# Operating Systems

Isfahan University of Technology
Electrical and Computer Engineering Department
1400-1 semester

Zeinab Zali

Session 11: CPU Scheduling Algorithms
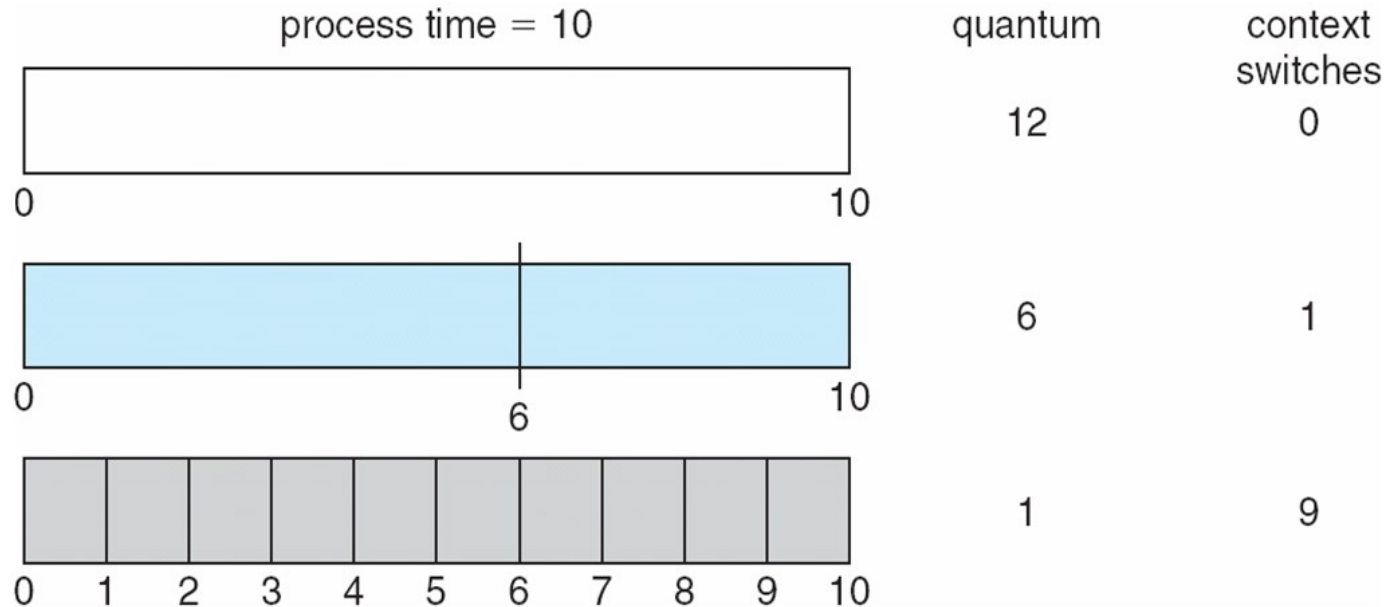
# Round Robin (RR)

- Each process gets a small unit of CPU time (**time quantum** $q$), usually 10-100 milliseconds. After this time has elapsed, the process is <span style="color:red">preempted</span> 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.

- **What is the maximum response time?**

  - **No process waits more than ($n$-1)$q$ time units.**

- Timer interrupts every quantum to schedule next process

- Performance

  - $q$ large $\Rightarrow$ FIFO

  - $q$ small $\Rightarrow$ $q$ must be large with respect to context switch, otherwise overhead is too high

# Time Quantum and Context Switch Time



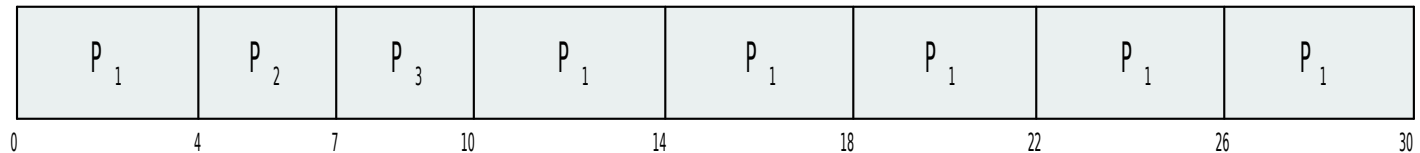q usually 10ms to 100ms, context switch < 10 usec

# 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

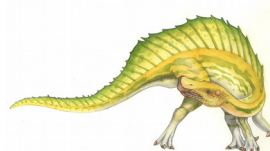■ Compare average waiting time, turnaround time and response time with SJF

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

# RR Properties

- Preemptive

- No starvation

- Suitable for time sharing systems

- Low throughput with short quantum

- Same as FCFS with a quantum larger than CPU bursts

# Priority Scheduling

- A priority number (integer) is associated with each process

- The CPU is allocated to the process with the 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
  - When shut downing IBM 7094 at MIT in 1973, there is a waiting job from 1967!

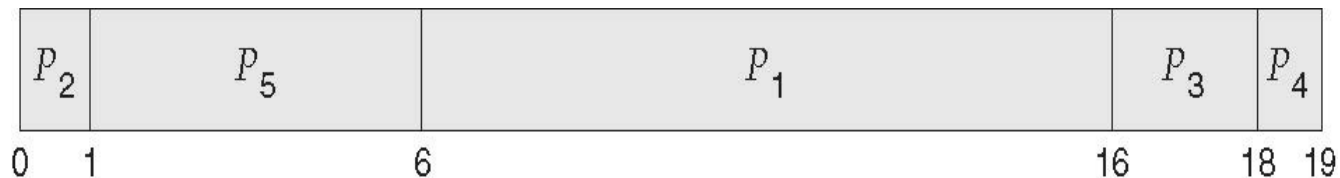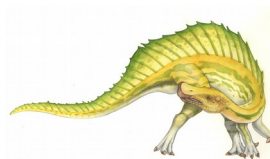- Solution ≡ **Aging** – as time progresses increase the priority of the process

# 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

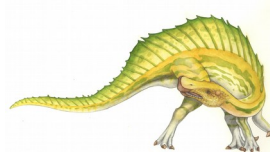| $P_2$ | $P_5$ | $P_1$ | $P_3$ | $P_4$ |
|-------|-------|-------|-------|-------|

0   1        6                          16      18  19

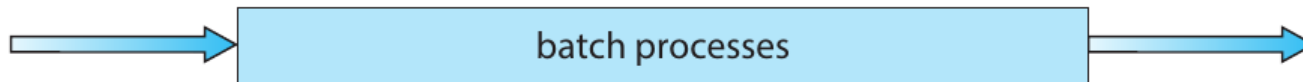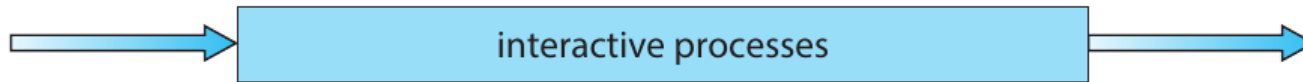- Average waiting time = 8.2 msec
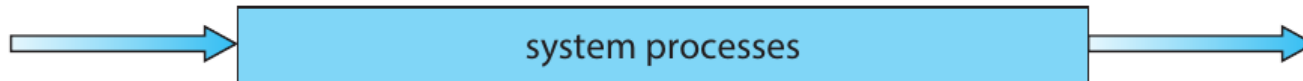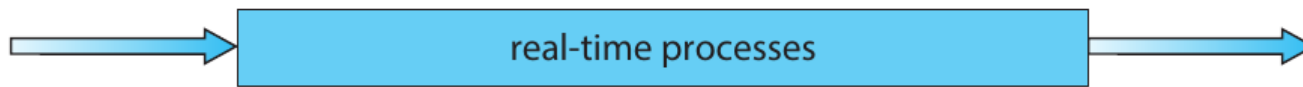
# Multilevel Queue

- Ready queue is partitioned into separate queues, eg:
  - **foreground** (interactive)
  - **background** (batch)
- Process permanently in a given queue
- 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 among its processes; i.e., 80% to foreground in RR, 20% to background in FCFS
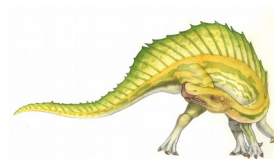
# Multilevel Queue Scheduling

highest priority



lowest priority

# Multilevel Feedback Queue

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

- Multilevel-feedback-queue scheduler defined by the following parameters:

  - number of queues

  - scheduling algorithms for each queue

  - method used to determine when to upgrade a process

  - method used to determine when to demote a process

  - method used to determine which queue a process will enter when that process needs service
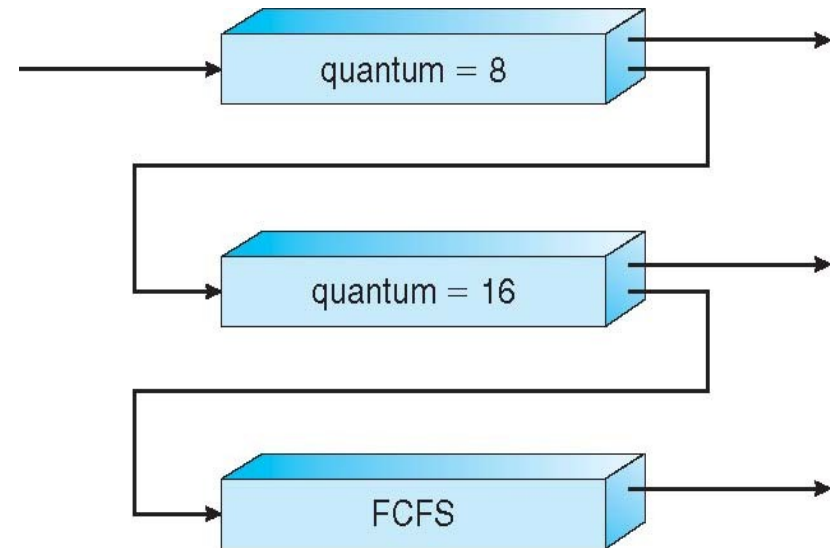
# Example of Multilevel Feedback Queue

- Three queues:
  - $Q_0$ – RR with time quantum 8 milliseconds
  - $Q_1$ – RR time quantum 16 milliseconds
  - $Q_2$ – FCFS

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



quantum = 8

quantum = 16

FCFS

# Real-Time CPU Scheduling

- Can present obvious challenges

- **Soft real-time systems** – no guarantee as to when critical real-time process will be scheduled

- **Hard real-time systems** – task must be serviced by its deadline

- Real-time characteristics:
  - Periodic (p value)
  - Processing time (t)
  - Deadline (d)

- Some algorithms:
  - Priority-based

Linux command to check process priority and scheduling
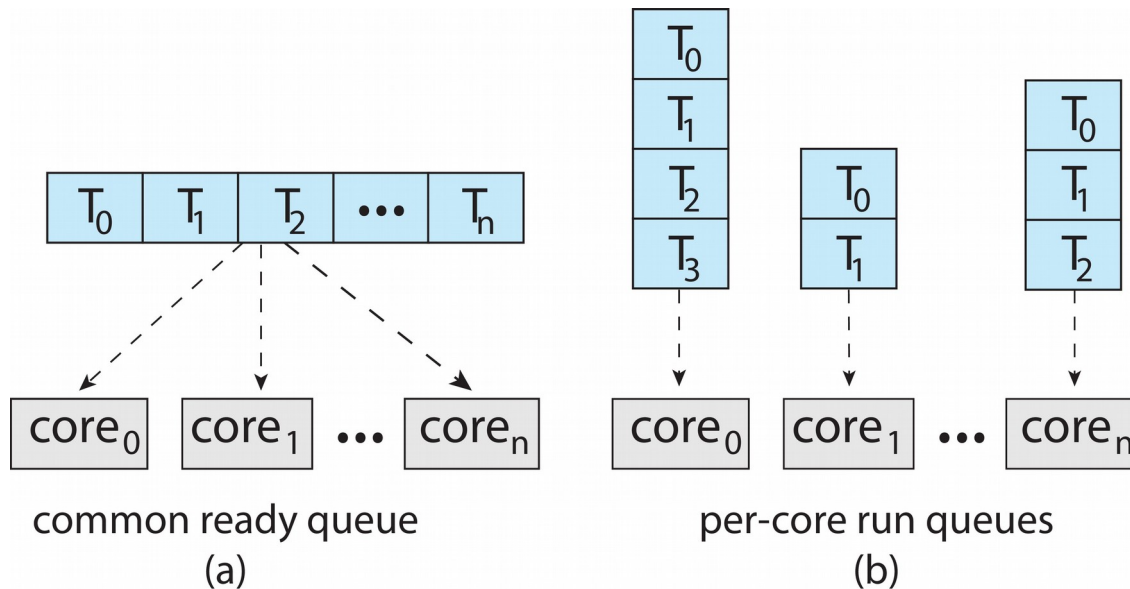**chrt**

# Multiple-Processor Scheduling

- CPU scheduling more complex when multiple CPUs are available

- **Asymmetric multiprocessing** – only one processor accesses the system data structures, alleviating the need for data sharing

- **Symmetric multiprocessing** (**SMP**) – each processor is self-scheduling, all processes in common ready queue, or each has its own private queue of ready processes

  - Currently, most common

- **Processor affinity** – process has affinity for processor on which it is currently running

  - **soft affinity**

  - **hard affinity**
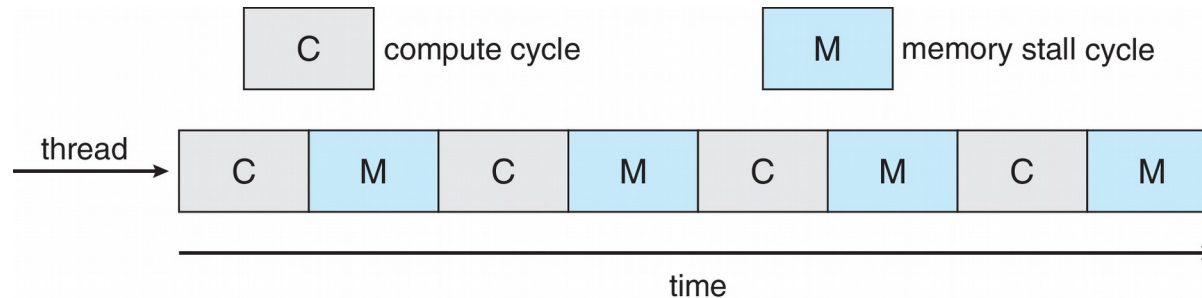
# Multiple-Processor Scheduling

- 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
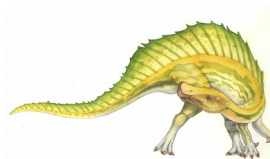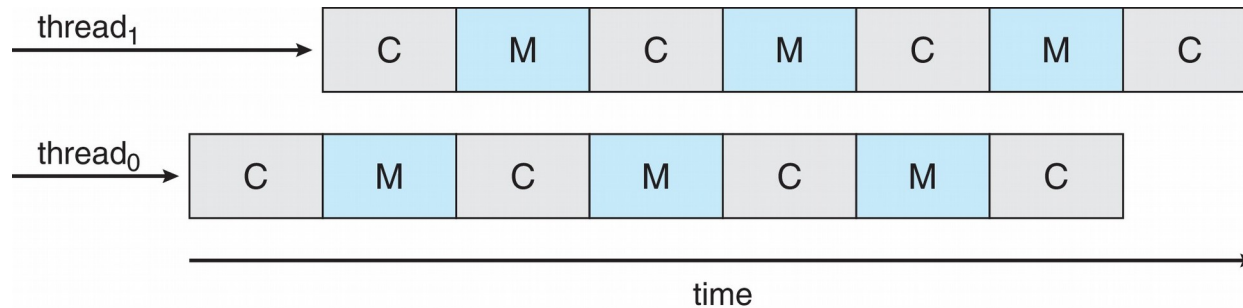(b)

# 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
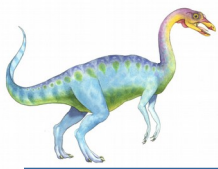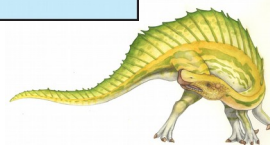
- Figure

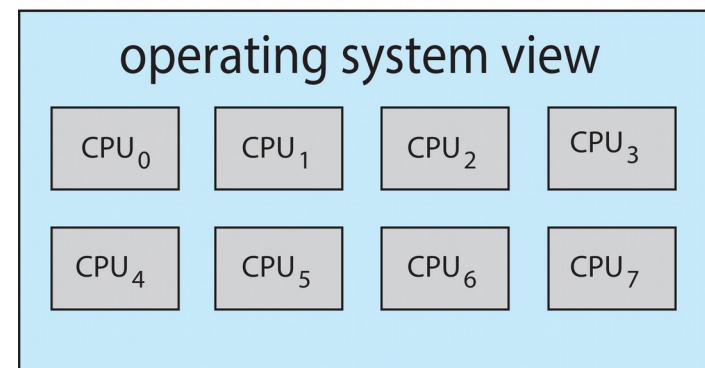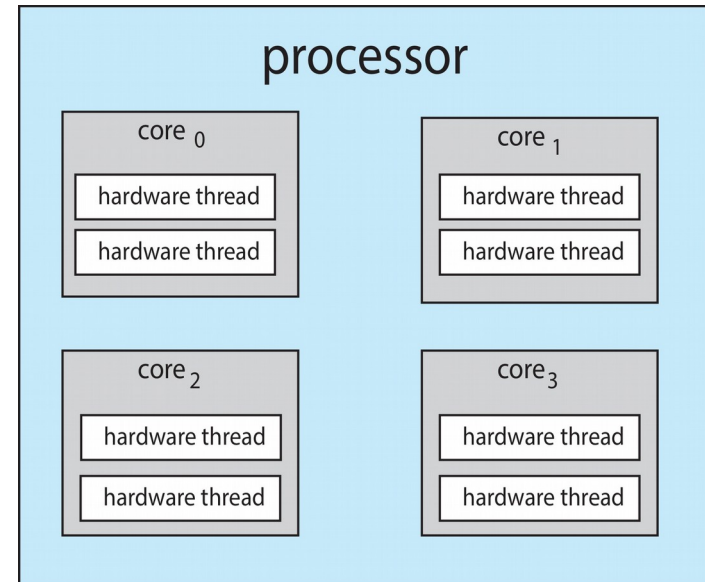# Multithreaded Multicore System

- 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**.)

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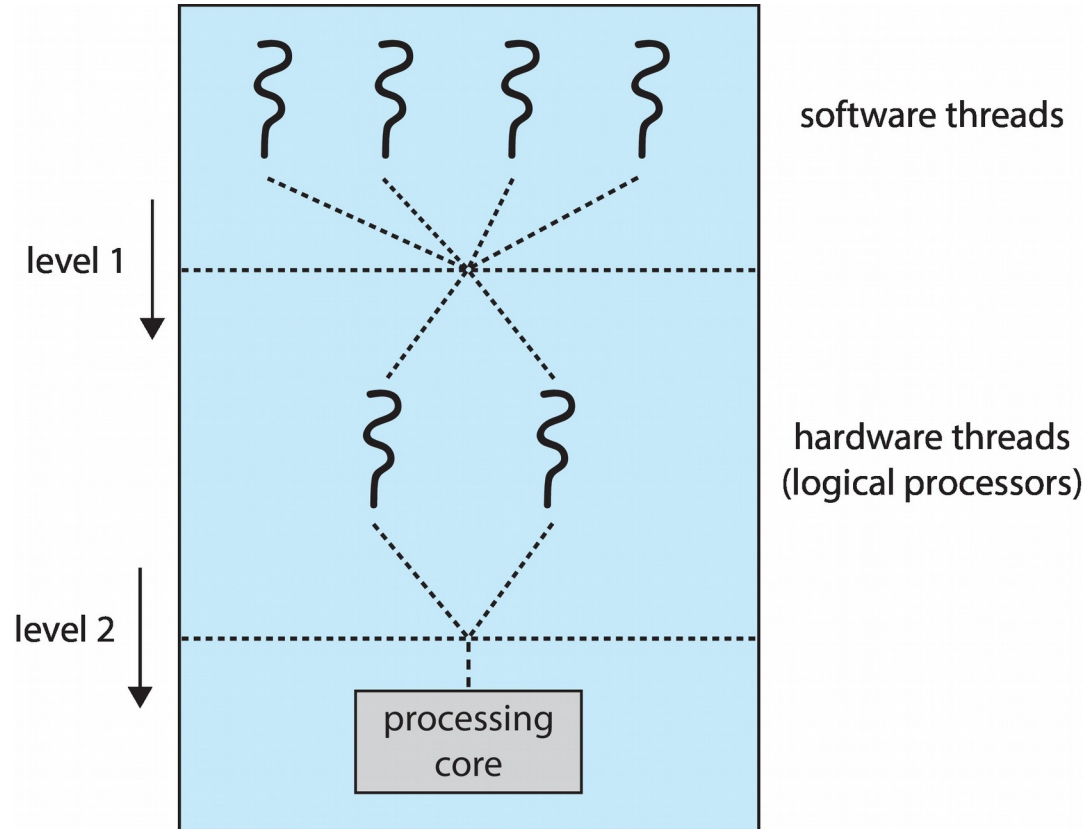
# Multithreaded Multicore System

- 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.

software threads

level 1

hardware threads
(logical processors)

level 2

processing
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

# 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.

Linux command to check process affinity
**taskset**

# Thread Scheduling

- Distinction between user-level and kernel-level threads

- **When threads supported, threads scheduled, not processes**

- Many-to-one and many-to-many models, thread library schedules user-level threads to run on LWP

  - Known as **process-contention scope** (**PCS**) since scheduling competition is within the process

  - Typically done via priority set by programmer

- Kernel thread scheduled onto available CPU is **system-contention scope** (**SCS**) – competition among all threads in system
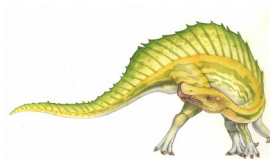
# Pthread Scheduling

- API allows specifying either PCS or SCS during thread creation

  - **PTHREAD_SCOPE_PROCESS** schedules threads using PCS scheduling

  - **PTHREAD_SCOPE_SYSTEM** schedules threads using SCS scheduling

- Can be limited by OS – Linux and Mac OS X only allow PTHREAD_SCOPE_SYSTEM

- Scheduling policies

  - FIFO

  - RR

  - OTHER

# Linux scheduling

- pthread_attr_init(&attr)
- Config Scope
  - pthread_attr_getscope(&attr,&scope)

    pthread_attr_setscope(&attr,PTHREAD_SCOPE_PROCESS)

- Config scheduling algorithm
  - pthread_attr_getschedpolicy(&attr,&policy)

    pthread_attr_setschedpolicy(&attr,SCHED_FIFO)

- Shell commands
  - chrt

# Linux Scheduling Through Version 2.5

- Prior to kernel version 2.5, ran variation of standard UNIX scheduling algorithm

- Version 2.5 moved to constant order $O(1)$ scheduling time

  - Preemptive, priority based

  - Two priority ranges: time-sharing and real-time

  - **Real-time** range from 0 to 99 and **nice** value from 100 to 140

  - **Higher priority gets larger q**

  - Worked well, but poor response times for interactive processes
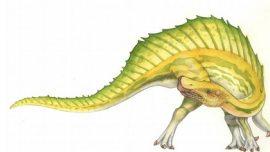
# Linux Scheduling in Version 2.6.23 +
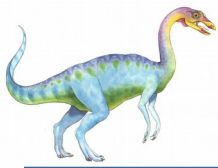
- ***Completely Fair Scheduler*** (CFS)

- **Scheduling classes**
  - 2 scheduling classes included, others can be added
    1. default
    2. real-time

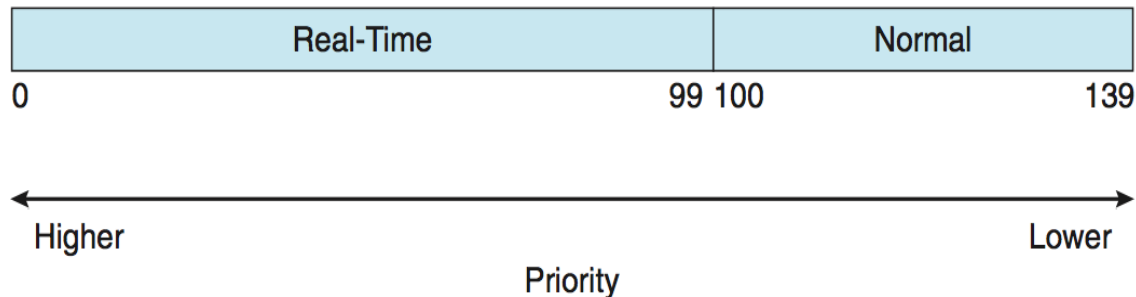- Quantum calculated based on **nice value** from -20 to +19
  - Lower value is higher priority
  - Calculates **target latency** – interval of time during which task should run at least once
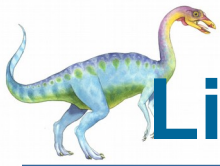  - Target latency can increase if say number of active tasks increases

# 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
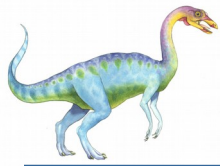- Nice value of +19 maps to priority 139

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

Higher ←——————————————————————————→ Lower
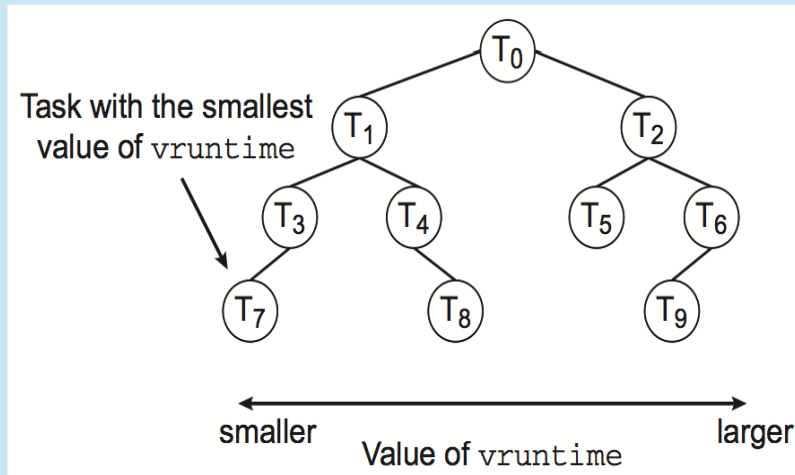
Priority

# Linux Scheduling in Version 2.6.23 +

- 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

  - For runtime= 200, what is **virtual runtime** for normal, high, low priority?

- To decide next task to run, scheduler picks task with lowest virtual run time

- Which one has higher priority?

  - IO bound or CPU bound process?

# 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.

# Algorithm Evaluation

- How to select CPU-scheduling algorithm for an OS?

- Determine criteria, then evaluate algorithms

- **Deterministic modeling**

- **Queueing Models**

- **Simulations**

- **Implementation**