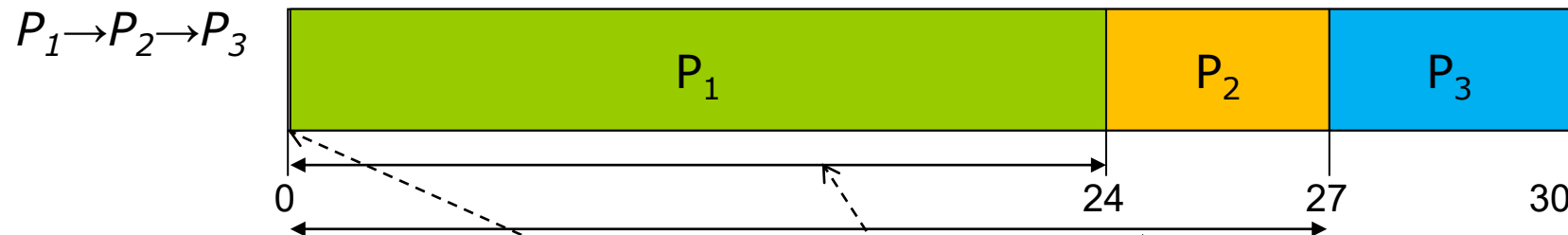
- **Non-preemptive** scheduling

First-Come, First-Served (FCFS) Scheduling

All arrived at
time 0 in order
 P_1, P_2, P_3

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

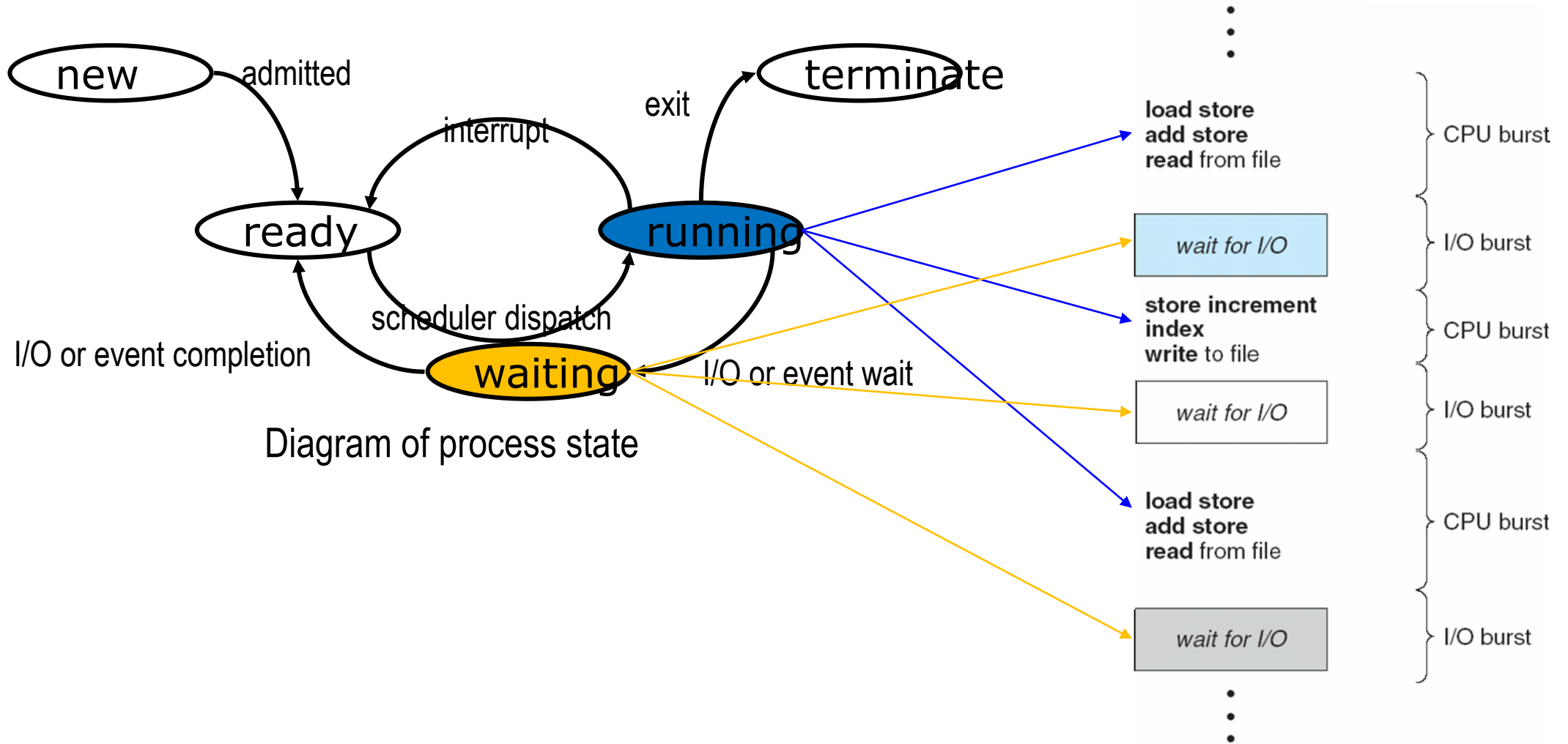
- Suppose that the processes arrive in the order: P_1, P_2, P_3
The **Gantt Chart** for the schedule is:



- Waiting time** for $P_1 = 0$; $P_2 = 24$; $P_3 = 27$
 - Average waiting time** (WT): $(0 + 24 + 27)/3 = \underline{17}$
 - Turnaround time (TT): $(24 + 27 + 30)/3 = 27$

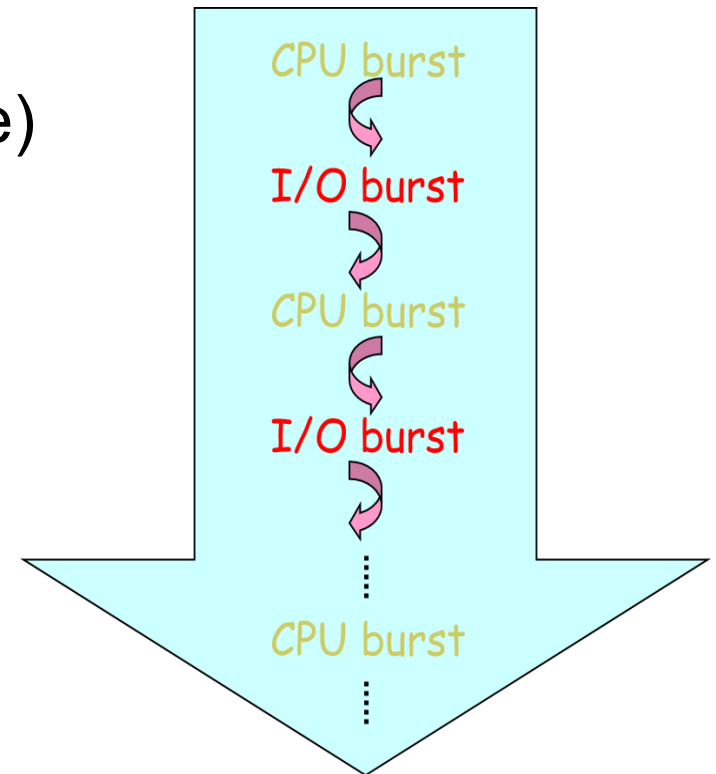
Can we do better?

CPU & I/O Bursts

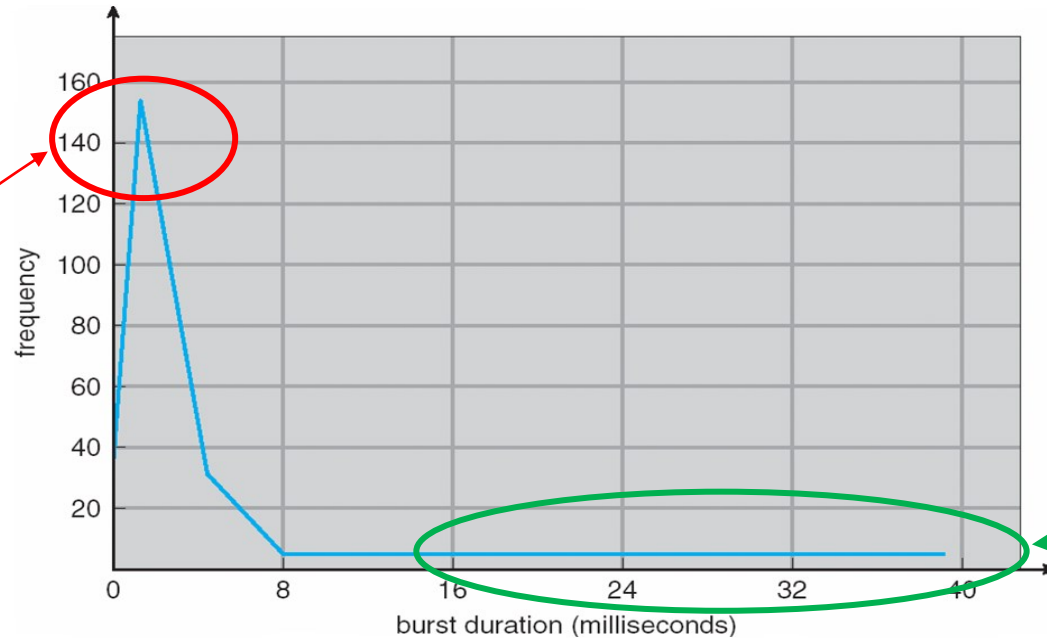


CPU & I/O Bursts

- Process execution consists of cycles of CPU bursts and I/O bursts
 - **CPU burst:** CPU executions (running state)
 - ▶ e.g., sorting a million-entry array in RAM
 - **I/O burst:** I/O wait (waiting state)
 - ▶ e.g., disk read/write, networking, printing
- **CPU burst time** is an important factor(criteria) for scheduling algorithms



Histogram of CPU-burst Durations



CPU bursts distribution

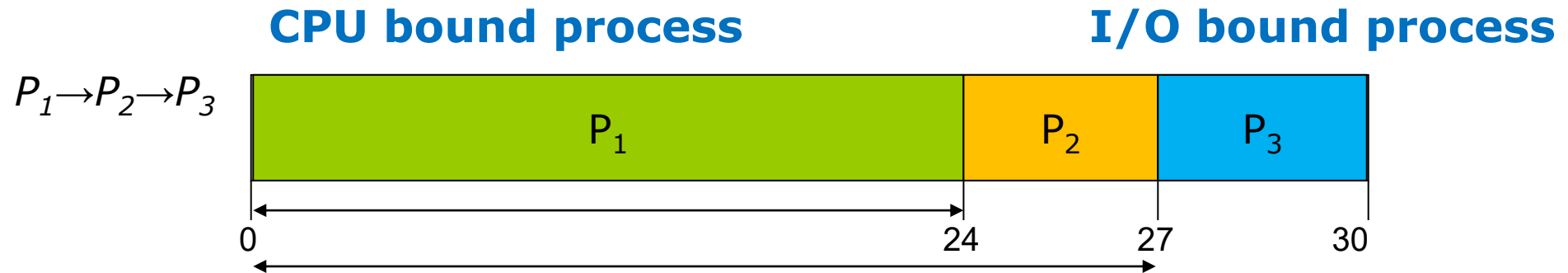
- Large number of short CPU bursts and small number of long CPU bursts
I/O bound process **CPU bound process**

FCFS: Convoy Effect



■ Convoy Effect

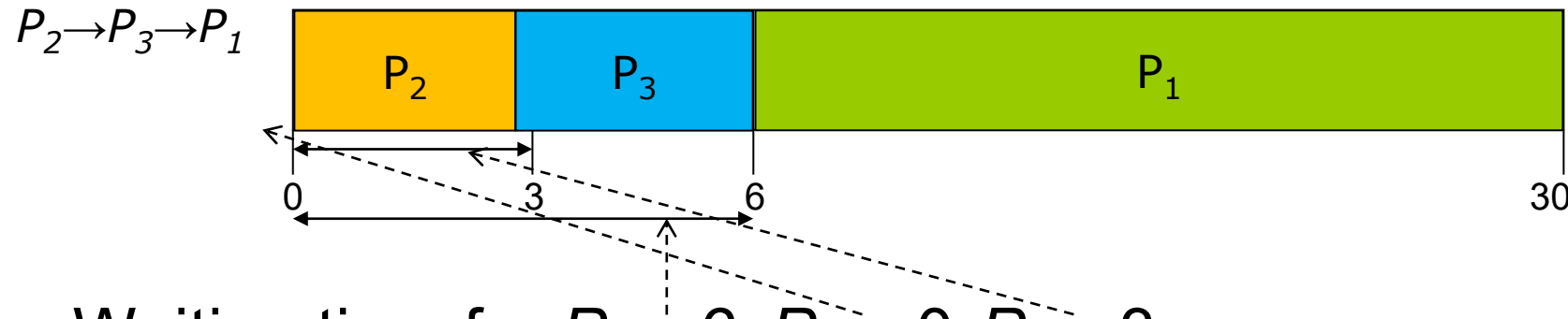
- Larger waiting time for I/O bound process
 - ▶ CPU-bound jobs will hold CPU until exit or I/O



FCFS Scheduling

Suppose that the processes arrive in the order: P_2 , P_3 , P_1

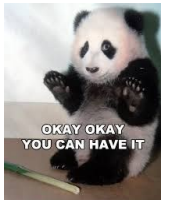
- The Gantt chart for the schedule is:



- Waiting time for $P_1 = 6$; $P_2 = 0$; $P_3 = 3$
 - Average waiting time (WT): $(6 + 0 + 3)/3 = \underline{3}$
 - Turnaround time (TT): $(30 + 3 + 6)/3 = 13$
- Lesson: scheduling algorithm can reduce WT and TT

Shortest-Job-First (SJF) Scheduling

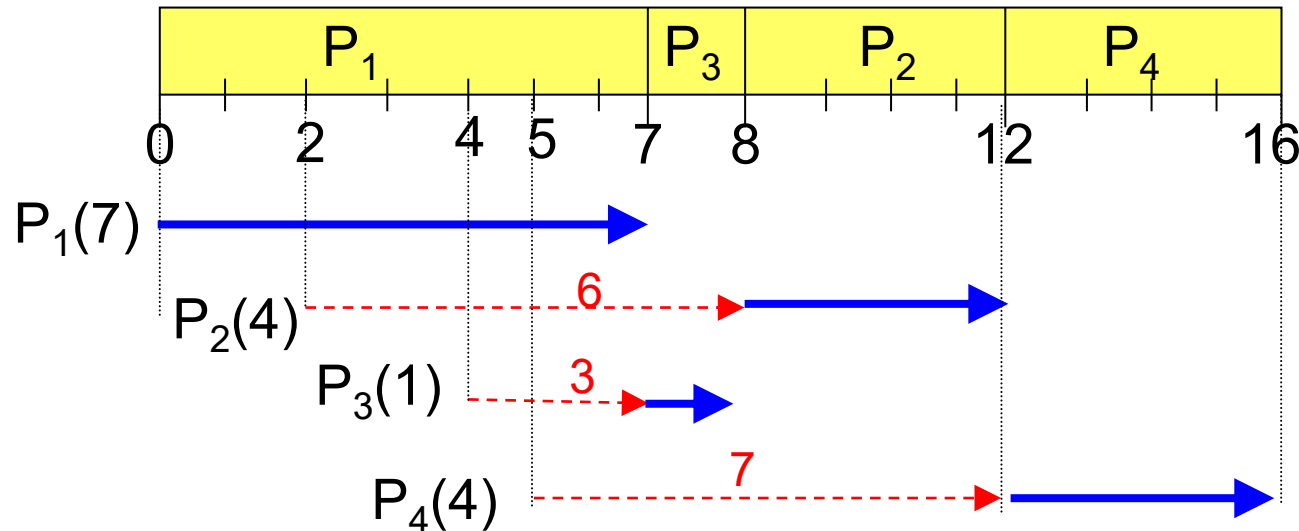
- Schedule the process with the shortest time
- Gives minimum average waiting time
- Two schemes:
 - **Nonpreemptive SJF**
 - ▶ Once CPU given to the process it cannot be preempted until completes its CPU burst
 - **Preemptive SJF**
 - ▶ If a new process arrives with CPU burst length less than remaining time of current executing process, preempt.
 - ▶ This scheme is known as the **Shortest-Remaining-Time-First (SRTF)**



Example of Non-Preemptive SJF

<u>Process</u>	<u>Arrival Time</u>	<u>Burst Time</u>
P_1	0	7
P_2	2	4
P_3	4	1
P_4	5	4

- SJF (non-preemptive)



P_1 's waiting time = 0

P_2 's waiting time = 6

P_3 's waiting time = 3

P_4 's waiting time = 7

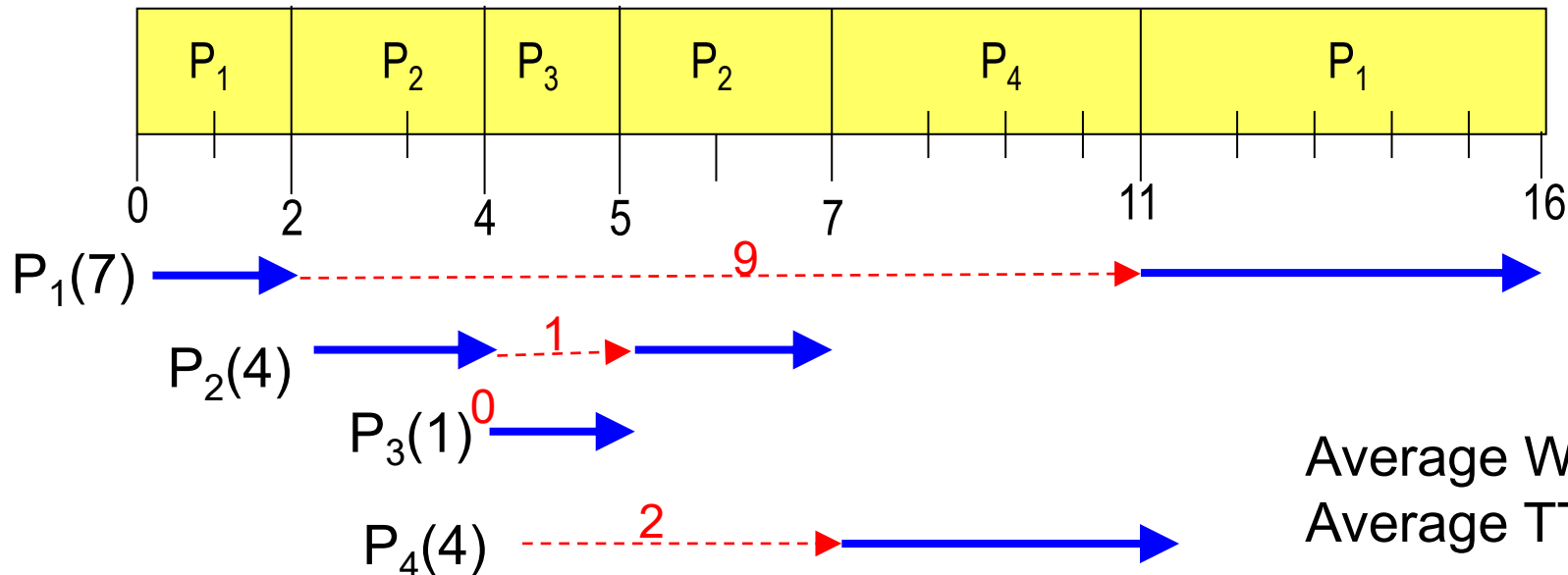
$$\text{Average WT} = (0 + 6 + 3 + 7)/4 = 4$$

$$\text{Average TT} = (7 + 10 + 4 + 11)/4 = 8$$

Example of Preemptive SJF (SRTF)

Process	Arrival Time	Burst Time
P_1	0	7
P_2	2	4
P_3	4	1
P_4	5	4

- SRTF (preemptive SJF)**



P_1 's wating time = 9

P_2 's wating time = 1

P_3 's wating time = 0

P_4 's wating time = 2

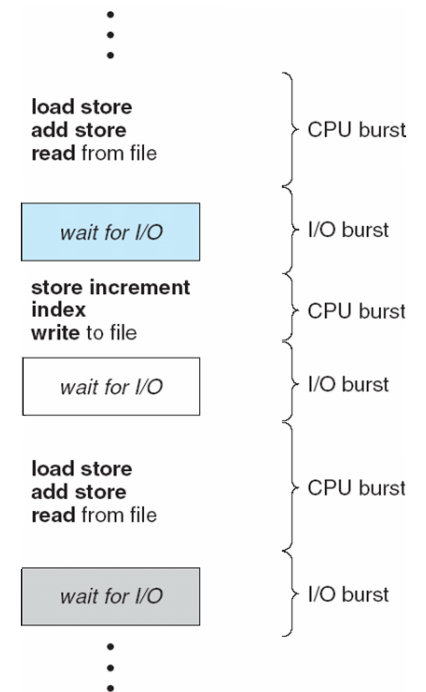
$$\text{Average WT} = (9 + 1 + 0 + 2)/4 = 3$$

$$\text{Average TT} = (16 + 5 + 1 + 6)/4 = 7$$

CPU Burst Prediction

- How do we know the next CPU burst length?
 - We cannot know for sure
 - But we may **predict** the CPU burst length
 - ▶ How? – “History is a mirror to the future”
 - ▶ The next CPU burst is generally predicted as an **exponential moving average** of the measured lengths of previous CPU bursts:
 - From n CPU bursts from the past: (t_1, t_2, \dots, t_n)
 - Predict the future $(n+1)^{\text{th}}$ CPU burst: (τ_{n+1})

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



CPU Burst Prediction

- Predicted value for the next CPU burst using *exponential moving average*

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

- ▶ τ_{n+1} : predicted value for the next CPU burst.
- ▶ t_n : length of the nth CPU burst, $0 \leq \alpha \leq 1$,

- $0 \leq \alpha \leq 1$ controls the weight of recent and past

- ▶ What does it mean

- if $\alpha = 0$? : $\tau_{n+1} = \tau_n$
- if $\alpha = 1$? : $\tau_{n+1} = t_n$

- If we expand the above equation,

CPU bursts in time order: $(t_1, t_2, \dots, t_{n-1}, t_n)$

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

$$\dots = \alpha t_n + \alpha(1 - \alpha)t_{n-1} + \alpha(1 - \alpha)^2 t_{n-2} + \alpha(1 - \alpha)^3 t_{n-3} + \dots$$

CPU Burst Prediction

- Example

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

P3's previous three CPU burst samples are: 8, 16, 8 (in time order; 시간순서). What should be P3's next predicted CPU burst value of P3? Assume we use *exponential moving average* with $\alpha=0.5$ and initial average $\tau_0=8$. [OS midterm exam'19]

$$\tau_1 = 0.5 \times t_0 + 0.5 \times \tau_0 = 8$$

$$\tau_2 = 0.5 \times t_1 + 0.5 \times \tau_1 = 12$$

$$\tau_3 = 0.5 \times t_2 + 0.5 \times \tau_2 = 10$$

P3's CPU burst = 10

SJF Scheduling (Cont.)

■ Pros

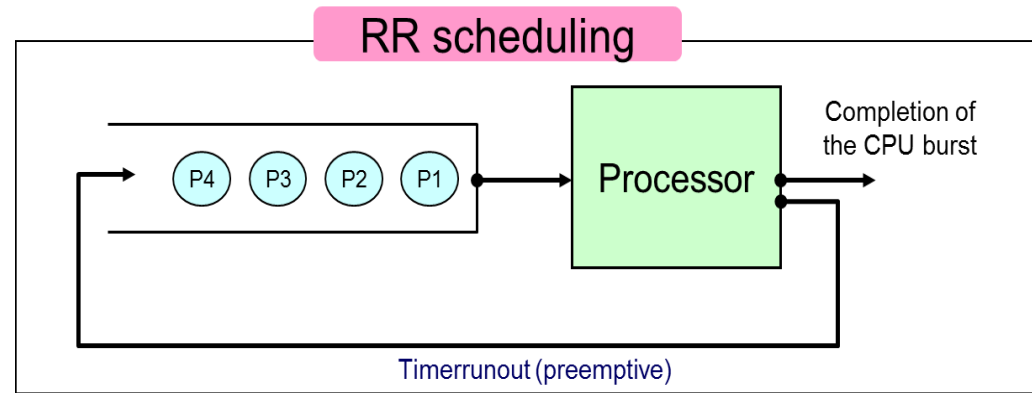
- Gives minimum average waiting time for a given set of processes
- Minimizes the number of processes in the system
 - ▶ Less average person waits in line (=waiting time), the shorter the line (=number of processes)

■ Cons

- **CPU burst time** can only be predicted
 - ▶ We can only approximate SJF scheduling
 - ▶ Thus, may not be optimal if prediction is inaccurate.

Round Robin (RR)

■ RR (Round-Robin) scheduling



■ User Timer

- Preempt job after **time quantum (time slice)**
- When preempted move back to Ready queue



RR (Round-Robin) scheduling

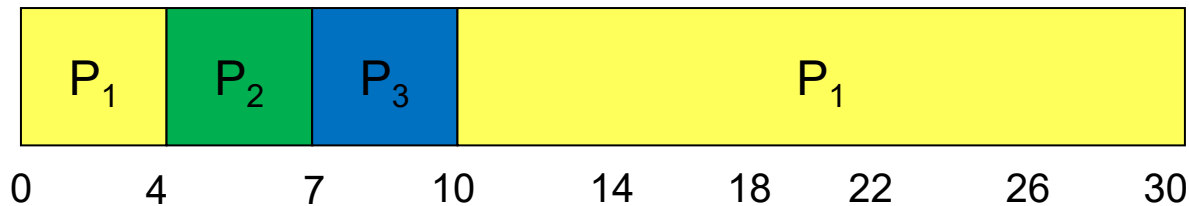
- RR (Round-Robin) scheduling
 - **Time quantum** (=time slice) for each process
 - ▶ System parameter (generally from 10 to 100 ms in length)
 - The CPU scheduler sets a **timer** to 1 **time quantum**
 - ▶ An interrupt will occur after 1 time quantum
 - ▶ The (running) process that has exhausted his time quantum releases the CPU and goes to the ready state (timerrunout)
 - ▶ A context switch occurs!

Example of RR with Time Quantum = 4

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

convoy effect case

- The Gantt chart for RR is:



- Typically, higher average **waiting time** than SJF, but fair and no starvation

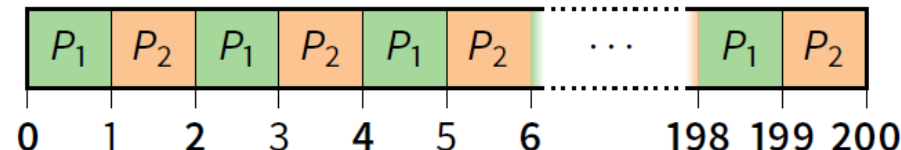
RR (Round-Robin) scheduling

■ Pros

- Low **response time** for all jobs
- Low average **waiting time** for I/O bound jobs
- No **starvation** for CPU bound jobs

■ Cons

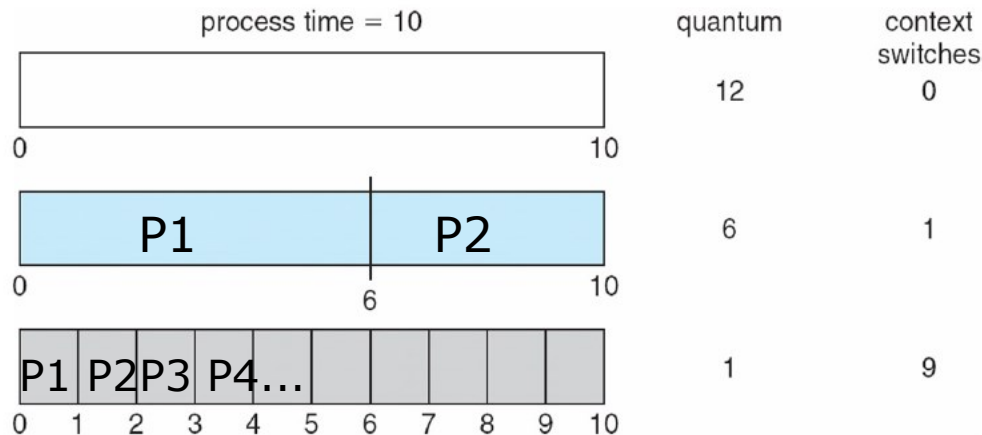
- Bad for *same-sized* jobs
 - ▶ Assume 2 jobs of time=100 each:
 - What would average turnaround time be with RR?
 - How does that compare to FCFS?
 - ▶ Even if context switches were free...
- Context switching overhead due to preemptions



Round Robin Scheduling

■ Performance

- quantum (q) extremely large = FCFS
- q small \Rightarrow **time-sharing** $\Rightarrow q$ must be large with respect to **context switch time**, otherwise overhead is too high.



➤ Number of Context Switching

- quantum = 12 \rightarrow FCFS
- quantum = 6 \rightarrow 2 quanta, 1 Context Switching
- quantum = 1 \rightarrow 10 quanta, 9 Context Switching

In practice,
- Quantum: 10~100ms
- Context switch time:
<10us
(0.1~0.01% overhead)

Priority Scheduling

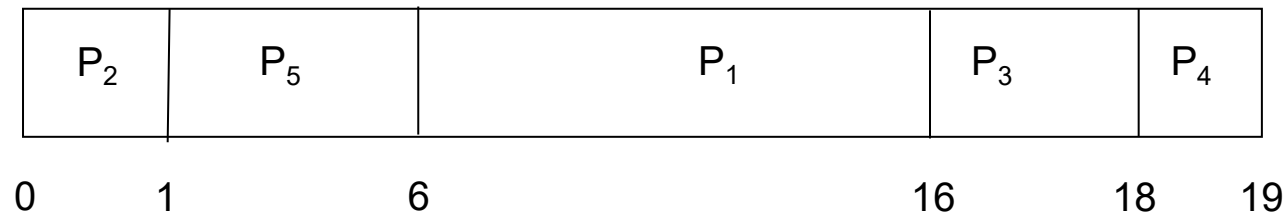
- A numeric **priority** is associated with each process
 - sometimes smaller number means higher priority (e.g., Unix/BSD)
 - or sometimes smaller number means lower priority (e.g., Pintos)
 - but you will be clearly given the definition in your problems
- The CPU is allocated to the process with the **highest priority**
 - Preemptive or Nonpreemptive
- Note **Shortest Job First (SJF)** is a **priority scheduling** where priority = predicted next CPU burst time
- Problem
 - **Starvation** – low priority processes may never execute
- Solution
 - **Aging** – as time progresses increase the priority of the process

Example of Priority Scheduling

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

Lower number is higher priority

Priority scheduling Gantt Chart



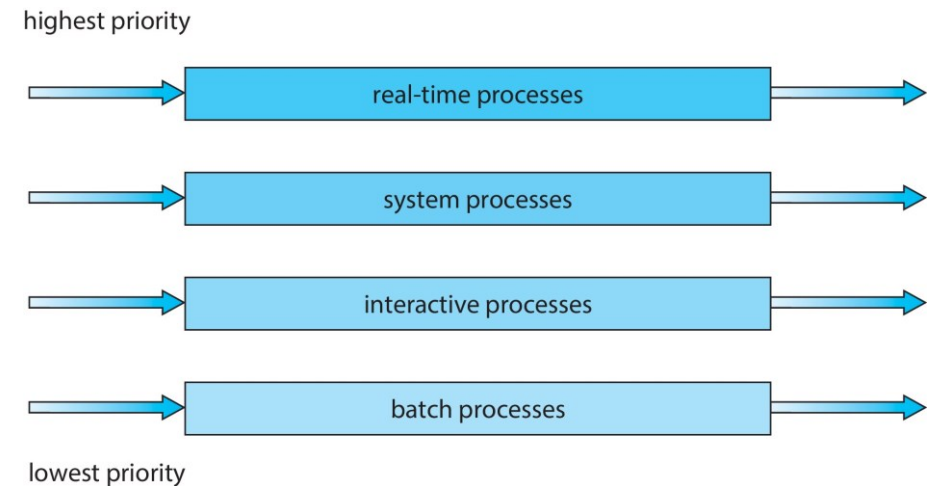
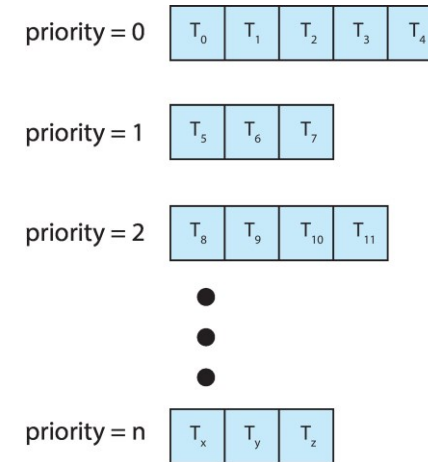
Average waiting time = 8.2 msec

- SJF : Special case of the general priority scheduling
- FCFS : Equal priority

Multilevel Queue

■ Multilevel Queue:

- Priority + round-robin scheduling
- Ready queue is partitioned into separate queues, eg:
- Kernel runs process on highest-priority non-empty queue
- Round-robins among processes on same queue



Midterm Exam'16

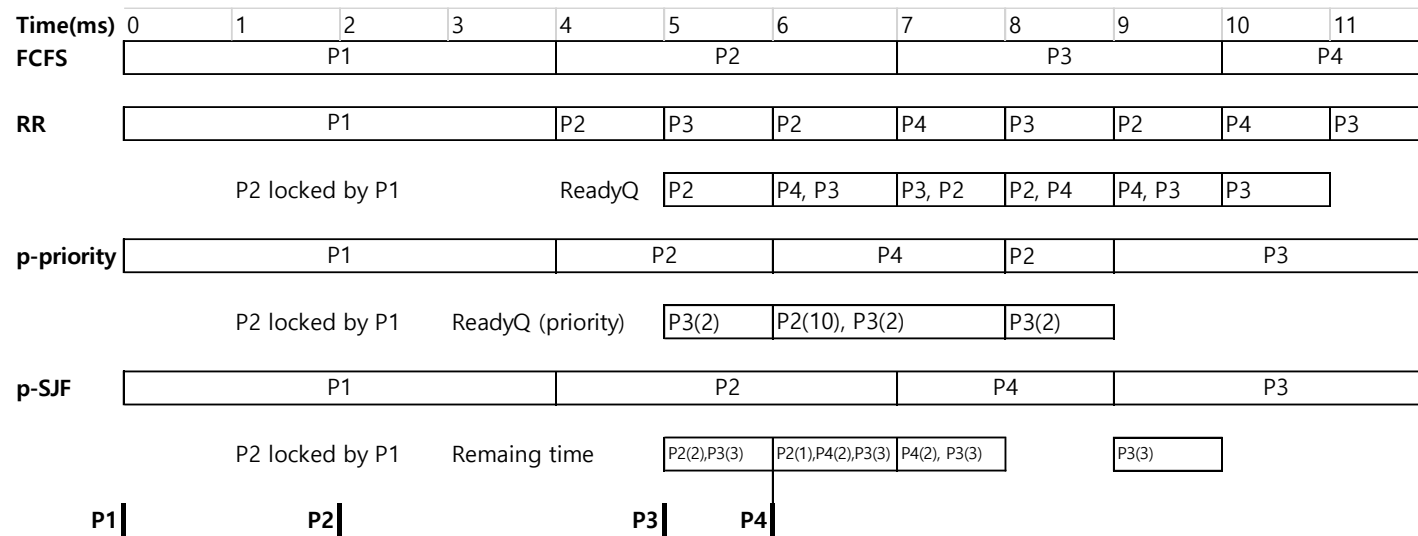
4. Process Scheduling: Consider the right processes, arrival times, waiting, priority, and CPU processing requirements. For each scheduling algorithm, draw the **GANTT chart** and compute the **average waiting time** for each scheduling algorithm (**FCFS**, **RR**, **preemptive-priority**, and **preemptive shortest -job-first (SJF)**)

Process	Arrival Time	Priority	Waits on Lock held by	CPU burst (ms)
P1	0	1	None	4
P2	2	10	P1	3
P3	5	2	None	3
P4	6	20	None	2

with the thread that is running on the CPU (for time-quantum (slice) based algorithms, assume a 1ms time quantum). [Total 2+5+5+5 = 17pts] **Notes:**

- For RR and Priority, assume that a newly arriving thread can run at the beginning of its arrival time, if the scheduling policy allows it. (스케줄링 규칙이 허용한다면, 도착 후 바로 실행 가능)
- A larger integer priority number indicates higher priority. If two threads have same priority, then FCFS is applied.
- If thread A tries to acquire a lock, but the lock is held by thread B, then thread A will be sent to WAITING state, and return to READY when thread A finishes the CPU burst. If thread A waits for a lock (formerly) held by thread B and B has already finished, it does not wait and acquires the lock immediately.
- You must show your work (계산과정) when computing each average waiting time.

Solution



FCFS: $(0+2+2+4)/4 = 2\text{ms}$

RR: $(0+5+4+3)/4 = 3\text{ms}$

p-priority: $(0+4+4+0)/4 = 2\text{ms}$

p-SJF: $(0+2+4+1)/4 = 1.75\text{ms}$

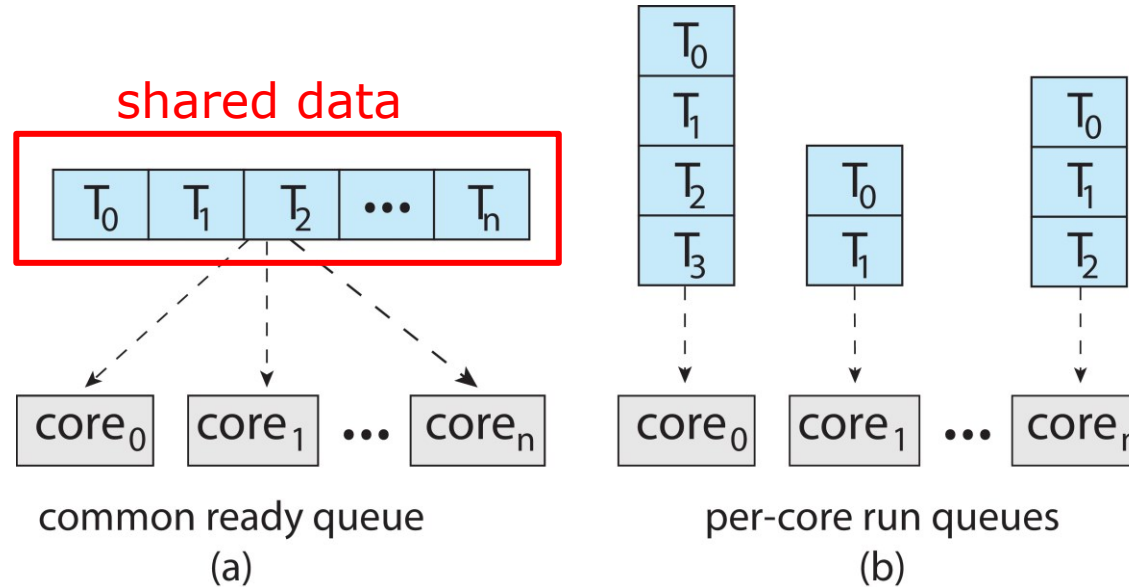
(1/4/4/4 pts for correct Gantt chart, 1pt for each correct computation/waiting time)

Chapter 5: CPU Scheduling

- Basic Concepts
- Scheduling Criteria
- Scheduling Algorithms
- **Advanced Scheduling**
 - **Multiple-Processor Scheduling**
 - Real-Time CPU Scheduling

Multiple-Processor Scheduling

- Symmetric multiprocessing (SMP) is where each processor is self scheduling.
 - a. All threads may be in a common ready queue
 - b. Each processor may have its own private queue of threads

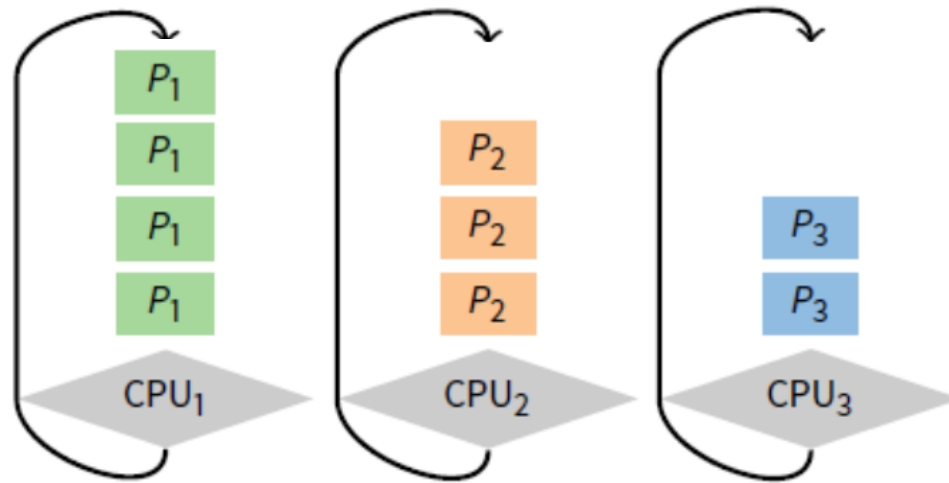


Race condition!

Unbalanced load!

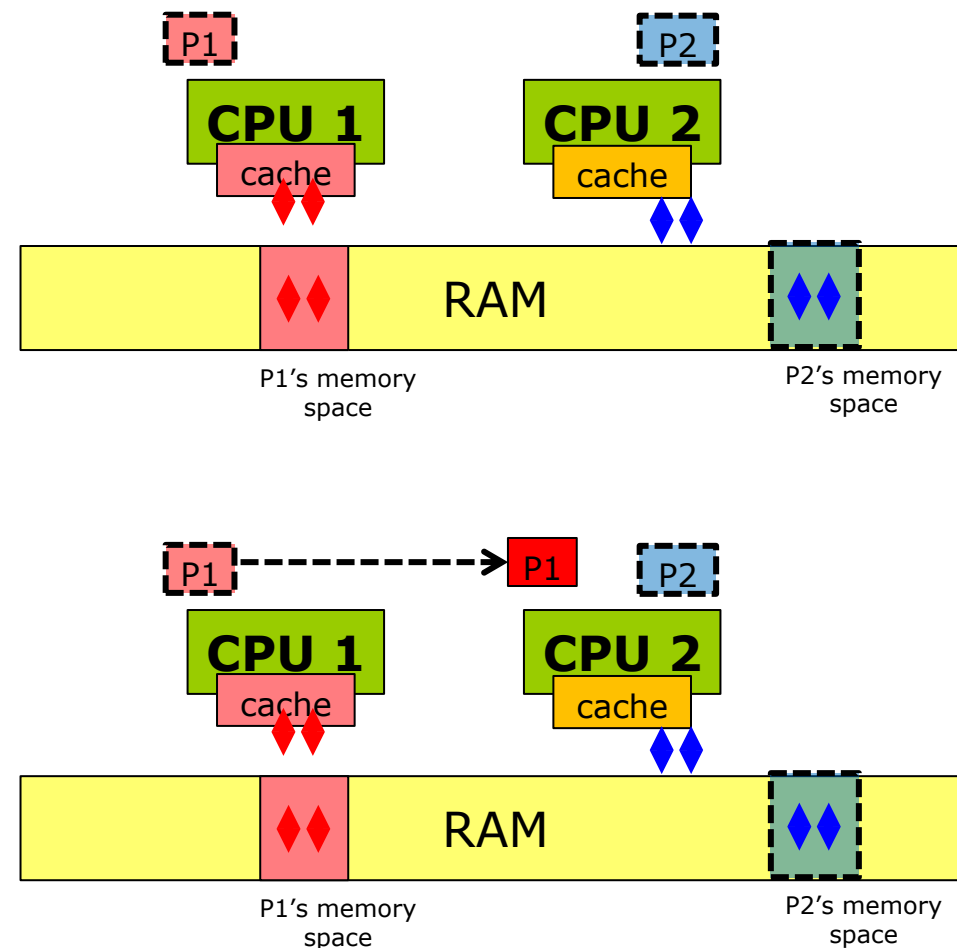
Load balancing

- **Load balancing** needed for per-core ready queues
 - Keep the workload between processors balanced
 - Migration
 - ▶ Move tasks from a busy processor to an idle processor



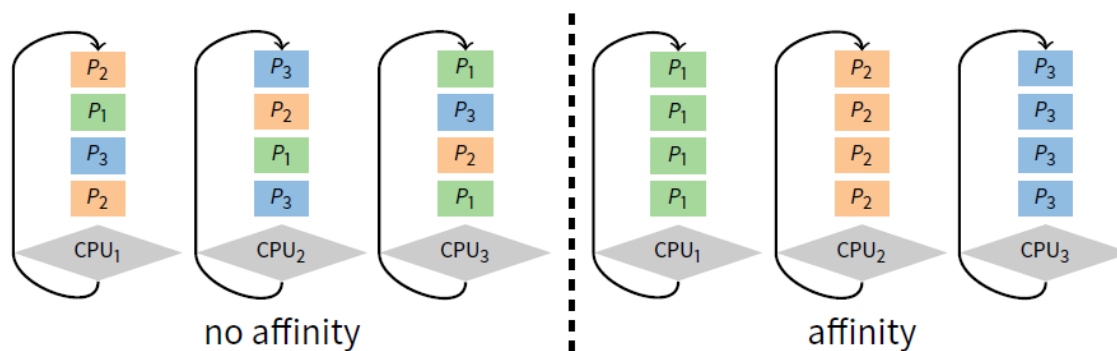
Processor Affinity

- When a thread has been running on one processor, the cache contents of that processor stores the memory accesses by that thread. (aka **warm cache**)
 - A thread has **processor affinity**
 - A lot of **cache hits**!
- Moving processes between processors has costs
 - The new processor cache must be repopulated - more **cache misses**!
 - e.g., common ready queue
 - e.g., load balancing



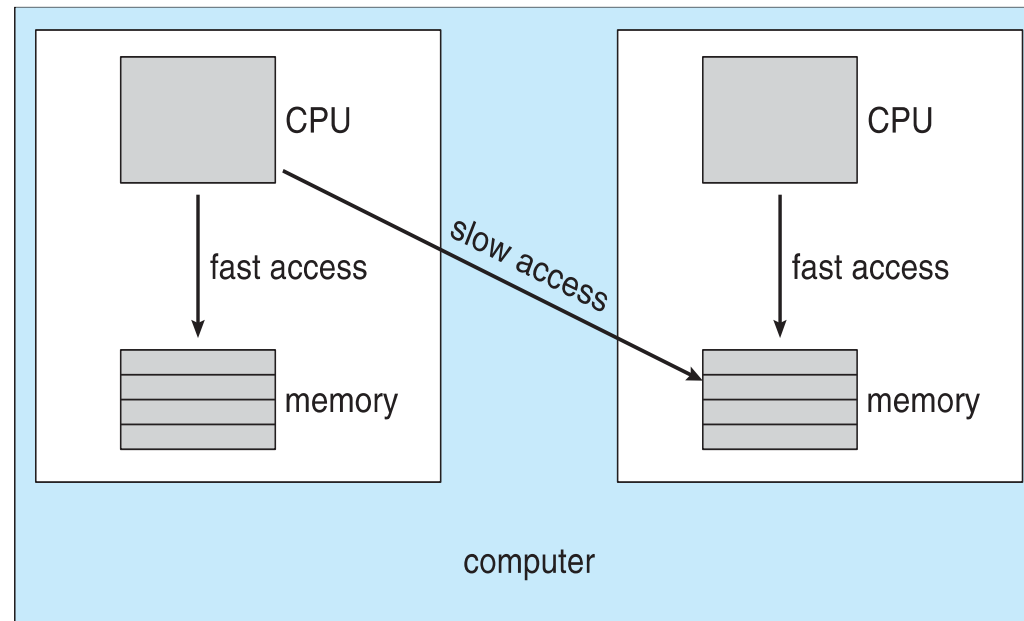
Processor Affinity

- To avoid **cache misses**, try to keep processes on same processor
= **Processor Affinity**
- Load balancing vs. Affinity Scheduling
 - Scheduler needs to balance between two



NUMA and CPU Scheduling

- If the operating system is Non-uniform memory access (**NUMA**)-**aware**, it will assign memory close to the CPU the thread is running on.
- NUMA systems also need **Processor affinity** for faster memory access!



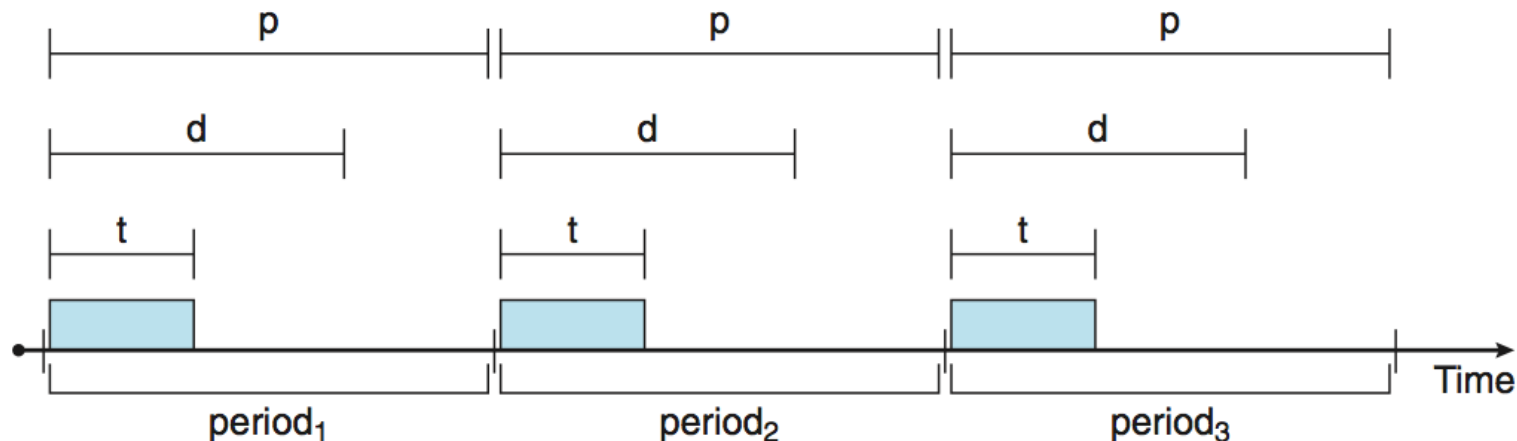
Chapter 5: CPU Scheduling

- Basic Concepts
- Scheduling Criteria
- Scheduling Algorithms
- **Advanced Scheduling**
 - Multiple-Processor Scheduling
 - **Real-Time CPU Scheduling**



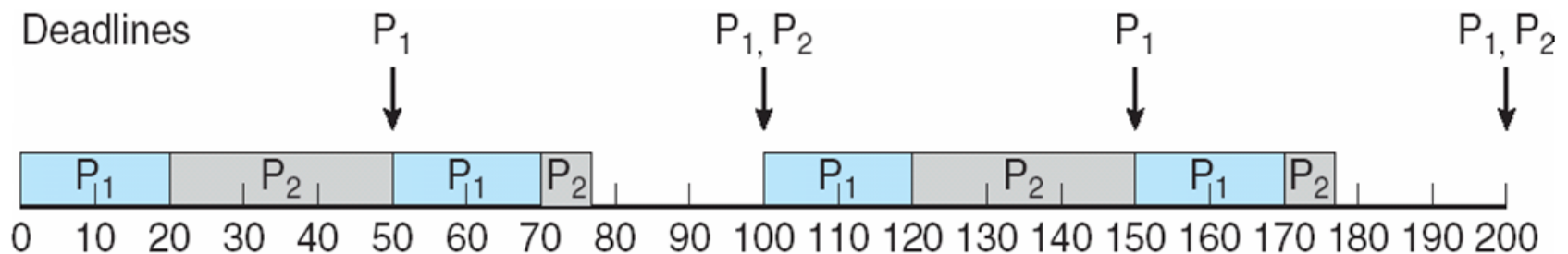
Real-time scheduling

- **Real-time scheduling** applications must provide ability to meet deadlines
- Processes have new characteristics: **periodic** ones require CPU at constant intervals
 - Has fixed **processing time t** , **deadline d** , **period p**
 - $0 \leq t \leq d \leq p$
 - **Rate** of periodic task is $1/p$



Rate Monotonic Scheduling

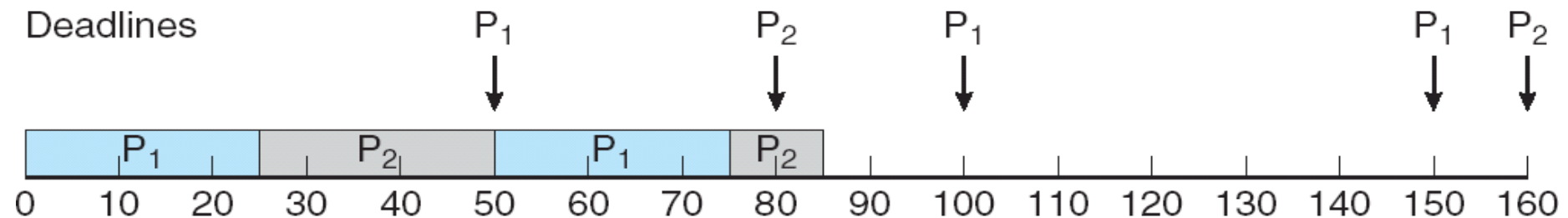
- A priority is assigned based on the inverse of its period (=rate)
 - Shorter periods = higher priority
 - Longer periods = lower priority
- Example
 - periods: $p_1 = 50$, $p_2 = 100$, CPU burst $t_1 = 20$, $t_2 = 35$
 - deadline = complete CPU burst before start of next period
 - P_1 is assigned a higher priority than P_2 .



Missed Deadlines with Rate Monotonic Scheduling

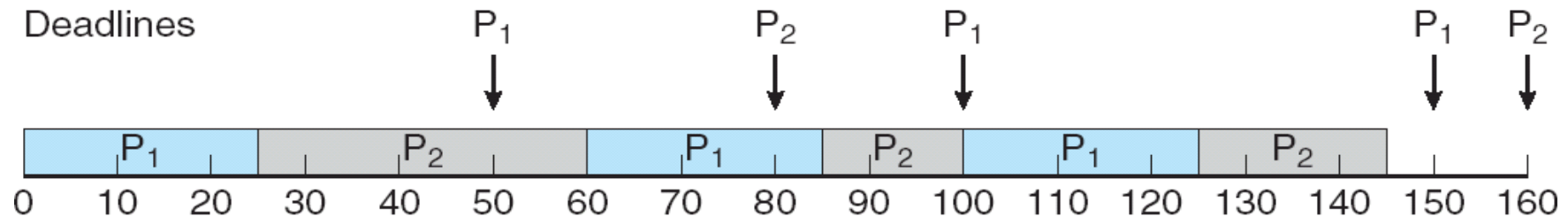
periods: $p_1 = 50$, $p_2 = 80$, CPU burst $t_1 = 25$,
 $t_2 = 35$

deadline = complete CPU burst before start
of next period



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



- periods: $p_1 = 50$, $p_2 = 80$, CPU burst $t_1 = 25$, $t_2 = 35$
- deadline = complete CPU burst before start of next period



<https://www.searchenginejournal.com/google-expands-rich-results-for-qa-pages-in-search/281355/>