# Chapter 5:  CPU Scheduling

# Outline

- Basic Concepts

- Scheduling Criteria

- Scheduling Algorithms

- Thread Scheduling

- Multi-Processor Scheduling

- Real-Time CPU Scheduling

- Operating Systems Examples
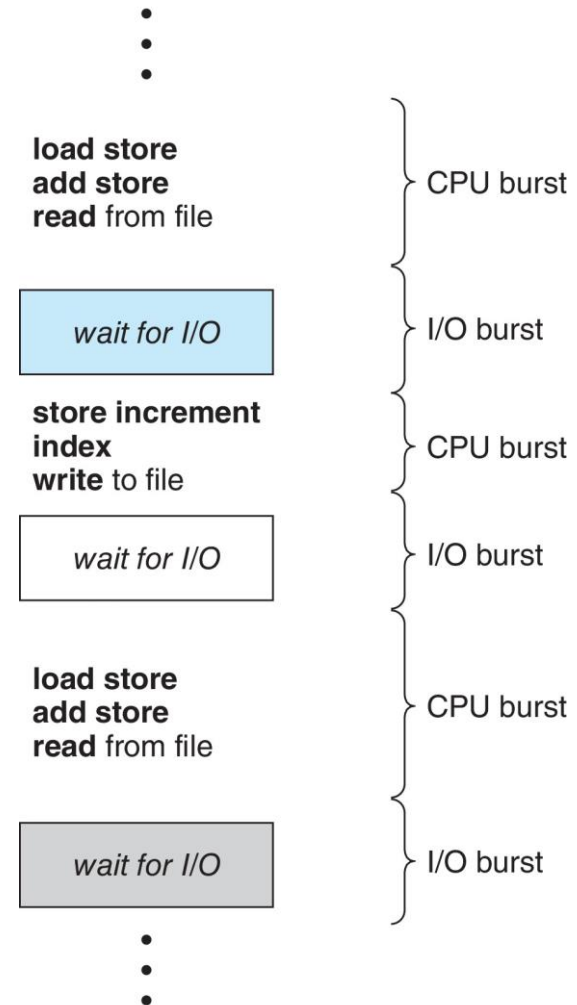
- Algorithm Evaluation

# Objectives

- Describe various CPU scheduling algorithms

- Assess CPU scheduling algorithms based on scheduling criteria

- Explain the issues related to multiprocessor and multicore scheduling

- Describe various real-time scheduling algorithms

- Describe the scheduling algorithms used in the Windows, Linux, and Solaris operating systems

- Apply modeling and simulations to evaluate CPU scheduling algorithms

# 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**
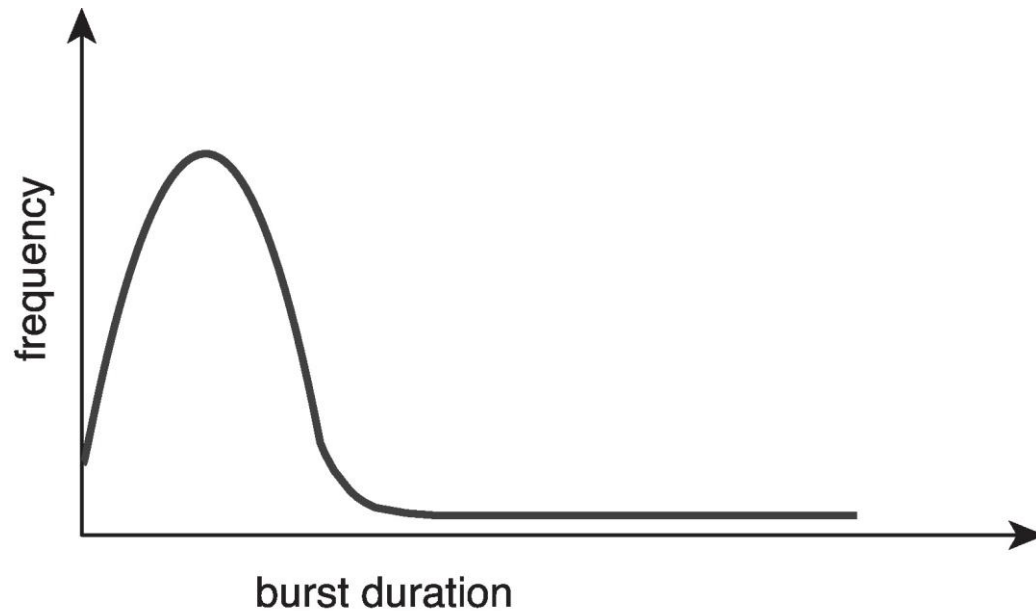
- CPU burst distribution is of main concern

```
   •
   •
   •
load store
add store      } CPU burst
read from file

wait for I/O   } I/O burst

store increment
index          } CPU burst
write to file

wait for I/O   } I/O burst

load store
add store      } CPU burst
read from file

wait for I/O   } I/O burst
   •
   •
   •
```

# Histogram of CPU-burst Times

Large number of short bursts

Small number of longer bursts

# CPU Scheduler

- The **CPU scheduler** selects from among the processes in ready queue, and allocates a CPU core to one of them
  - Queue may be ordered in various ways
- 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
- For situations 1 and 4, there is no choice in terms of scheduling. A new process (if one exists in the ready queue) must be selected for execution.
- For situations 2 and 3, however, there is  a choice.

# Preemptive and Nonpreemptive Scheduling

- When scheduling takes place only under circumstances 1 and 4, the scheduling scheme is **nonpreemptive**.

- Otherwise, it is **preemptive**.

- Under Nonpreemptive scheduling, once the CPU has been allocated to a process, the process keeps the CPU until it releases it either by terminating or by switching to the waiting state.

- Virtually all modern operating systems including Windows, MacOS, Linux, and UNIX use preemptive scheduling algorithms.

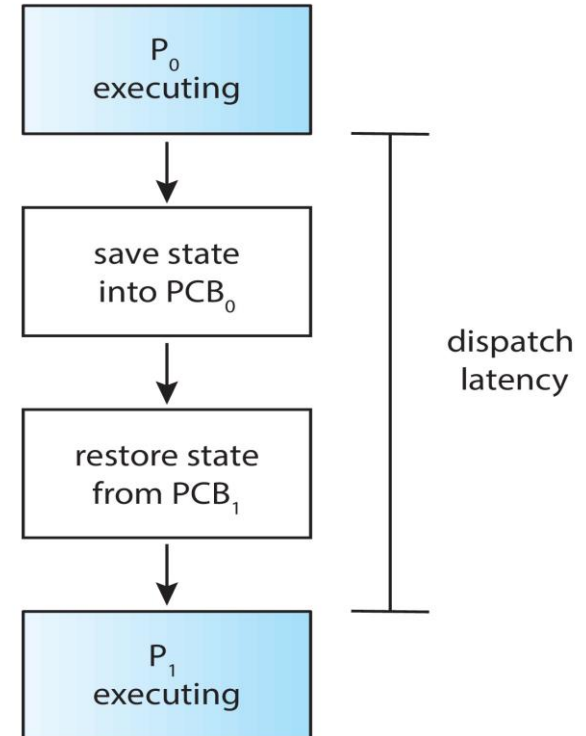# Preemptive Scheduling and Race Conditions

- Preemptive scheduling can result in race conditions when data are shared among several processes.

- Consider the case of two processes that share data. While one process is updating the data, it is preempted so that the second process can run. The second process then tries to read the data, which are in an inconsistent state.

- This issue will be explored in detail in Chapter 6.

# Dispatcher

- Dispatcher module gives control of the CPU to the process selected by the CPU 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.
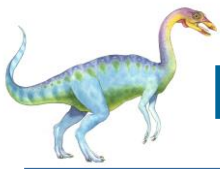
# Scheduling Algorithm Optimization Criteria

- Max CPU utilization

- Max throughput

- Min turnaround time

- Min waiting time

- Min response time

# 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:

| $P_1$ | $P_2$ | $P_3$ |
|-------|-------|-------|

0　　　　　　　　　　　　　　　　　24　　27　　30

- 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:

| P$_2$ | P$_3$ | P$_1$ |
|---|---|---|
| 0 | 3 | 6 |

(0 — 3 — 6 — 30)

- 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

- Preemptive version called **shortest-remaining-time-first**

- How do we determine the length of the next CPU burst?

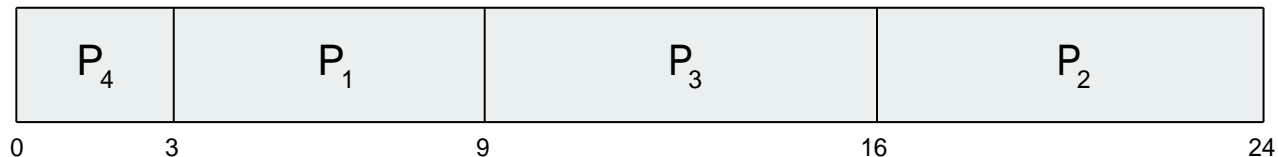  - Could ask the user

  - Estimate

# Example of SJF

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

- SJF scheduling chart

| $P_4$ | $P_1$ | $P_3$ | $P_2$ |
|:-----:|:-----:|:-----:|:-----:|

0       3               9                       16                      24

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

# Determining 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

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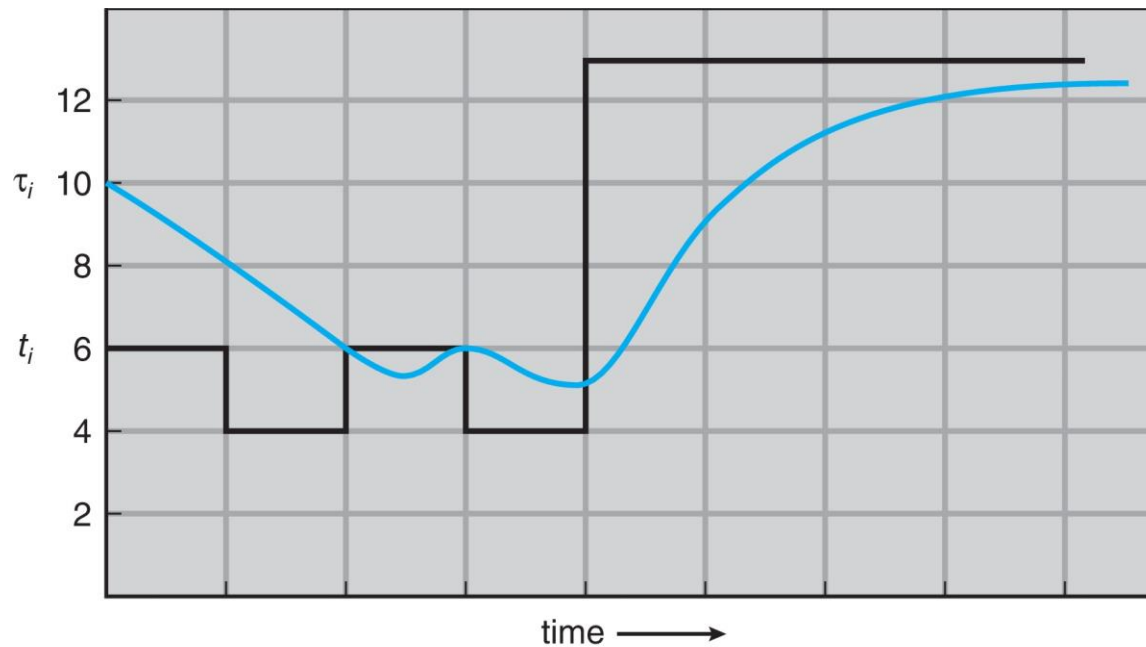  3. $\alpha, 0 \leq \alpha \leq 1$
  4. Define:
  $$\tau_{n+1} = \alpha\, t_n + (1 - \alpha)\tau_n.$$

- Commonly, α set to ½

| 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, we get:

$$\tau_{n+1} = \alpha\, t_n + (1 - \alpha)\alpha\, t_{n-1} + \ldots$$
$$+ (1 - \alpha)^j \alpha\, t_{n-j} + \ldots$$
$$+ (1 - \alpha)^{n+1} \tau_0$$

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

# Shortest Remaining Time First Scheduling

- Preemptive version of SJN

- Whenever a new process arrives in the ready queue, the decision on which process to schedule next is redone using the SJN algorithm.

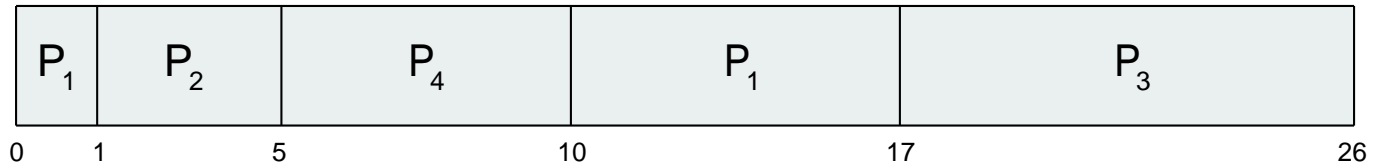- Is SRT more "optimal" than SJN in terms of the minimum average waiting time for a given set of processes?

# Example of 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

| P₁ | P₂ | P₄ | P₁ | P₃ |
|----|----|----|----|----|

0    1         5              10              17                    26

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

# 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

  - $q$ large $\Rightarrow$ FIFO (FCFS)

  - $q$ small $\Rightarrow$ RR

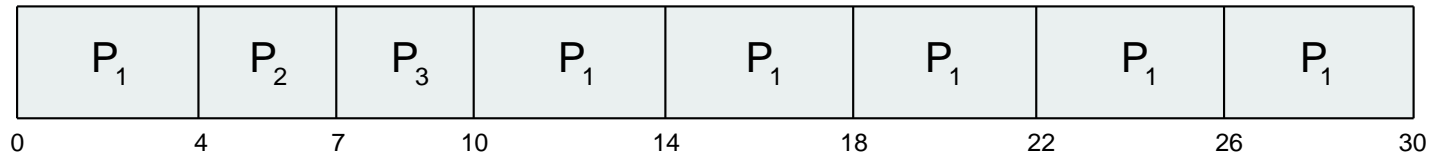- Note that q must be large with respect to context switch, otherwise overhead is too high

# Example of RR with Time Quantum = 4

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

- The Gantt chart is:

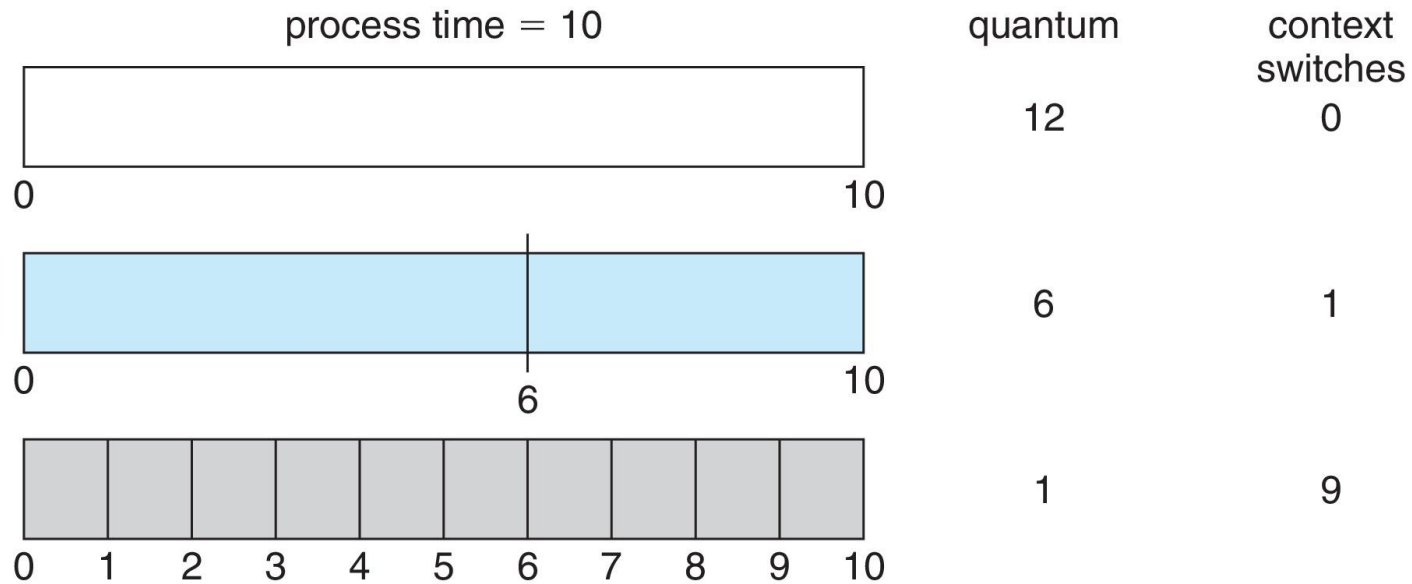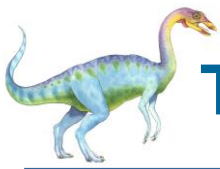| $P_1$ | $P_2$ | $P_3$ | $P_1$ | $P_1$ | $P_1$ | $P_1$ | $P_1$ |
|-------|-------|-------|-------|-------|-------|-------|-------|

0    4    7    10    14    18    22    26    30

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

- q should be large compared to context switch time

  - q usually 10 milliseconds to 100 milliseconds,
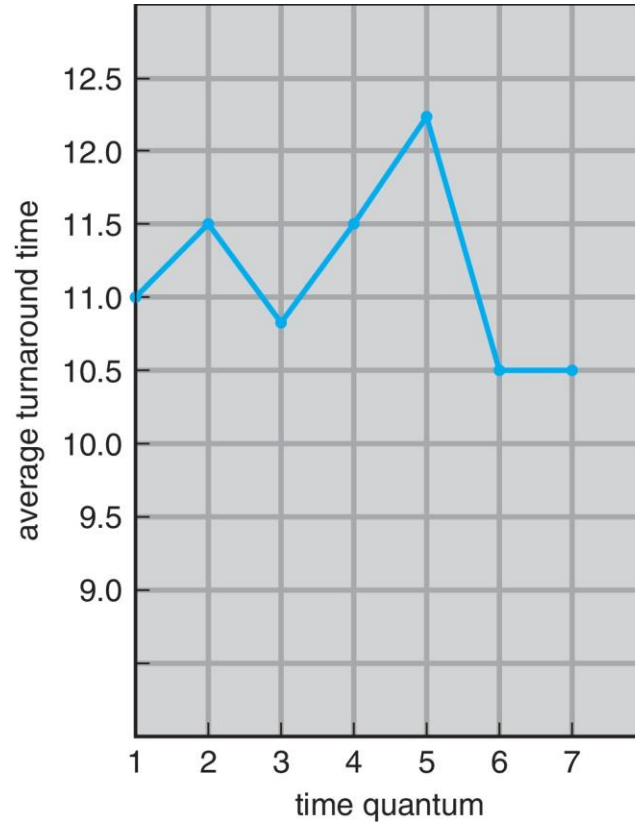  - Context switch < 10 microseconds

# Time Quantum and Context Switch Time

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

80% of CPU bursts
should be shorter than q

# 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)
  - Preemptive
  - Nonpreemptive

- SJF is priority scheduling where priority is the inverse of predicted next CPU burst time

- Problem ≡ **Starvation** – low priority processes may never execute

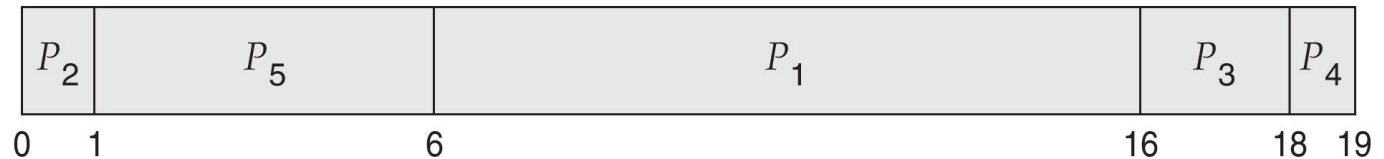- Solution ≡ **Aging** – as time progresses increase the priority of the process

# 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 |

- Priority scheduling Gantt Chart

| $P_2$ | $P_5$ | $P_1$ | $P_3$ | $P_4$ |
|:---:|:---:|:---:|:---:|:---:|

0   1           6                              16      18  19
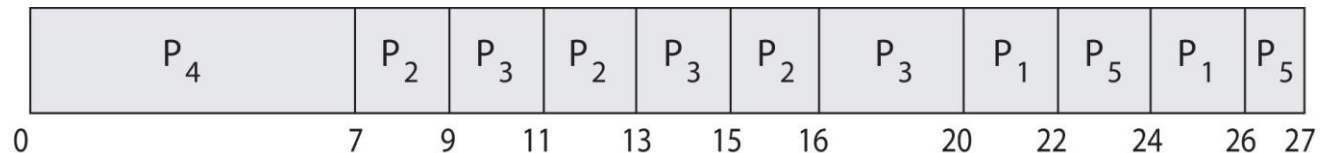
- Average waiting time = 8.2

# Priority Scheduling w/ Round-Robin

- Run the process with the highest priority. Processes with the same priority run round-robin

- Example:

| Process | Burst Time | Priority |
|---------|-----------|----------|
| $P_1$ | 4 | 3 |
| $P_2$ | 5 | 2 |
| $P_3$ | 8 | 2 |
| $P_4$ | 7 | 1 |
| $P_5$ | 3 | 3 |

- Gantt Chart with time quantum = 2

| $P_4$ | $P_2$ | $P_3$ | $P_2$ | $P_3$ | $P_2$ | $P_3$ | $P_1$ | $P_5$ | $P_1$ | $P_5$ |
|---|---|---|---|---|---|---|---|---|---|---|

0       7   9   11   13   15   16      20   22   24   26 27

# Multilevel Queue
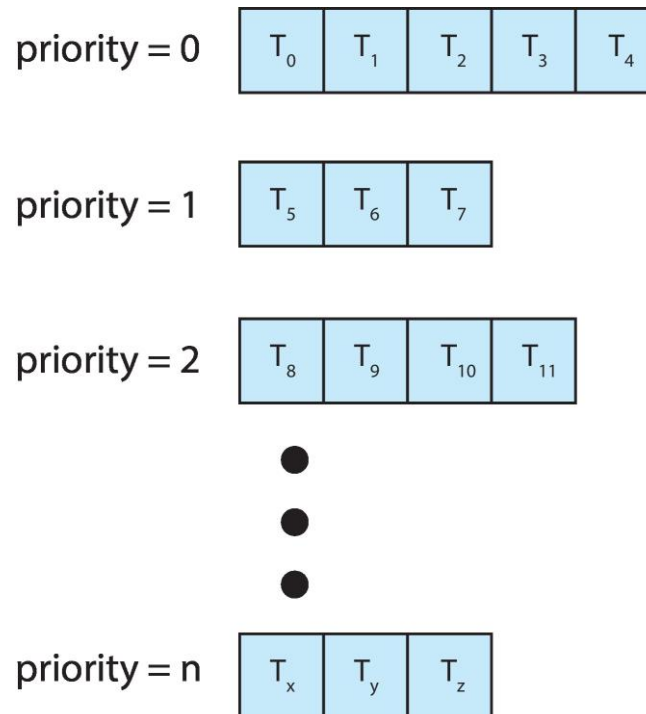
- The ready queue consists of multiple queues

- Multilevel queue scheduler defined by the following parameters:

  - Number of queues

  - Scheduling algorithms for each queue

  - Method used to determine which queue a process will enter when that process needs service

  - Scheduling among the queues

# Multilevel Queue
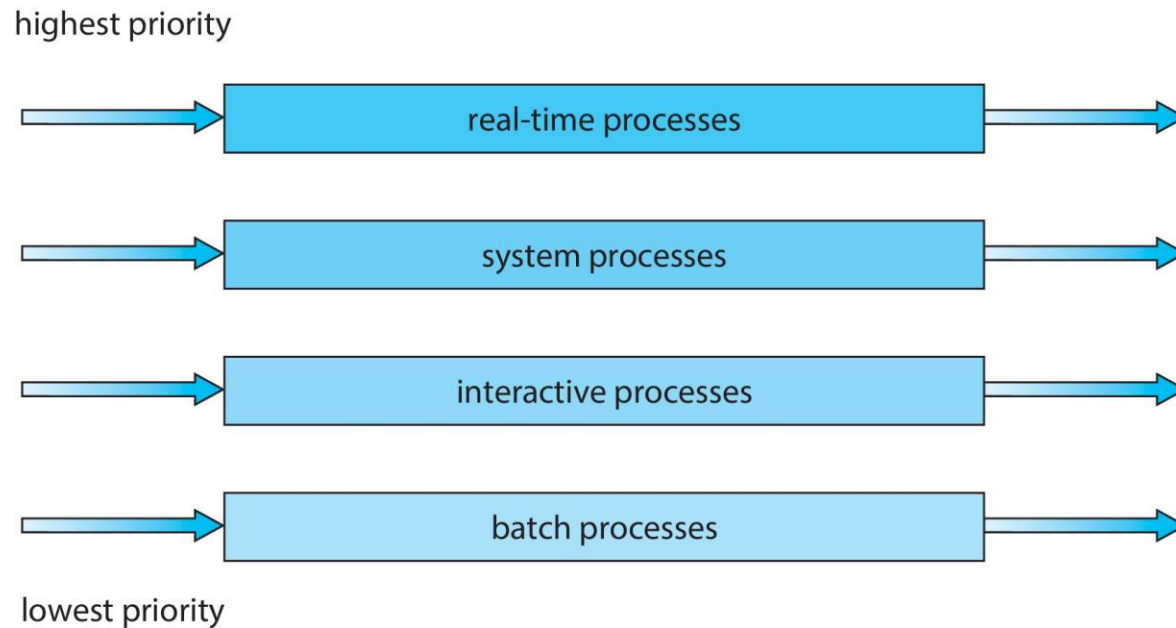
- With priority scheduling, have separate queues for each priority.
- Schedule the process in the highest-priority queue!

| | | | | | |
|---|---|---|---|---|---|
| priority = 0 | $T_0$ | $T_1$ | $T_2$ | $T_3$ | $T_4$ |

| | | | | |
|---|---|---|---|---|
| priority = 1 | $T_5$ | $T_6$ | $T_7$ | |

| | | | | |
|---|---|---|---|---|
| priority = 2 | $T_8$ | $T_9$ | $T_{10}$ | $T_{11}$ |

●
●
●

| | | | |
|---|---|---|---|
| priority = n | $T_x$ | $T_y$ | $T_z$ |

# Multilevel Queue

- Prioritization based upon process type

highest priority

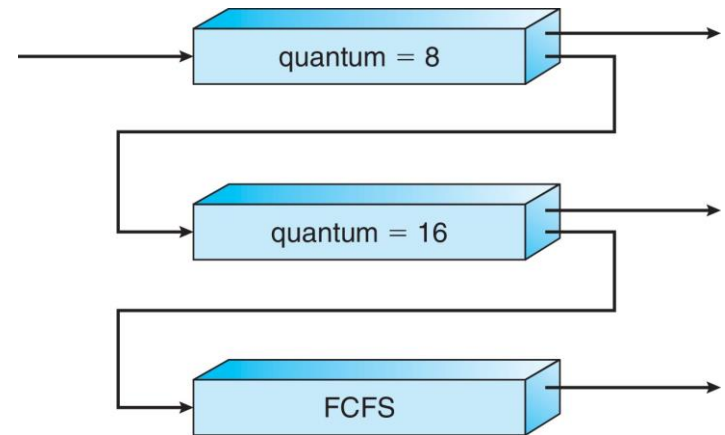| | |
|---|---|
| → | real-time processes → |
| → | system processes → |
| → | interactive processes → |
| → | batch processes → |

lowest priority

# Multilevel Feedback Queue

- Typically, when multilevel queue scheduling algorithm is used, processes are permanently assigned to a queue when they enter the system

  - low scheduling overhead, but it is inflexible

- A process can move between the various queues in Multilevel Feedback queue.

- 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

- Aging can be implemented using multilevel feedback queue

# 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 process enters queue $Q_0$ which is served in RR
    - ▸ When it gains CPU, the process receives 8 milliseconds
    - ▸ If it does not finish in 8 milliseconds, the process is moved to queue $Q_1$
  - At $Q_1$ job is again served in RR and receives 16 additional milliseconds
    - ▸ If it still does not complete, it is preempted and moved to queue $Q_2$

- This scheme leaves I/O-bound and interactive processes—which are typically characterized by short CPU bursts —in the higher-priority queues.

quantum = 8

quantum = 16

FCFS

# 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 (light weight process) **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 macOS (one-to-one model) only allow PTHREAD_SCOPE_SYSTEM

# Pthread Scheduling API

```c
#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

```c
/* 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);
}
```
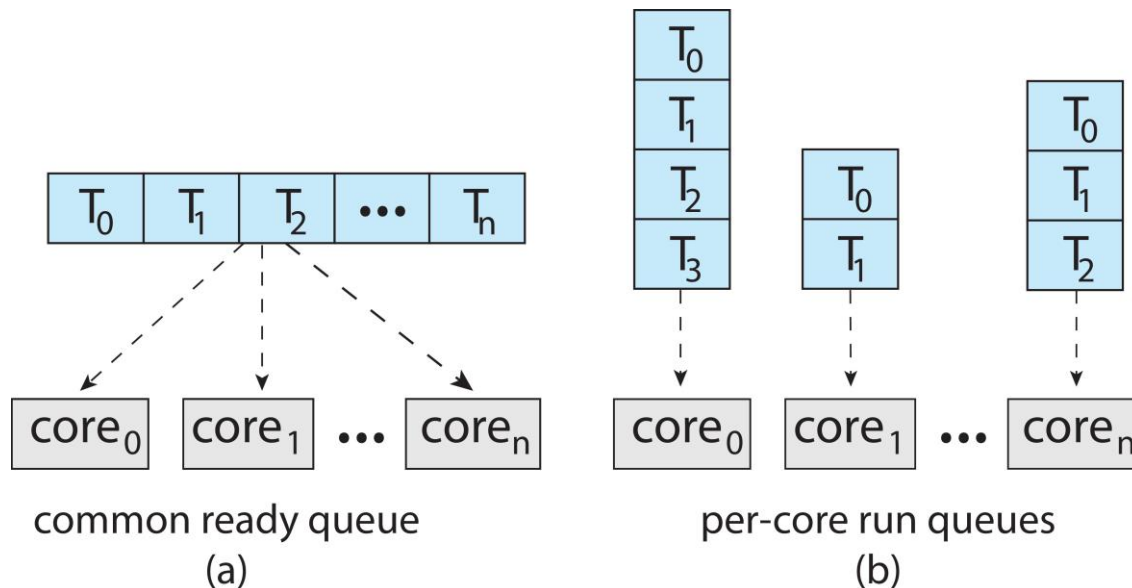
# Multiple-Processor Scheduling

- CPU scheduling more complex when multiple CPUs are available

- Multiprocess may be any one of the following architectures:

  - Multicore CPUs

  - Multithreaded cores

  - NUMA (Non-uniform Memory Access) systems

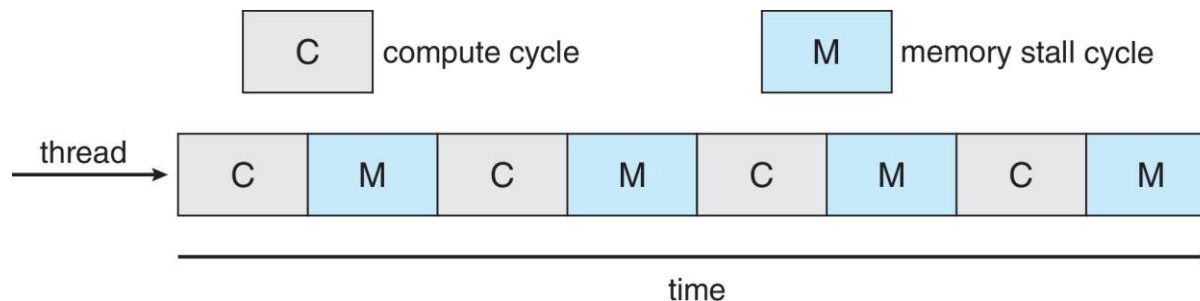  - Heterogeneous multiprocessing

# Multiple-Processor Scheduling

- **Symmetric multiprocessing** (SMP) is where each processor is self scheduling.

  - **asymmetric multiprocessing** = one master server + other processors execute only user code

- All threads may be in a common ready queue (a) [Race condition]

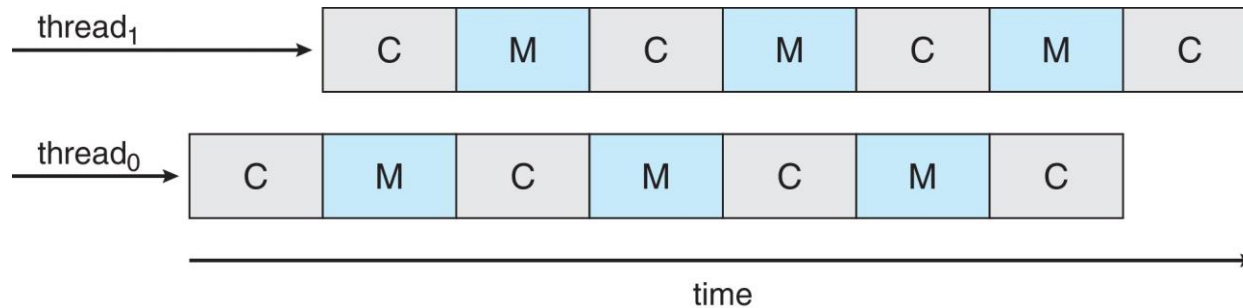- Each processor may have its own private queue of threads (b) [Mostly used]



common ready queue
(a)

per-core run queues
(b)

# Multicore Processors

- Recent trend to place multiple processor cores on same physical chip

  - Each core appears to the operating system to be a separate logical CPU

- 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
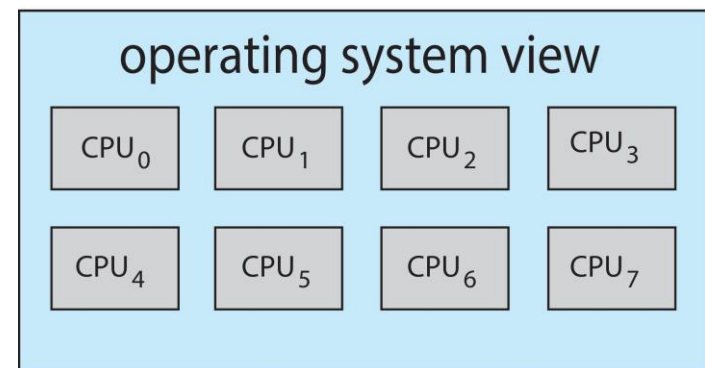
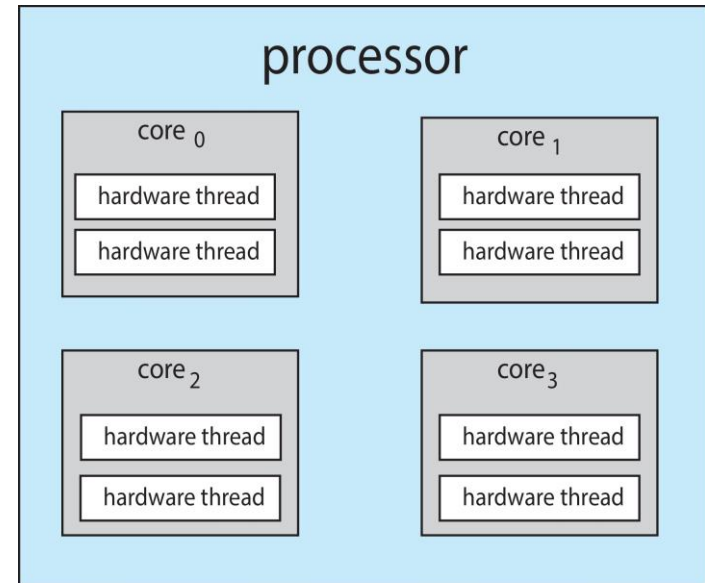- Figure

# Multithreaded Multicore System

- Each core has > 1 hardware threads.

- If one thread has a memory stall, switch to another thread!

- Figure

# Multithreaded Multicore System

- **Chip-multithreading** (CMT) assigns each core multiple hardware threads. (Intel refers to this as **hyperthreading**.)

  - i7—support two threads per core

- On a quad-core system with 2 hardware threads per core, the operating system sees 8 logical processors.
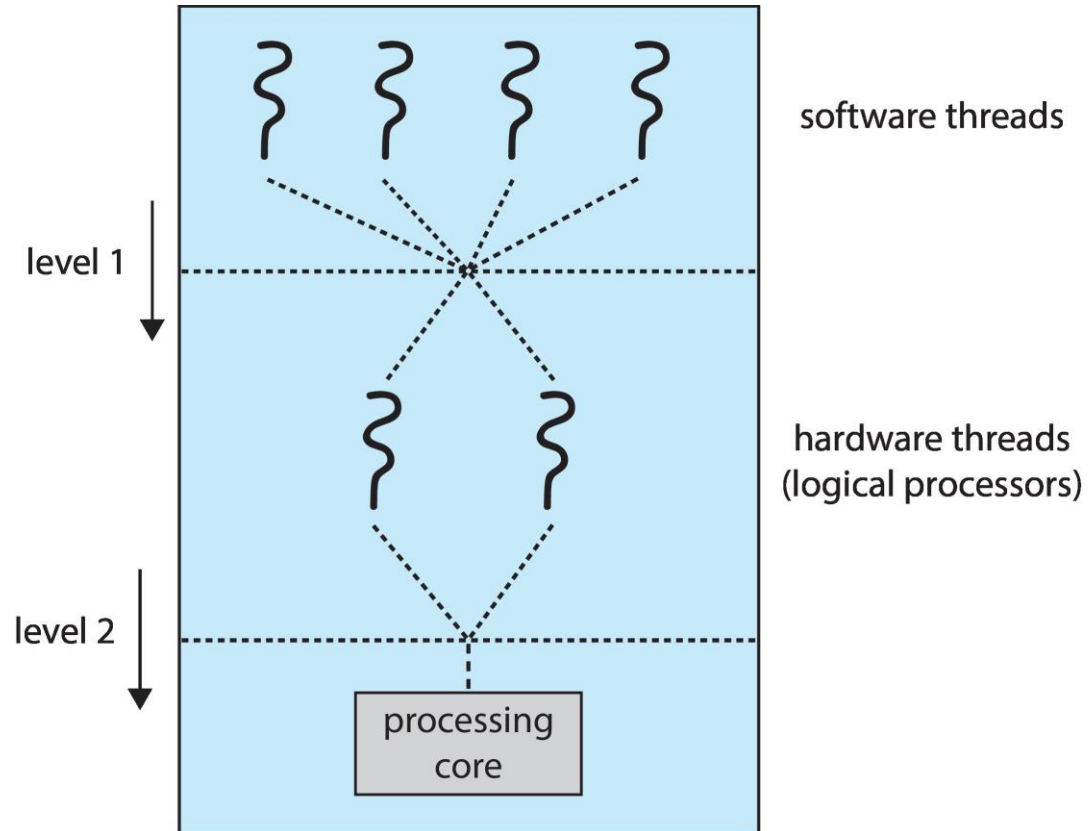
# Multithreaded Multicore System

- Two levels of scheduling:

  1. The operating system deciding which software thread to run on a logical CPU

  2. How each core decides which hardware thread to run on the physical core.

# 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

- Push and pull migration need not be mutually exclusive and are, in fact, often implemented in parallel on load-balancing systems

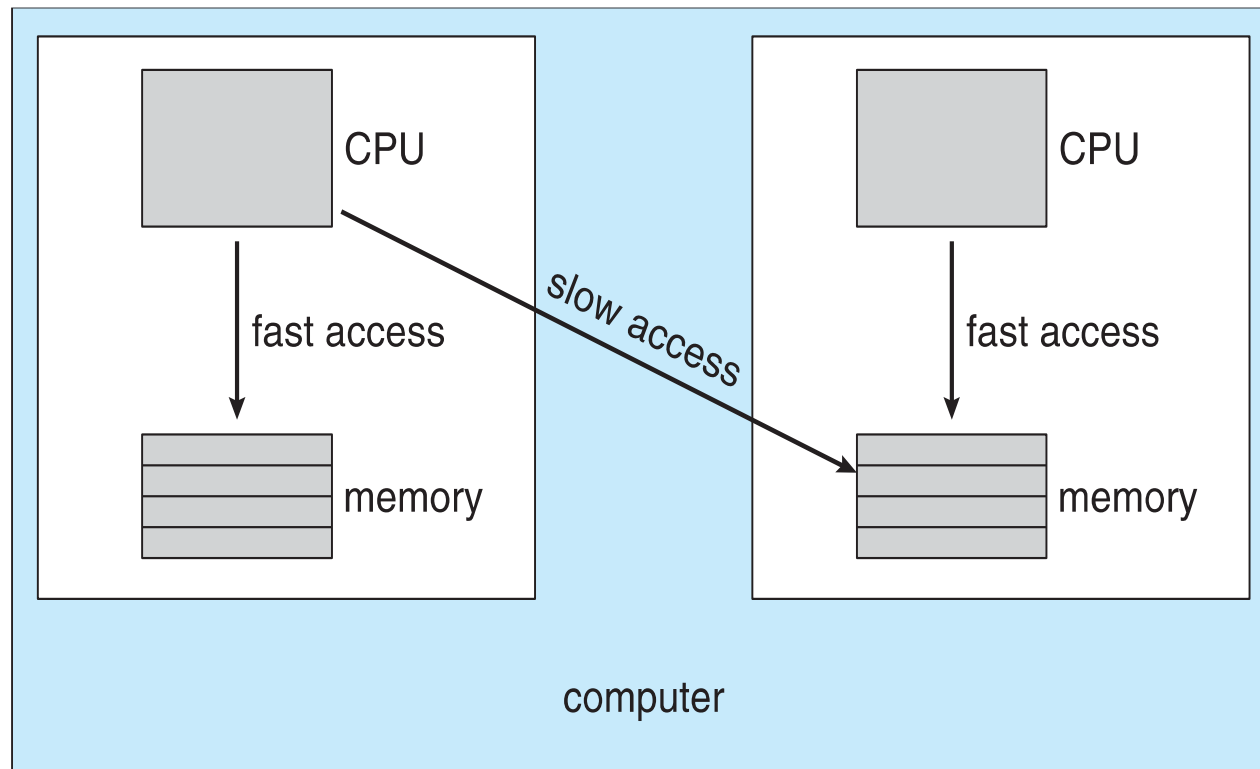# Multiple-Processor Scheduling – Processor Affinity

- When a thread has been running on one processor, the cache contents of that processor stores the memory accesses by that thread.

- most OS with SMP support try to avoid migrating a thread from one processor to another and instead attempt to keep a thread running on the same processor
  - We refer to this as a thread having affinity for a processor (i.e., "processor affinity")

- Load balancing may affect processor affinity as a thread may be moved from one processor to another to balance loads.
  - **Soft affinity** – the operating system attempts to keep a thread running on the same processor, but no guarantees.
  - **Hard affinity** – allows a process to specify a set of processors it may run on.

- Linux implements soft affinity; sched_setaffinity() system call can be used to modify it to hard affinity.

# NUMA and CPU Scheduling

If the operating system is **NUMA-aware**, it will assign memory closes to the CPU the thread is running on.
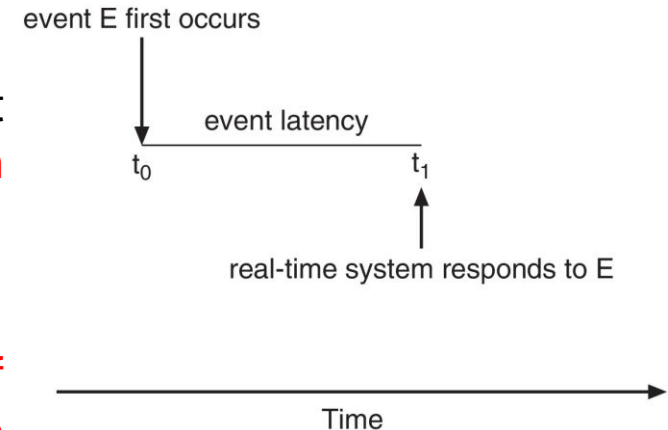
# Real-Time CPU Scheduling

- Can present obvious challenges

- **Soft real-time systems** – Critical real-time tasks have the highest priority, but no guarantee as to when tasks will be scheduled

- **Hard real-time systems –** task must be serviced by its deadline
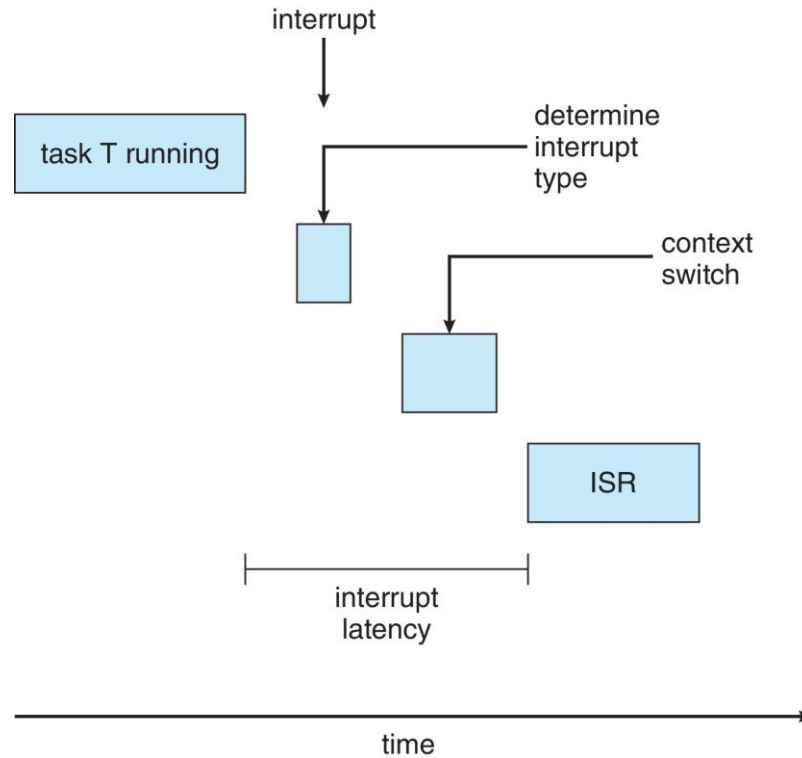  - service after the deadline has expired is the same as no service.

# Real-Time CPU Scheduling

- The system typically waits for an event in real time to occur.

  - Different events have different latency requirements

- Event latency – the amount of time that elapses from when an event occurs to when it is serviced.

- 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

event E first occurs

event latency

$t_0$            $t_1$
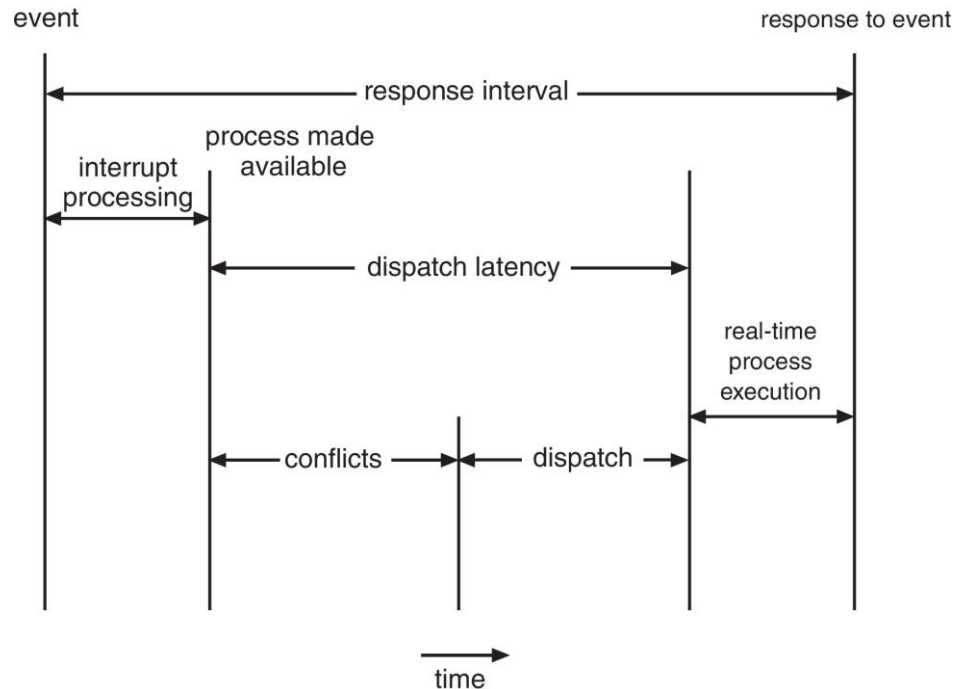
real-time system responds to E
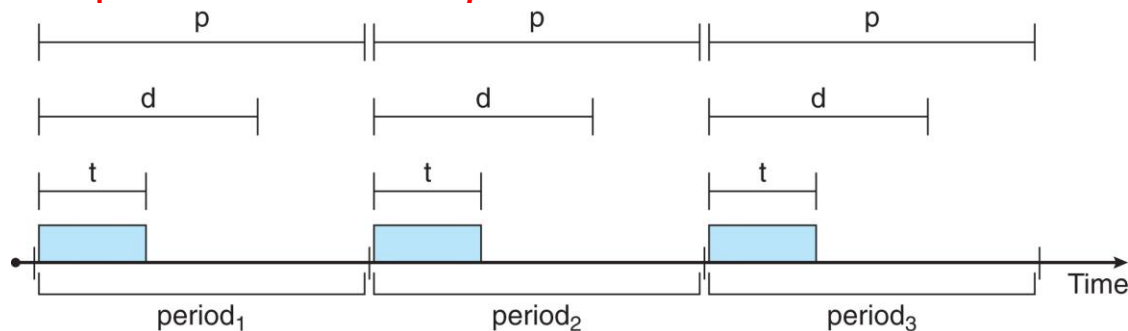
Time

# Interrupt Latency

# Dispatch Latency

- **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 scheduler 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 $1/p$

# Priority-based Scheduling

- Schedulers can take advantage of these characteristics and assign priorities according to a process's deadline or rate requirements.

- Process may have to announce its deadline requirements to the scheduler.

- Admission Control:

  - It either admits the process, guaranteeing that the process will complete on time.

  - or rejects the request.
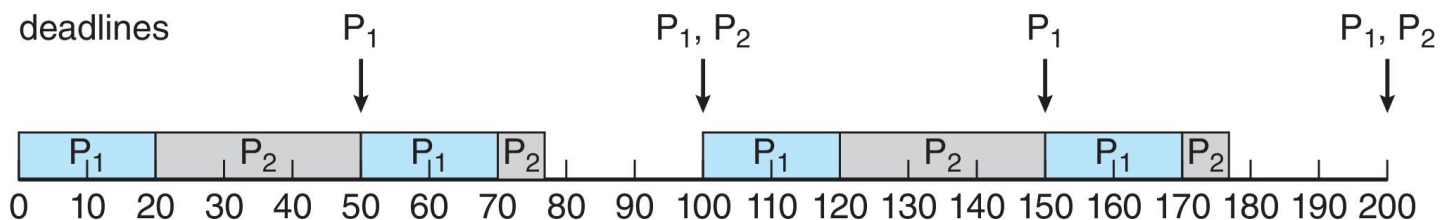
# Rate Monotonic Scheduling

- The rationale behind this policy is to assign a higher priority to tasks that require the CPU more often.

- **Assumption**: Processing time (t) of a periodic process is the same for each CPU burst.

- A priority is assigned based on the inverse of its period

- Shorter periods = higher priority;

- Longer periods = lower priority

# Rate Monotonic Scheduling: Example

- Two processes $P_1$ and $P_2$

- Periods: 50 and 100

- Processing Time: 20 and 35

- Deadline: Start of next period

- Can we schedule both the processes?

  - CPU utilization>1 (Not possible) <1 (??)

- Rate-monotonic scheduling is considered optimal

  - if a set of processes cannot be scheduled by this algorithm, it cannot be scheduled by any other algorithm that assigns static priorities.

# Rate Monotonic Scheduling: Example

- Two processes $P_1$ and $P_2$
- Periods: 50 and 80
- Processing Time: 25 and 35
- Deadline: Start of next period
- Observation??

# Rate Monotonic Scheduling: Example

- Two processes $P_1$ and $P_2$

- Periods: 50 and 80

- Processing Time: 25 and 35

- Deadline: Start of next period

- Observation??

  - Despite being optimal, then, rate-monotonic scheduling has a limitation: CPU utilization is bounded

  - The worst-case CPU utilization for scheduling $N$ processes

$$N(2^{1/N} - 1)$$

# Other Scheduling Algorithms

- Earliest Deadline First Scheduling (EDF)

- Proportional Share Scheduling

- Read by yourself !!

# 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:

  - Both functions return nonzero values if an error occurs.

  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

```c
#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");
    }
```

```
    /* 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);

}
```

# Operating System Examples

- Linux scheduling

- Windows scheduling

- Solaris scheduling

# Linux Scheduling Through Version 2.5

- Prior to kernel version 2.5, ran variation of standard UNIX scheduling algorithm
  - These algorithm was not designed with SMP systems in mind
- resulted in poor performance for systems with a large number of runnable processes.
- With Version 2.5 of the kernel included a scheduling algorithm— known as *O*(1)
- Also included support for
  - SMP systems,
    - ▸ Processor affinity
    - ▸ load balancing between processors
- However, poor response times for the interactive processes.
- In version 2.6 kernel, the scheduler was again revised
  - ***Completely Fair Scheduler*** (CFS)

# 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

  - Two scheduling classes included, others can be added
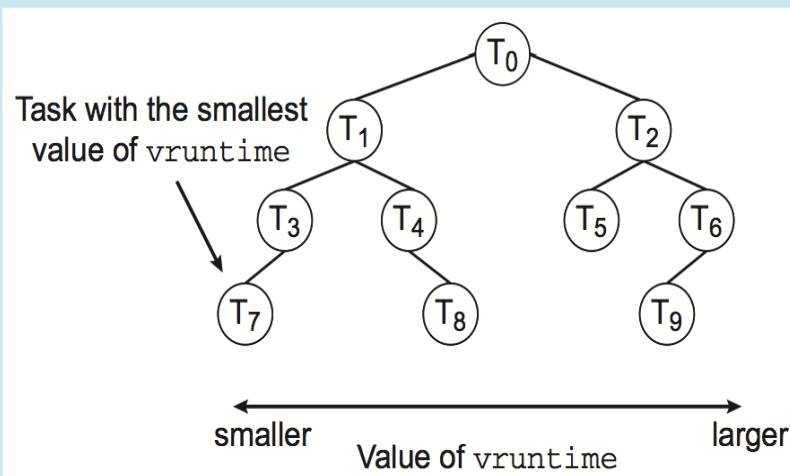
    1. default

    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 runtime** 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

# CFS Performance

The Linux CFS scheduler provides an efficient algorithm for selecting which task to run next. Each runnable task is placed in a red-black tree—a balanced binary search tree whose key is based on the value of `vruntime`. This tree is shown below:
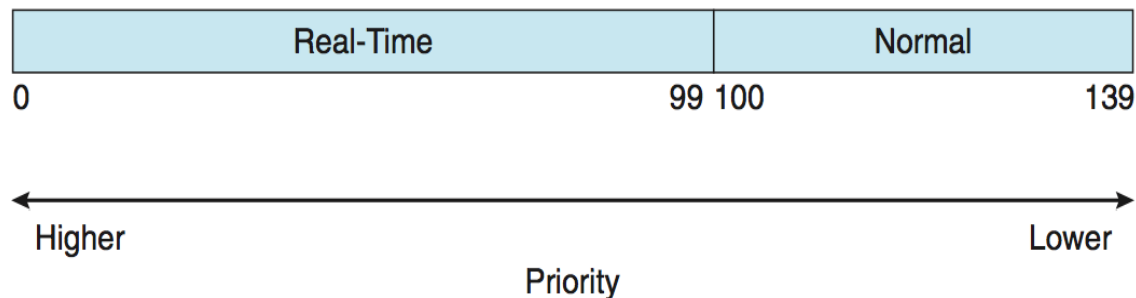


When a task becomes runnable, it is added to the tree. If a task on the tree is not runnable (for example, if it is blocked while waiting for I/O), it is removed. Generally speaking, tasks that have been given less processing time (smaller values of `vruntime`) are toward the left side of the tree, and tasks that have been given more processing time are on the right side. According to the properties of a binary search tree, the leftmost node has the smallest key value, which for the sake of the CFS scheduler means that it is the task with the highest priority. Because the red-black tree is balanced, navigating it to discover the leftmost node will require $O(lg N)$ operations (where $N$ is the number of nodes in the tree). However, for efficiency reasons, the Linux scheduler caches this value in the variable `rb_leftmost`, and thus determining which task to run next requires only retrieving the cached value.

# Linux Scheduling (Cont.)

- Any task scheduled using either the SCHED_FIFO or the SCHED_RR real-time policy runs at a higher priority than normal

- Linux uses two separate priority ranges, one for real-time tasks and a second for normal tasks.

- Real-time scheduling according to POSIX.1b
  - Real-time tasks have static priorities

- Real-time plus normal map into global priority scheme

- Nice value of -20 maps to global priority 100

- Nice value of +19 maps to priority 139

| Real-Time | Normal |
|-----------|--------|
| 0                                99 | 100                          139 |

Higher ← → Lower

Priority

# Windows Scheduling

- Read by yourself.

# End of Chapter 5