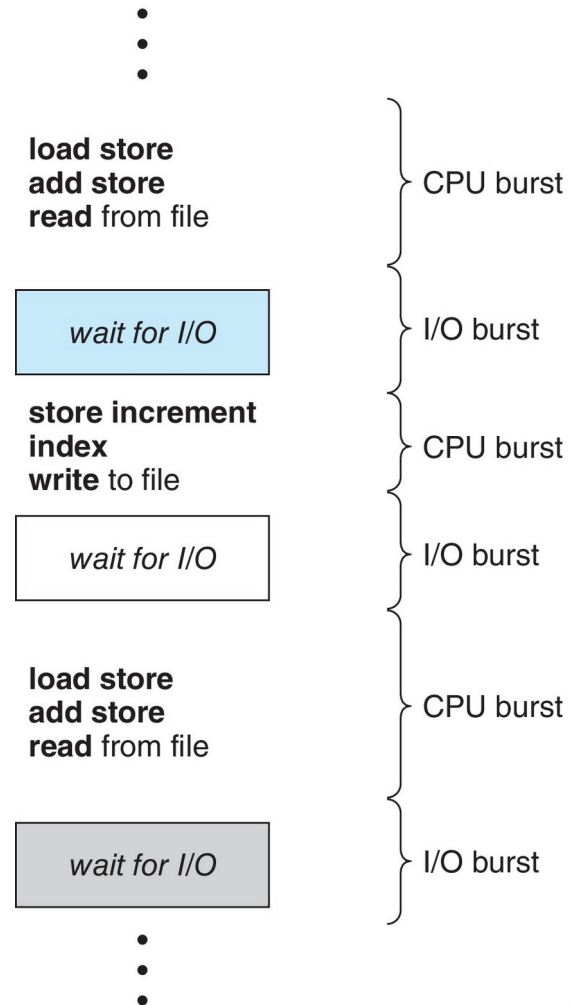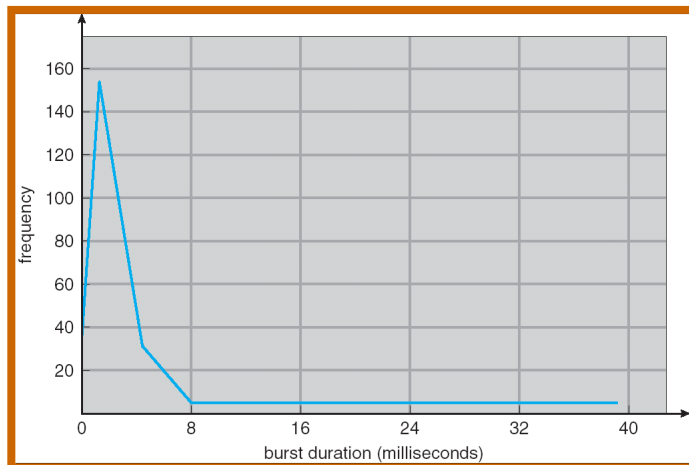# Chapter 5:  CPU Scheduling

# Outline

- Basic Concepts
- Scheduling Criteria
- Scheduling Algorithms
- Thread Scheduling
- Multi-Processor Scheduling
- Real-Time CPU Scheduling
- Operating Systems Examples
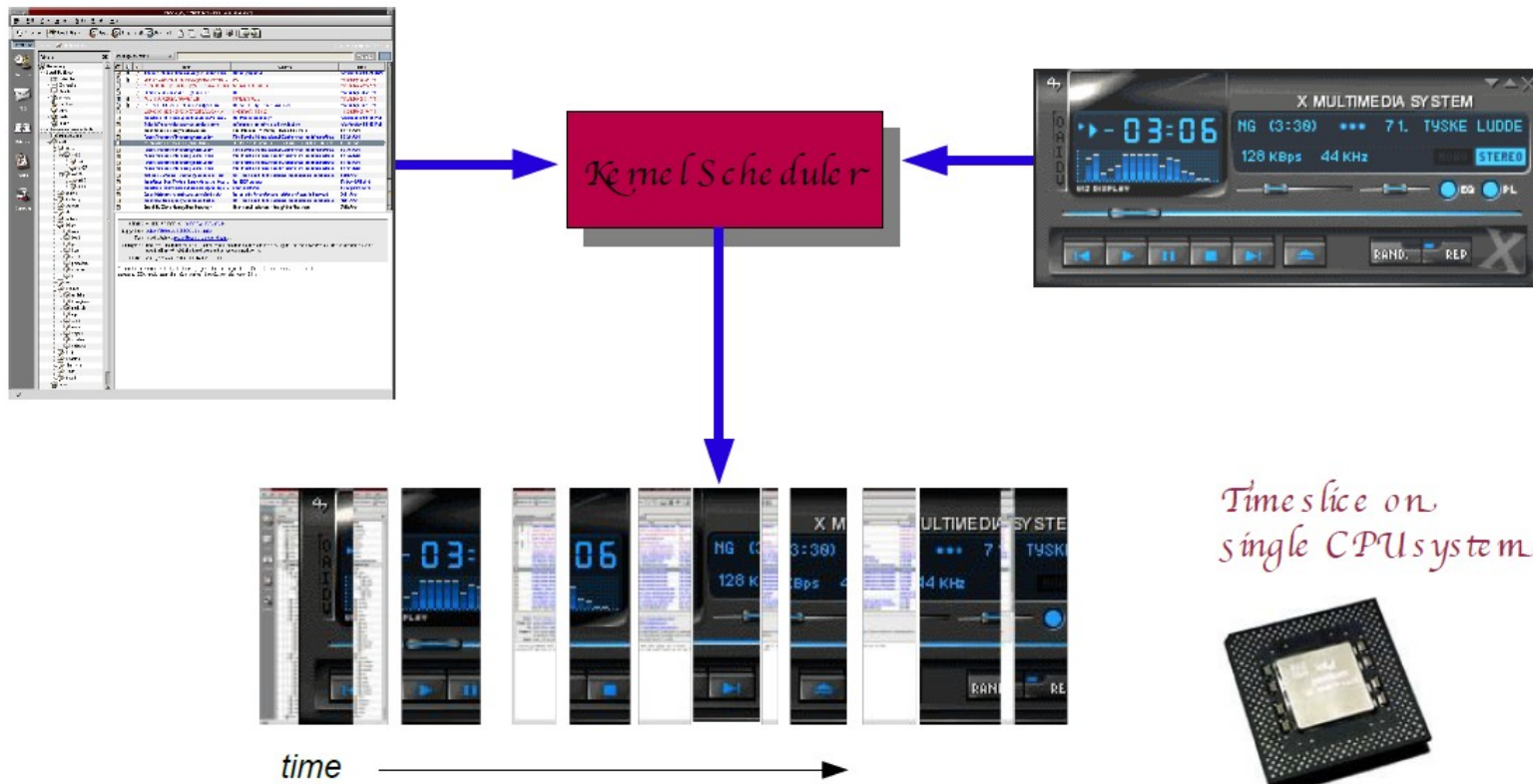- Algorithm Evaluation

# 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
    •
    •
    •
```

# Basic Concepts

- CPU Scheduler (Short term scheduler) switches among tasks frequently, i.e., 100 times/sec



Kernel Scheduler

time

Time slice on single CPU system

@ 2007 Matt Welsh

# 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**. (early version of WINDOWS)

- 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.

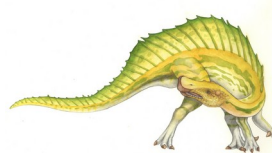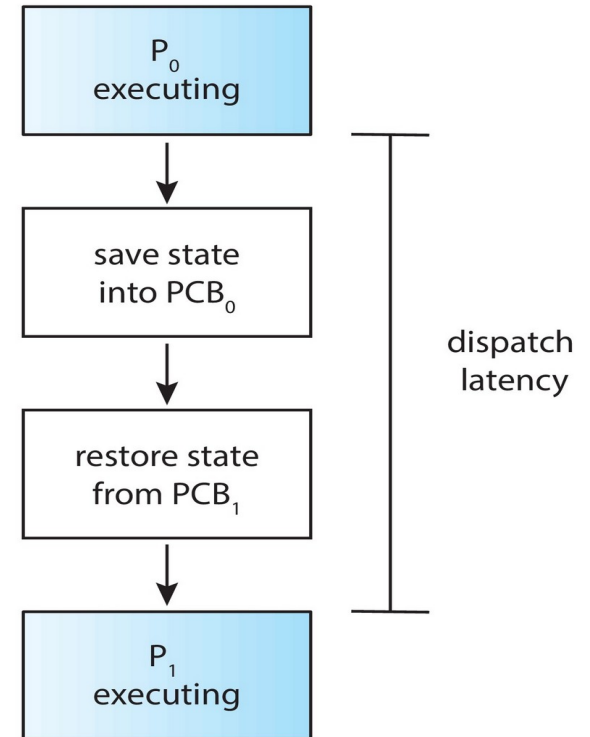# Preemptive Scheduling and Race Conditions

- 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 (is a process)

- 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. The actual time will depends on the hardware as well as the OS.

# Scheduling Criteria                    Optimization

- **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
  - Turnaround time = Finish time – Arrival time = Waiting time + CPU time

- **Waiting time** – amount of time a process has been waiting in the ready queue

- **Response time** – amount of time it takes from when a request was submitted until the first response is produced, **not** output  (for time-sharing environment) ( 1st segment of waiting time)

- 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

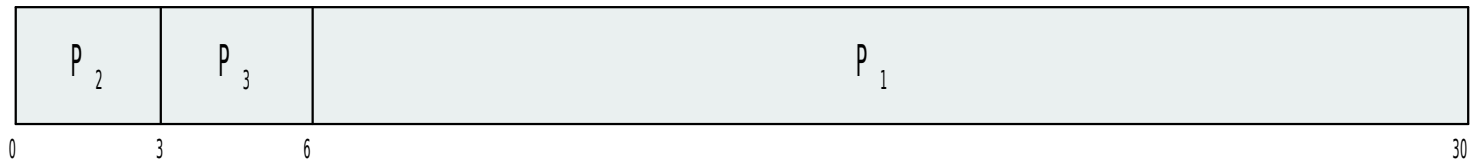# FCFS Scheduling (Cont.)

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                                                        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 Non-Preemptive SJF
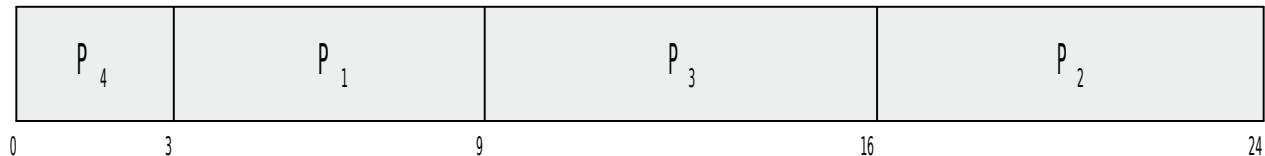
Process  Burst Time

$P_1$  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

# Example of Preemptive SJF

| Process | Arrival Time | Burst Time |
|---------|--------------|------------|
| $P_1$ | 0.0 | 7 |
| $P_2$ | 2.0 | 4 |
| $P_3$ | 4.0 | 1 |
| $P_4$ | 5.0 | 4 |

P1 = 7 − 2 = 5
P2 = 4 − 2 = 2, P1(5)
P3 = 5 − 4 = 1 (Finished), P2(2), P1(5)

Remaining time first:
P2(2): 7 − 5 = 2 (Finished), P1(5)
P4 = 11 − 7 = 4 (Finished), P1(5)
P1 = 16 − 11 = 5 (Finished)

- SJF (preemptive)

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

0   2     4   5     7         11             16

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

Turnaround time:
P1 = 16 − 0 = 16
P2 = 7 − 2 = 5
P3 = 5 − 4 = 1
P4 = 11 − 5 = 6

Waiting time:
P1 = 16 − 7 = 9
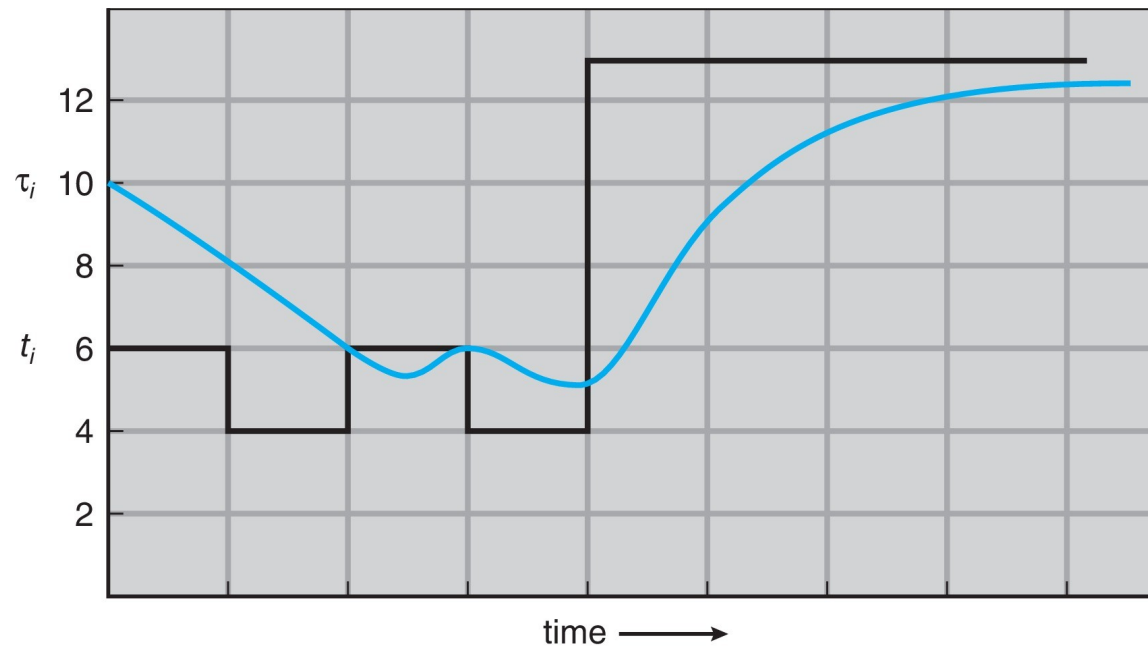P2 = 5 − 4 = 1
P3 = 1 − 1 = 0
P4 = 6 − 4 = 2

# 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:

- 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

# 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

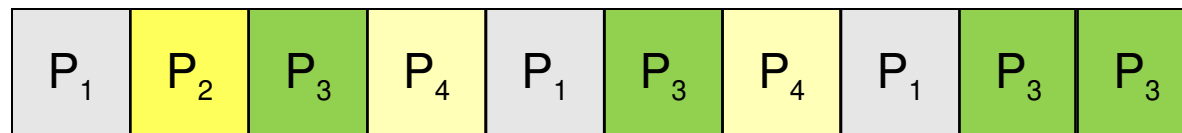# Example of RR with Time Quantum = 20

Process  Burst Time

$P_1$    53

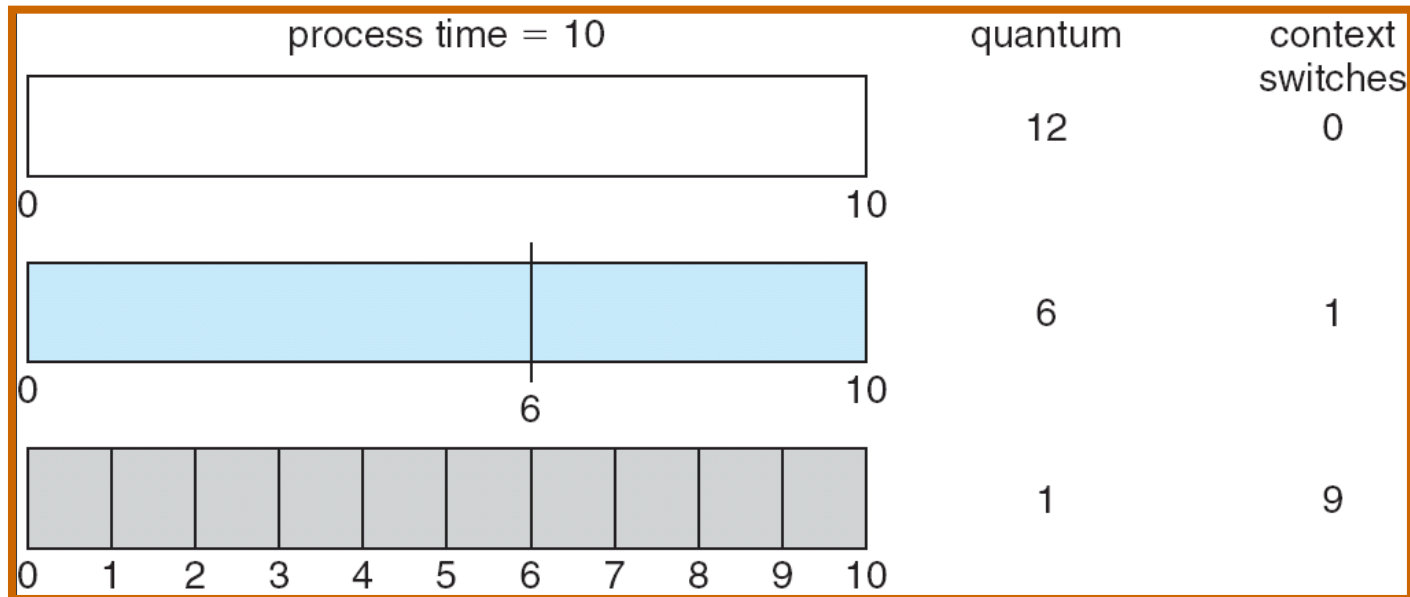$P_2$    17

$P_3$    68

$P_4$    24

- The Gantt chart is:

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

162

- Typically, higher average turnaround time than SJF, but better *response time.*
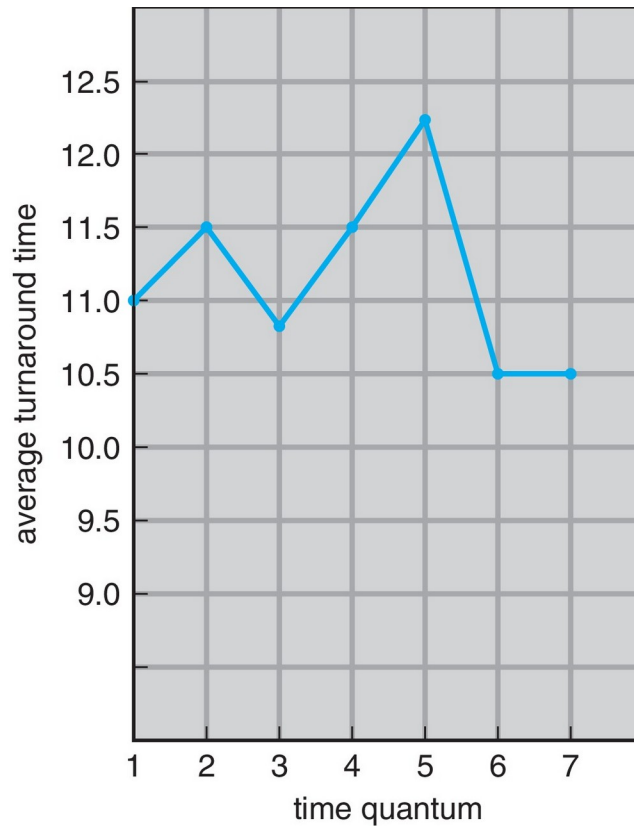
1. If the time quantum is extremely large, the RR policy is the same as the FCFS policy.
2. If the time quantum is small, it must be large with respect to context switch, otherwise overhead is too high.
3. q should be large compared to context switch time
   q usually 10 milliseconds to 100 milliseconds,
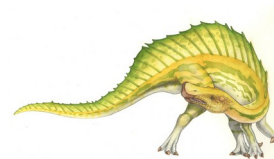   Context switch < 10 microseconds

# Turnaround Time Varies With The Time Quantum

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

# 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

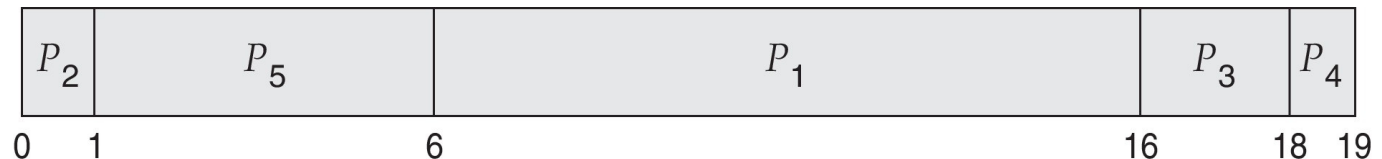- Solution ≡ **Aging** – as time progresses increase the priority of the process

# 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

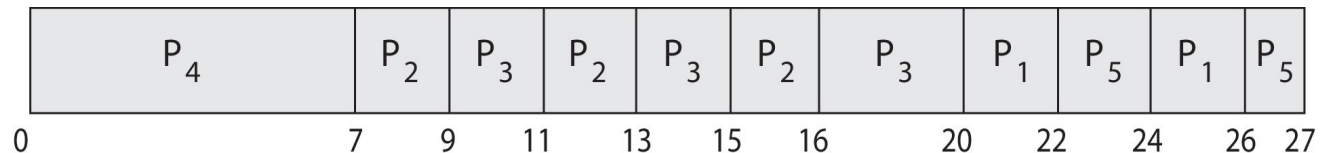- Average waiting time = 8.2

# 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
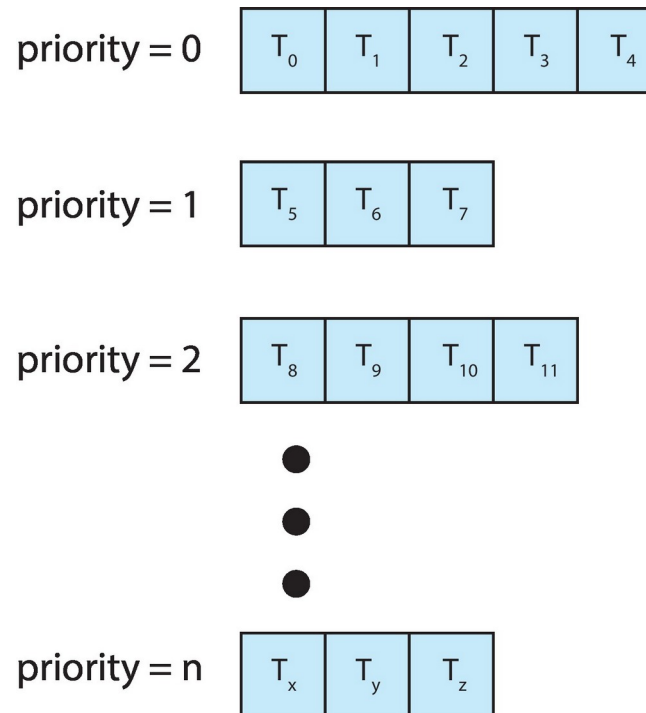
# Multilevel Queue

- 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
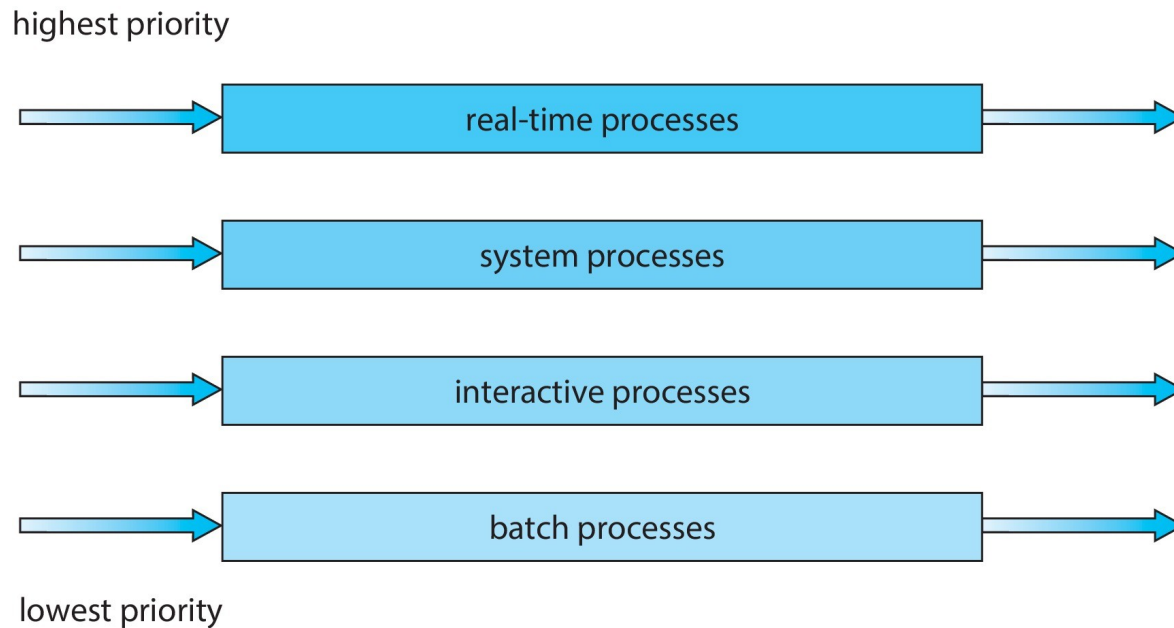
  - 20% to background in FCFS

# Multilevel Queue

- 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

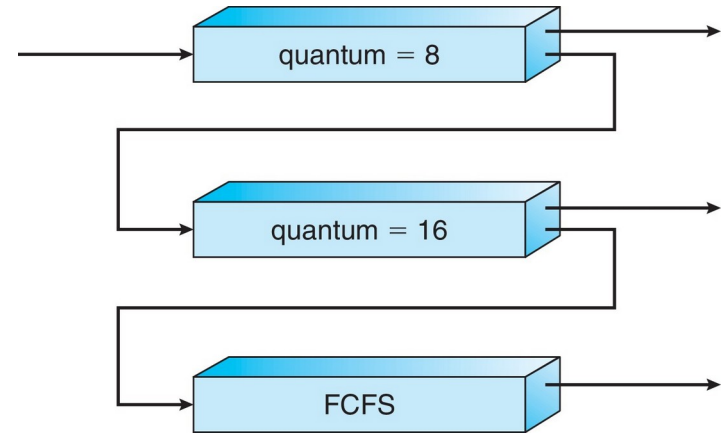# Multilevel Feedback Queue

- A process can move between the various queues.

- 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 $Q_2$

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 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

- 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 only allow PTHREAD_SCOPE_SYSTEM

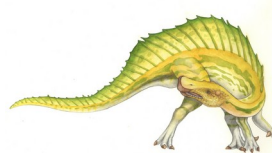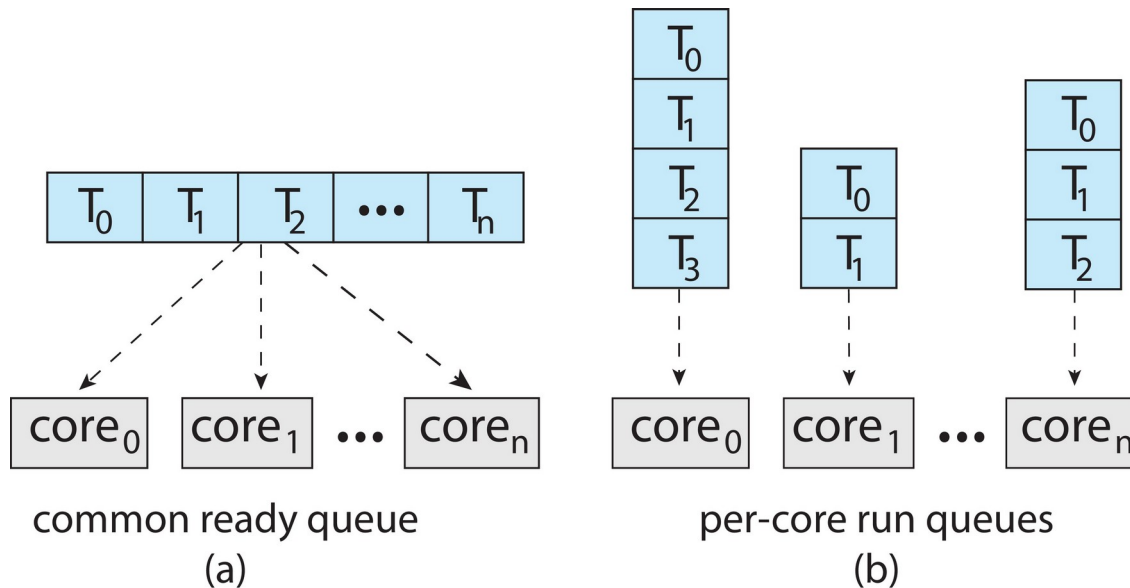# Multiple-Processor Scheduling

- CPU scheduling more complex when multiple CPUs are available

- Multiprocess may be any one of the following architectures:

  - Multicore CPUs

  - Multithreaded cores

  - NUMA systems

  - Heterogeneous multiprocessing

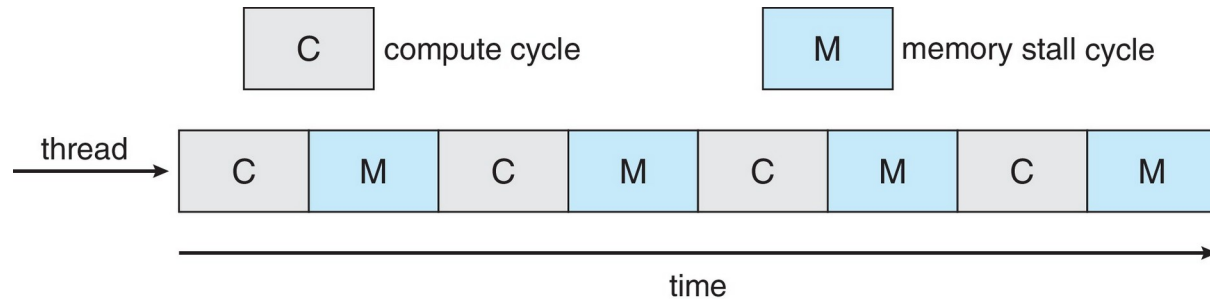# Multiple-Processor Scheduling

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

- All threads may be in a common ready queue (a)

- Each processor may have its own private queue of threads (b)

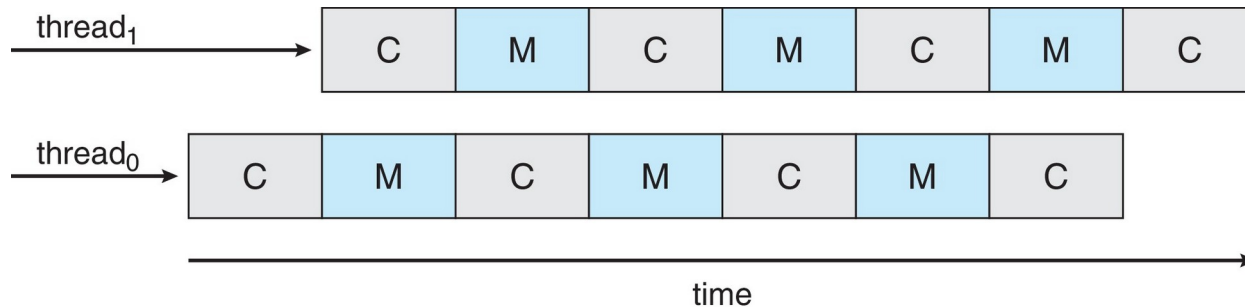common ready queue
(a)

per-core run queues
(b)

# 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
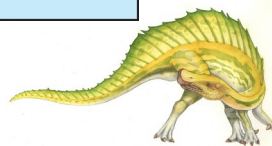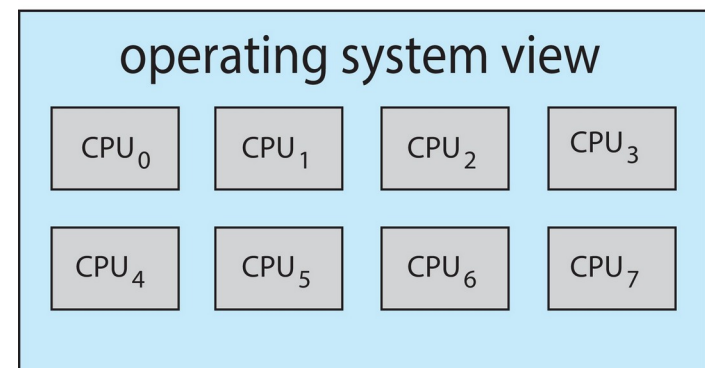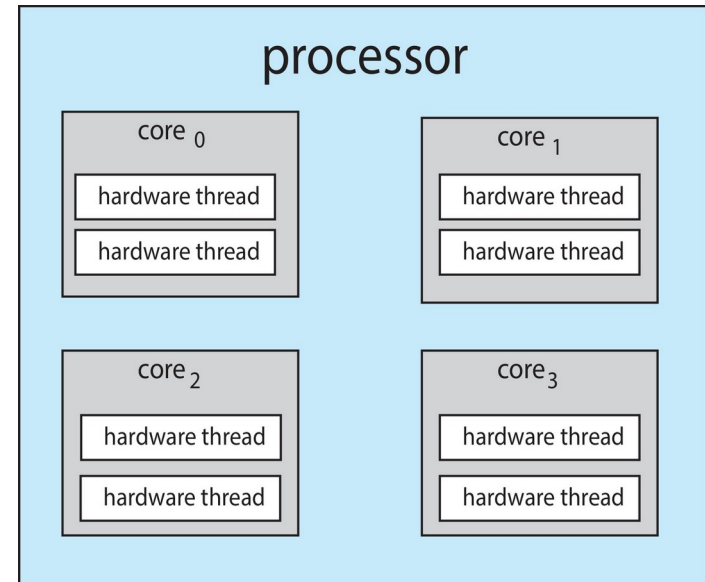
| C | compute cycle |   | M | memory stall cycle |

| thread → | C | M | C | M | C | M | C | M |

time

- Each core has > 1 hardware threads.

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

| thread$_1$ → | C | M | C | M | C | M | C |

| thread$_0$ → | C | M | C | M | C | M | C |

time

# Multithreaded Multicore System

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

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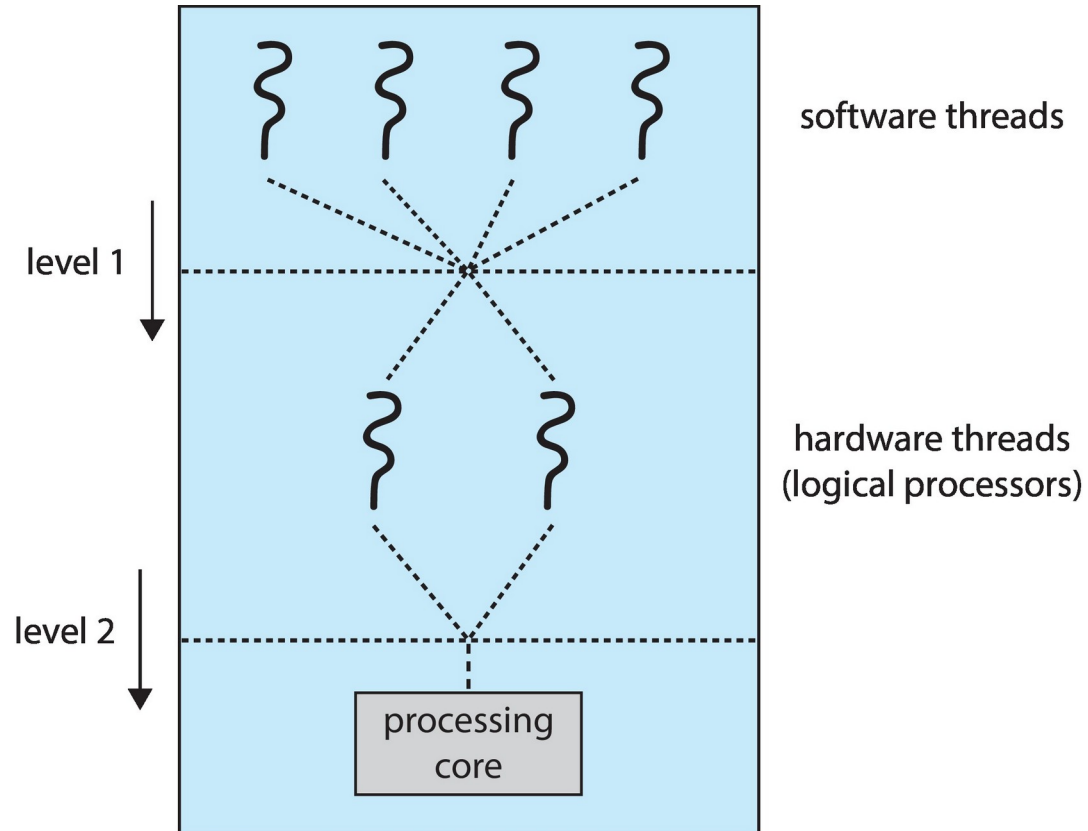
# Multithreaded Multicore System

- 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.

level 1

level 2

software threads

hardware threads
(logical processors)

processing
core

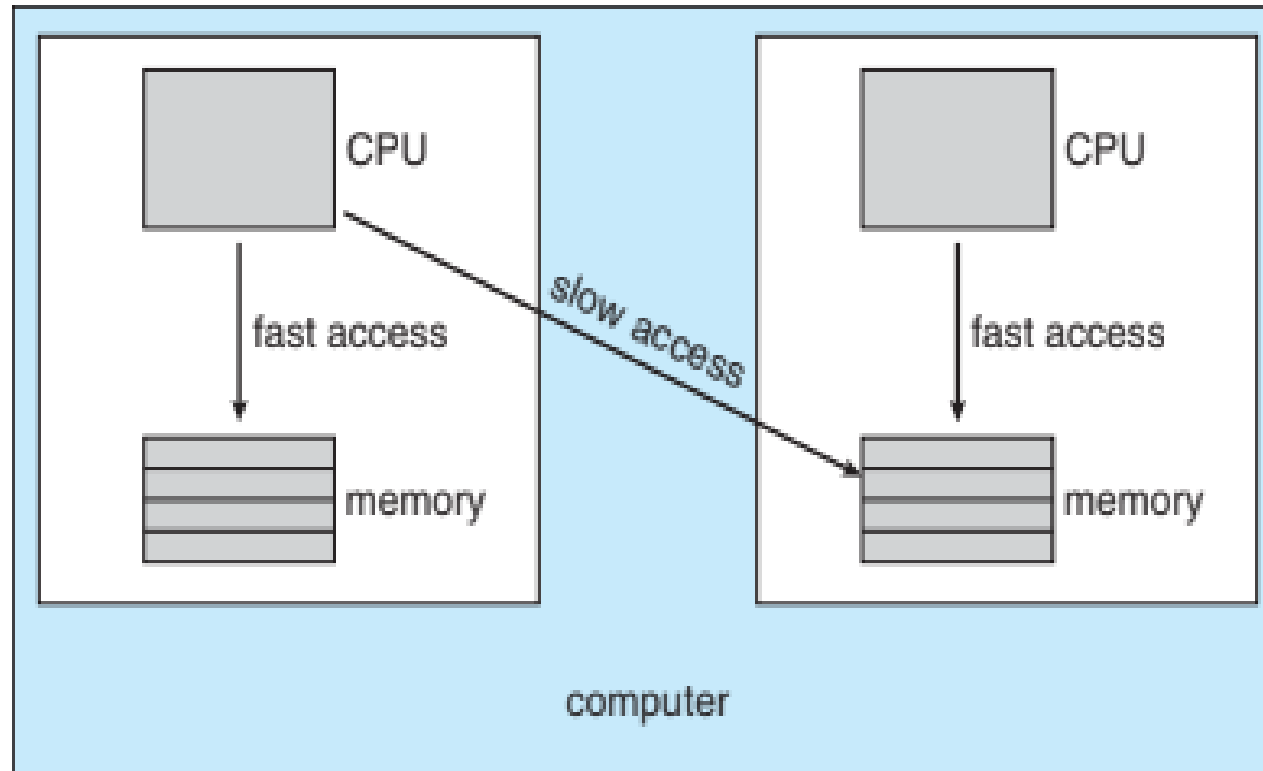# Multiple-Processor Scheduling – Processor Affinity

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

- 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, yet that thread loses the contents of what it had in the cache of the processor it was moved off of.

- **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.
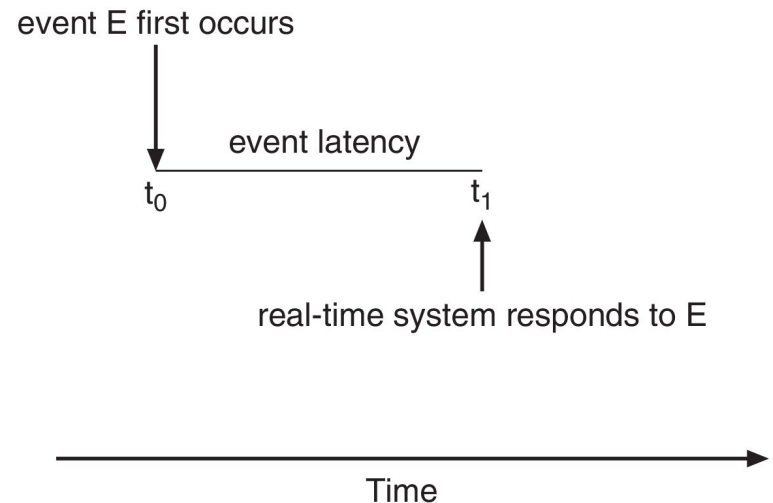
# NUMA and CPU Scheduling

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**

- 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

# Operating System Examples

- Linux scheduling
- Windows scheduling
- Solaris scheduling

# 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

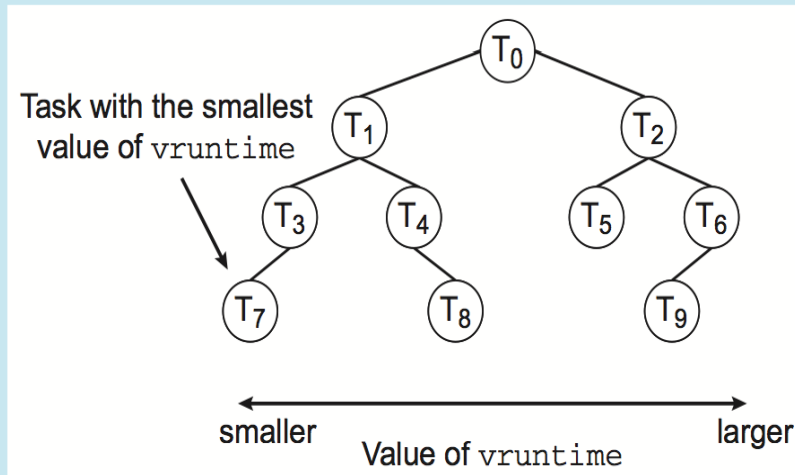# Linux Scheduling in Version 2.6.23 + (Cont.)

- 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

# 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:
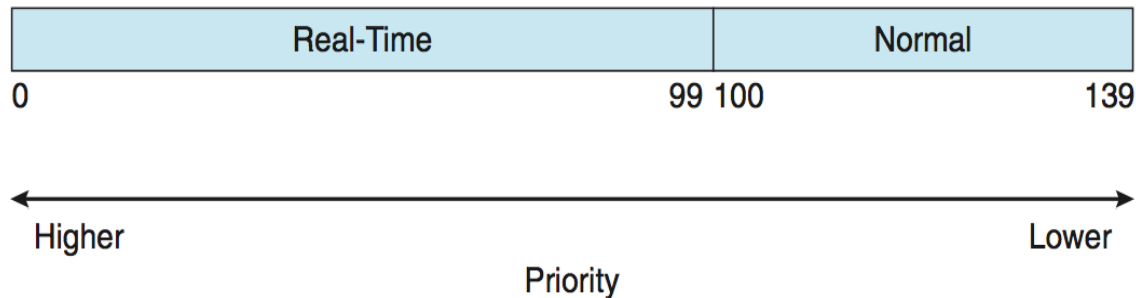


When a task becomes runnable, it is added to the tree. If a task on the tree is not runnable (for example, if it is blocked while waiting for I/O), it is removed. Generally speaking, tasks that have been given less processing time (smaller values of `vruntime`) are toward the left side of the tree, and 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.
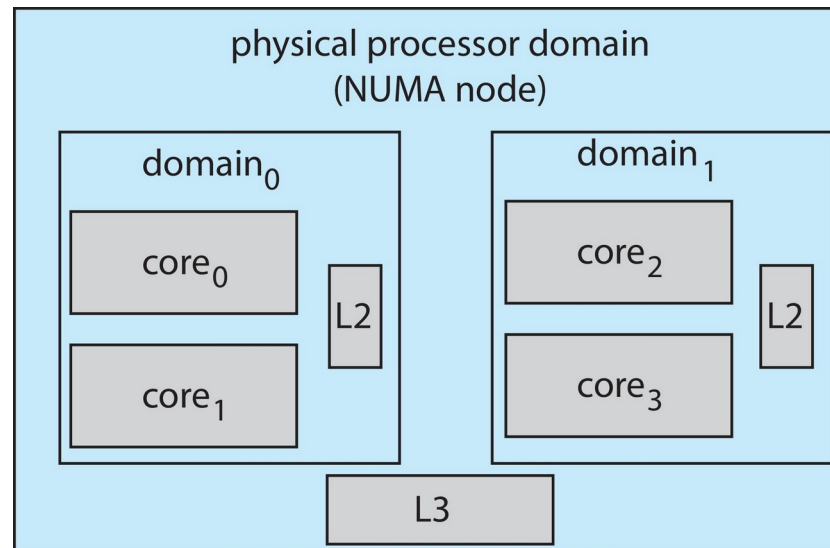
# Linux Scheduling (Cont.)

- 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

# Linux Scheduling (Cont.)

- Linux supports load balancing. It is also NUMA-aware.

- **Scheduling domain** is a set of CPU cores that can be balanced against one another.

- Domains are organized by what they share (i.e., cache memory.) Goal is to keep threads from migrating between domains.
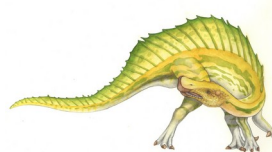
# Windows Scheduling

- Windows uses priority-based preemptive scheduling

- Highest-priority thread runs next

- **Dispatcher** is scheduler

- Thread runs until (1) blocks, (2) uses time slice, (3) preempted by higher-priority thread

- Real-time threads can preempt non-real-time

- 32-level priority scheme

- **Variable class** is 1-15, **real-time class** is 16-31

- Priority 0 is memory-management thread

- Queue for each priority

- If no run-able thread, runs **idle thread**

# Windows Priority Classes

- Win32 API identifies several priority classes to which a process can belong
  - REALTIME_PRIORITY_CLASS, HIGH_PRIORITY_CLASS, ABOVE_NORMAL_PRIORITY_CLASS,NORMAL_PRIORITY_CLASS, BELOW_NORMAL_PRIORITY_CLASS, IDLE_PRIORITY_CLASS
  - All are variable except REALTIME
- A thread within a given priority class has a relative priority
  - TIME_CRITICAL, HIGHEST, ABOVE_NORMAL, NORMAL, BELOW_NORMAL, LOWEST, IDLE
- Priority class and relative priority combine to give numeric priority
- Base priority is NORMAL within the class
- If quantum expires, priority lowered, but never below base

# Windows Priority Classes (Cont.)

- If wait occurs, priority boosted depending on what was waited for

- Foreground window given 3x priority boost

- Windows 7 added **user-mode scheduling** (**UMS**)

  - Applications create and manage threads independent of kernel

  - For large number of threads, much more efficient

  - UMS schedulers come from programming language libraries like C++ **Concurrent Runtime** (ConcRT) framework

|  | 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 |

# End of Chapter 5