# Operating Systems

## 07. Process Scheduling

Paul Krzyzanowski

Rutgers University

Spring 2015

# Running more than one process

- ## Batch systems
  - Run one job. When it finishes, run the next one, …

- ## Cooperative multitasking
  - Run a process until it makes a system call
    $\Rightarrow$ transfers control to the OS
  - OS can then decide to context switch and run another process

- ## Preemptive multitasking
  - OS programs a timer to generate an interrupt
  - Interrupt gives control back to the OS
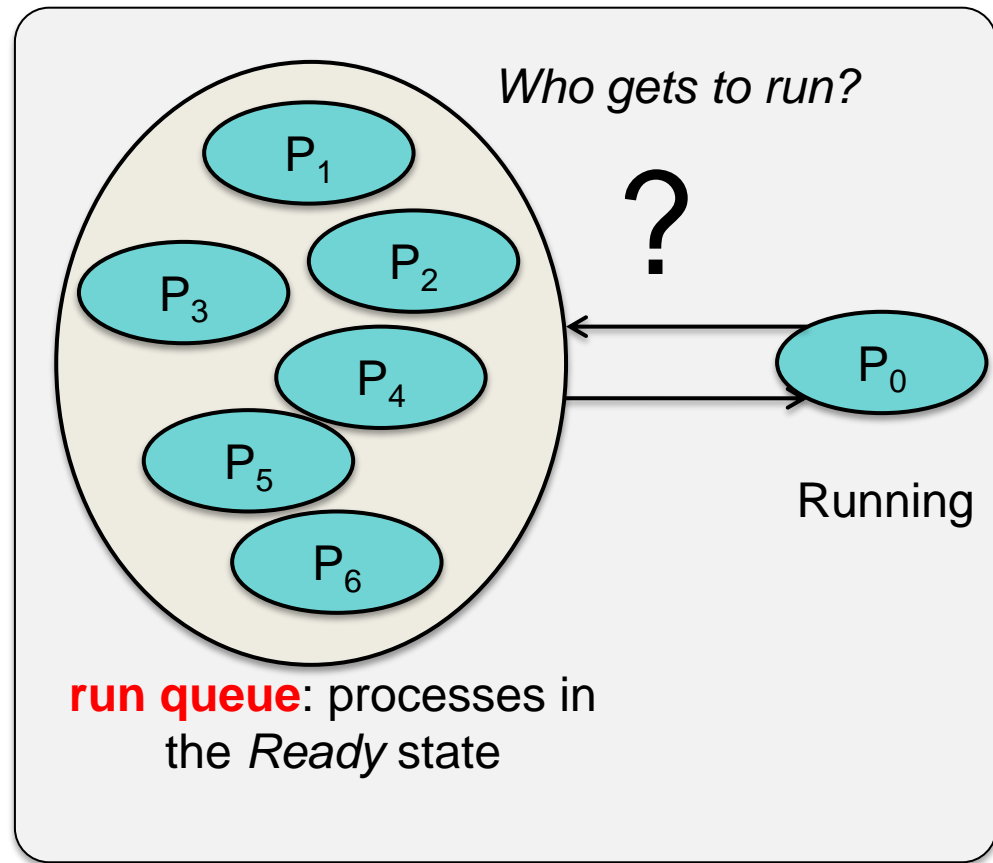    $\Rightarrow$ decides whether to context switch

# Process Scheduler

We have multiple tasks *ready* to run.
Which one should get to run?

## Scheduling algorithm:
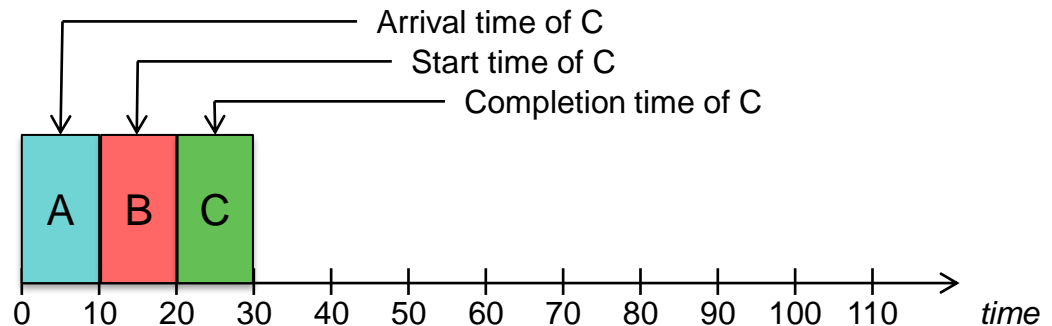
– **Policy**: Makes the decision of who gets to run

## Dispatcher:

– **Mechanism** to do the context switch



*Who gets to run?*

?

Running

**run queue**: processes in the *Ready* state

# First Come, First Served (FCFS)

- Run jobs to completion in the order they arrive

- Sounds fair?

Arrival time of C
Start time of C
Completion time of C

| A | B | C |

0  10  20  30  40  50  60  70  80  90  100  110    *time*

- Turnaround time: Time to complete a job since submitting it

- Turnaround time $= T_{completion} - T_{arrival}$
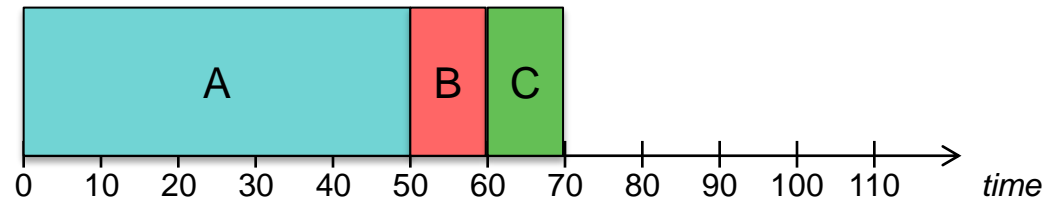
- Assume A, B, & C arrive at around the same time

$T_{turnaround}(A) = 10$

$T_{turnaround}(B) = 20$

$T_{turnaround}(C) = 30$

$T_{turnaround}(average) = (10+20+30) \div 3 = 20$

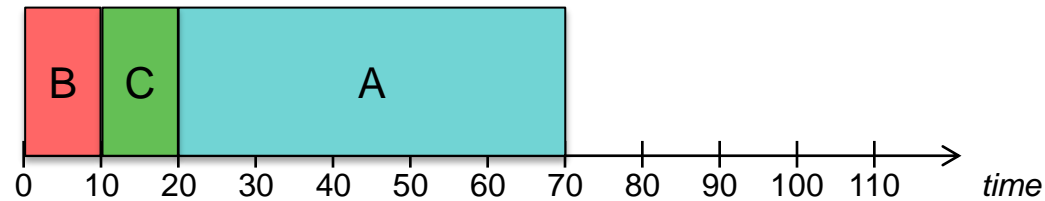# First Come, First Served



- What if A was a long-running job?

$$T_{turnaround}(A) = 50$$

$$T_{turnaround}(B) = 60$$

$$T_{turnaround}(C) = 70$$

$$T_{turnaround}(average) = (50+60+70) \div 3 = 60$$

# Shortest Job First (SJF)



- Let shortest jobs run first ⇒ optimizes turnaround time

$$T_{turnaround}(B) = 10$$

$$T_{turnaround}(C) = 20$$

$$T_{turnaround}(A) = 70$$
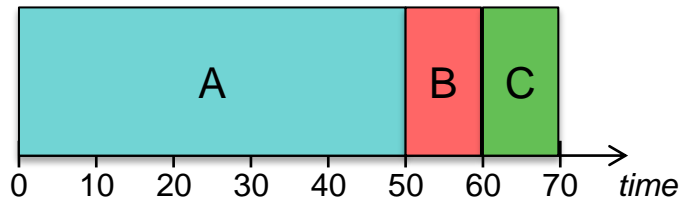
$$T_{turnaround}(average) = (10+20+70) \div 3 = 33.333 \; vs. \; 60$$

- 1.8x better than FCFS! (in this example)

- But if *B* and *C* arrive a bit after *A*, we're still out of luck

# Response time

- FCFS and SJF: non-preemptive schedulers
  - One job hold up all others!

- Let's consider response time
  - Response time = delay before a job starts to run
  - *Response time = $T_{arrival} - T_{run}$*



Average response time =
= (0 + 50 + 60)÷3 = 36.67

Average response time =
= (0 + 10 + 20)÷3 = 10

# Round Robin

- Let's add preemption
  - Let a job run for some time (time slice = quantum)
  - Then context switch and give someone else a turn



- If quantum = 2.5:
  - average response time = (0+2.5+5) ÷ 3 = 2.5 ⇒ Great!
  - average turnaround time = (70+27.5+30) ÷ 3 = 42.5
    - worse than SJF (33.3) but better than worst-case FCFS (60)
    - In general, Round Robin is not good for turnaround time

# Time slice (quantum) length

- <u>Short quantum</u>: increases overhead % of context switching

- <u>Long quantum</u>: reduces interactivity
  - Tasks are allowed to run longer before a context switch is forced
  - <u>Amortizes</u> overhead of context switch

- No perfect answer
  - Servers: higher emphasis on efficiency
    - Use a longer quantum to reduce overhead of context switches
    - But still need interactivity to schedule I/O and provide decent response
  - Interactive systems: higher emphasis on fast user response
    - Use a shorter quantum to have more context switches

- But…
  - Interactive and I/O-bound tasks rarely will use up their time slice

# What about I/O?

We ignored I/O so far

Most tasks fall into one of two categories:
1. Large # of short CPU bursts between I/O requests
2. Small # of long CPU bursts between I/O requests

| CPU | I/O | CPU | I/O | CPU | I/O | CPU | I/O |
|-----|-----|-----|-----|-----|-----|-----|-----|

*CPU Burst*

CPU idle

# Task Behavior

Interactive task: mostly short CPU bursts

| CPU | I/O | CPU | I/O | CPU | I/O | CPU | I/O |
|-----|-----|-----|-----|-----|-----|-----|-----|

CPU idle          CPU idle          CPU idle          CPU idle

Compute task: mostly long CPU bursts

| CPU | I/O | CPU | I/O |
|-----|-----|-----|-----|

CPU idle                                CPU idle

# Task Scheduling With I/O

Goal:
- Maximize use of CPU & improve throughput
- Let another task run when the current one is waiting on I/O



**Think of each CPU burst as an individual job that needs to be scheduled**

# Process Scheduling With a Mix of Processes

Improve CPU utilization (increase chance of CPU being busy)

- Some processes will use long stretches of CPU time
  - Preempt them periodically and let another process run
- More processes than CPUs: keep them in the *ready* list
- Perhaps all processes are waiting on I/O: nothing to run!

# When does the scheduler make decisions?

Four events may cause the scheduler to get called:

1. Current process goes from *running* to *blocked* state
2. Current process terminates
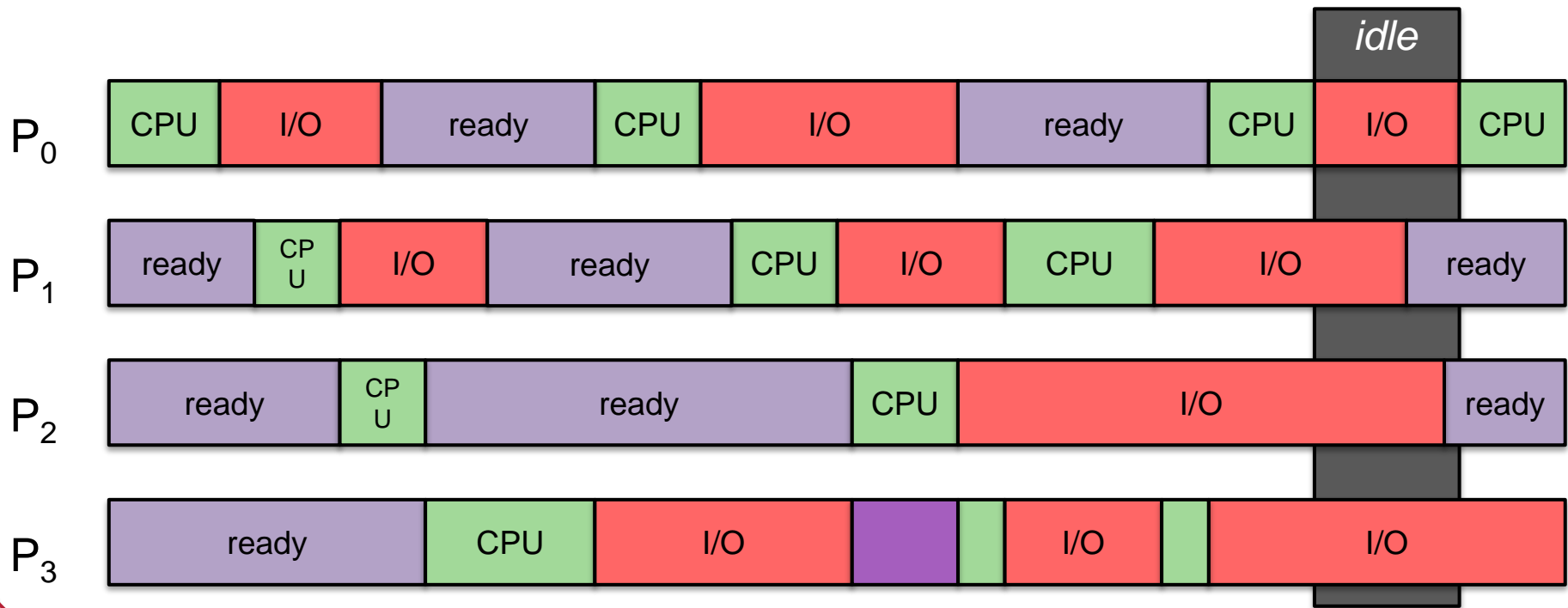3. Interrupt gives the scheduler a chance to move a process from *running* to *ready: scheduler decides it's time for someone else to run*
4. Current process goes from *blocked* to *ready*
   I/O is complete (including blocking events, such as semaphores)
   *This does not necessarily mean the currently running process will change*

- Preemptive scheduler

- Cooperative (non-preemptive) scheduler
  – CPU cannot be taken away unless a system call takes place or process exits

- Run-to-completion scheduler (old batch systems)

# Scheduling algorithm goals

| | |
|---|---|
| Be fair | (to processes? To users?) |
| Be efficient | Keep CPU busy … and don't spend a lot of time deciding! |
| Maximize throughput | Get as many processes to complete as quickly as possible |
| Minimize response time | Minimize time users must wait |
| Be predictable | Tasks should take about the same time to run & responsiveness should be similar when run multiple times |
| Minimize overhead | |
| Maximize resource use | Try to keep devices busy! |
| Avoid starvation | |
| Enforce priorities | |
| Degrade gracefully | |

# First-Come, First-Served (FCFS)



Ready queue

Running

$P_9 \rightarrow P_7 \rightarrow P_3 \rightarrow P_1 \rightarrow P_0 \rightarrow$ Completed

Waiting → Ready

$P_4$ $P_6$ $P_8$ $P_1$

$P_5$ $P_2$

Block on I/O

Waiting (blocked) queues

# First-Come, First-Served (FCFS)

- Non-preemptive

- A process with a long CPU burst will hold up other processes
  - I/O bound tasks may have completed I/O and are ready to run: poor device utilization
  - Poor average response time

# Round-Robin Scheduling

Preemptive Scheduling:
A Process can not run for longer than its assigned quantum (time slice)

*Ready queue*

*Running*

$P_9$ → $P_7$ → $P_3$ → $P_1$ → $P_0$ → *Completed*

*preempted*

*Waiting → Ready*

*Block on I/O*

$P_4$   $P_6$   $P_8$   $P_1$

$P_5$   $P_2$

*Waiting (blocked) queues*

# Round-Robin Scheduling

- Behavior depends on the quantum
  - <u>Long quantum</u> makes this similar to FCFS
  - <u>Short quantum</u> increases interactivity but increases the overhead % of context switching

- Advantages
  - Every process gets an equal share of the CPU
  - Easy to implement
  - Easy to compute average response time: $f$(# processes on list)

- Disadvantage
  - Giving every process an equal share isn't necessarily good
  - Highly interactive processes will get scheduled the same as CPU-bound processes

# Shortest Remaining Time First Scheduling

- Sort tasks by anticipated CPU burst time

- Schedule shortest ones first

- Optimize average response time

| Burst time | 2 | 2 | 10 | 3 | 8 | Total time = 25 |
|---|---|---|---|---|---|---|
| Process queue | E | D | C | B | A | |
| Total run time | 25 | 23 | 21 | 11 | 8 | Mean time = 17.6 |

*last* ⟵———————— *first*

| Burst time | 10 | 8 | 3 | 2 | 2 | Total time = 25 |
|---|---|---|---|---|---|---|
| Process queue | C | A | B | D | E | |
| Total run time | 25 | 15 | 7 | 4 | 2 | Mean time = 10.6 |

Mean completion time for a process falls by almost 40%!

# Shortest Remaining Time First Scheduling

- Biggest problem: *we're optimizing with data we don't have!*

- All we can do is estimate

- Exponential average – estimate of next CPU burst:

  $$e_{n+1} = \alpha t_n + (1 - \alpha)e_n$$

  average of all previous CPU bursts

  time of current CPU burst

  $\alpha$ is a weight factor to balance the weight of the last burst period vs. historic periods ($0 \leq \alpha \leq 1$)

  If $\alpha = 0$:     $e_{n+1} = e_n$   (recent history has no effect)
  If $\alpha = 1$:     $e_{n+1} = \alpha t_n$   (use only the last burst time)

# Shortest Remaining Time First Scheduling

- Advantage
  - Short-burst tasks run fast

- Disadvantages
  - Long-burst (CPU intensive) tasks get a long mean waiting time
    - Starvation risk!
  - Need to rely on ability to estimate CPU burst length

# Priority Scheduling

Round Robin assumes all processes are equally important

- Not true
  - Interactive tasks need high priority for good response
  - We might want non-interactive tasks to get the CPU less frequently: *this goal led us to SRTF*
  - Some tasks might be time critical
  - Users may have different status (e.g., administrator)

- Priority scheduling algorithm:
  - Each process has a priority number assigned to it
  - Pick the process with the highest priority
  - Processes with the same priority are scheduled round-robin

# Priority Scheduling – Assigning Priorities

- Priority assignments:
  - Internal: time limits, memory requirements, I/O:CPU ratio, …
  - External: assigned by administrators

- Static & dynamic priorities
  - Static priority: priority never changes
  - Dynamic priority: scheduler changes the priority during execution
    - Increase priority if it's I/O bound for better interactive performance or to increase device utilization
    - Decrease a priority to let lower-priority processes run
    - Example: use priorities to drive SJF/SRTF scheduling
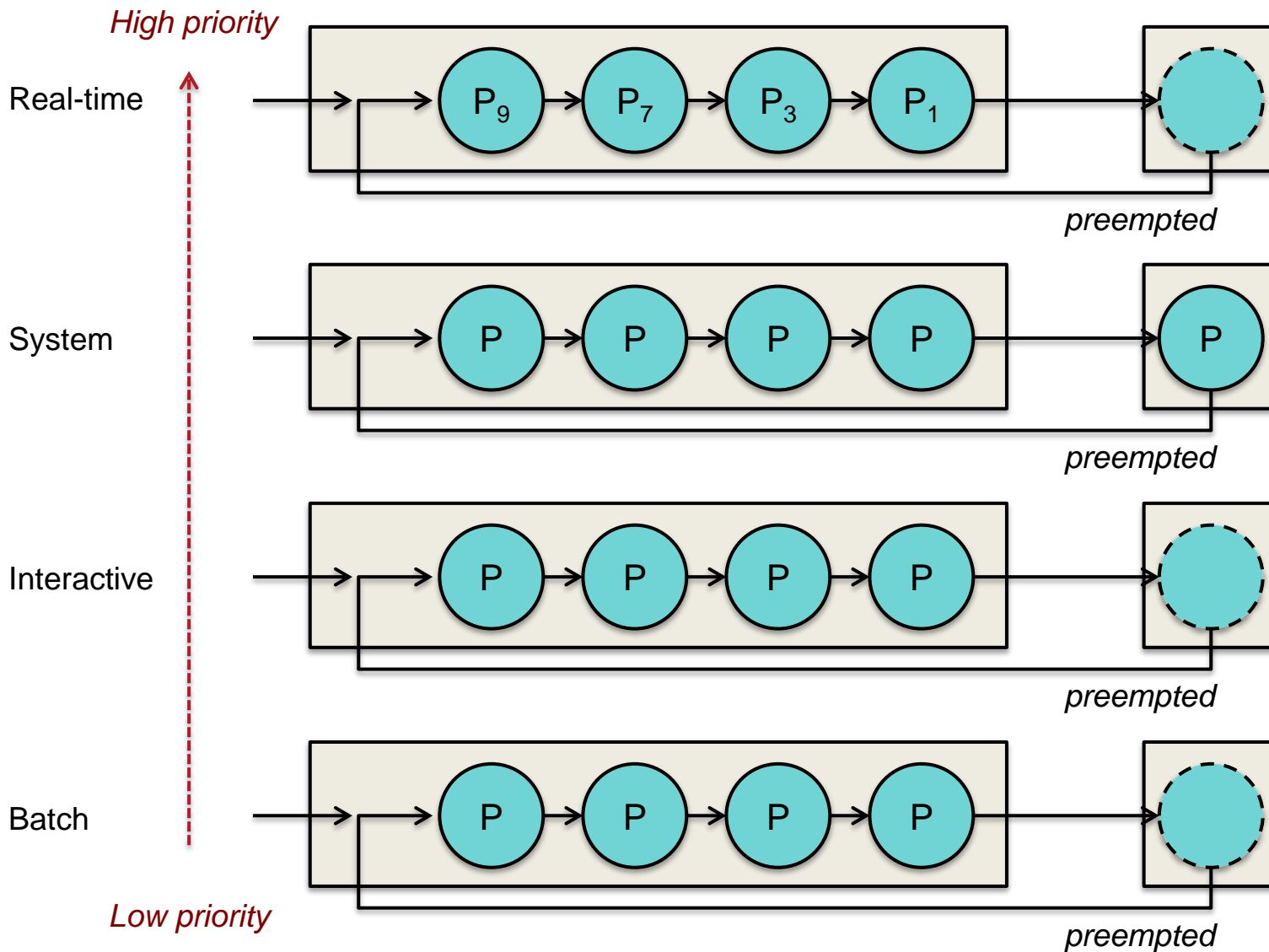
# Priority Scheduling – Problems

- ## Priority Inversion
  - A low-priority thread may not get scheduled, thereby preventing a high-priority thread that is holding a resource from making progress

- ## Starvation
  - A low priority thread may *never* get scheduled if there is always a high-priority thread ready to run

# Multilevel Queues

*Does each task need to have a unique priority level?*

- Priority classes: a ready queues for each priority level
  - Each priority class gets its own queue
  - Processes are permanently assigned to a specific queue
  - Examples: System processes, interactive processes, slow interactive processes, background non-interactive processes

- Implementation
  - Priority scheduler with queues per priority level
  - Each queue may have a different scheduling algorithm (usually round-robin)
  - Quantum may be increased at each lower priority level
    - Lower-priority processes tend to be compute bound

# Multilevel Queues

High priority

Real-time

$P_9$ → $P_7$ → $P_3$ → $P_1$

preempted

System

P → P → P → P → P

preempted

Interactive

P → P → P → P

preempted

Batch

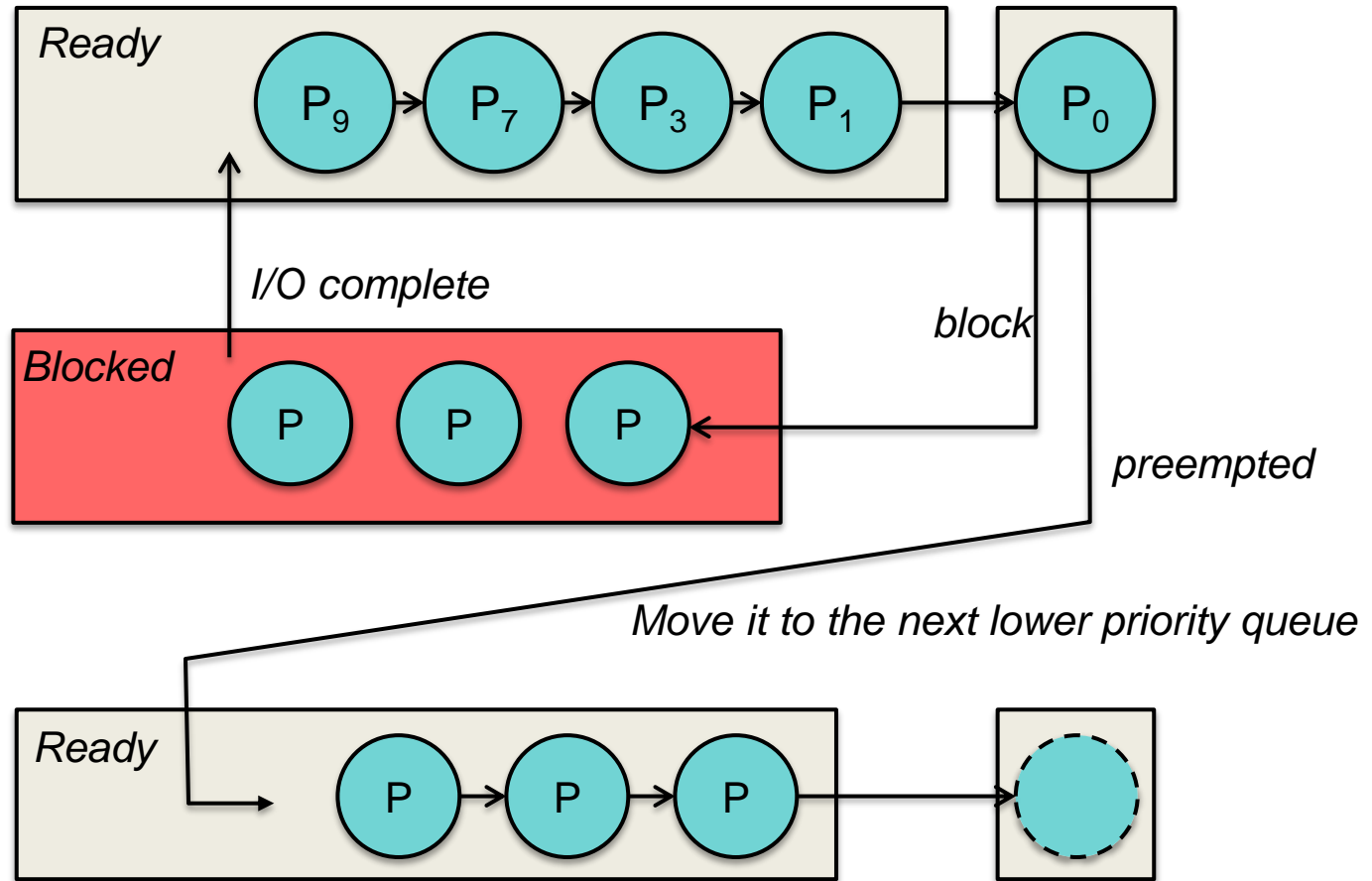P → P → P → P

preempted

Low priority

# Multilevel Feedback Queues

- ## Goals
  - Allow processes to move between priority queues based on feedback
    - Have the scheduler learn the behavior of each task and adjust priorities

- ## Separate processes based on CPU burst behavior
  - I/O-bound processes will end up on higher-priority queues

- ## Rules
  1. A new process gets the highest priority
  2. If a process does not finish its quantum (blocks on I/O) then it will stay at the same priority level (round robin) otherwise it moves to the next lower priority level

# Multilevel Feedback Queues

Pick the process from the head of the highest priority queue

*High priority*

*Ready*

$P_9$ → $P_7$ → $P_3$ → $P_1$ → $P_0$

*I/O complete*

*block*

*Blocked*

P   P   P

*preempted*

*Move it to the next lower priority queue*

*Ready*

P → P → P

*Low priority*

# Multilevel Feedback Queues

High priority

Ready   $P_9$ → $P_7$ → $P_3$ → $P_1$ → $P_0$

preempted

Ready   P → P → P → P →

preempted

Ready   $P_3$ → $P_1$ →

Low priority

# Example

One long-running process

# Example

- Suppose a highly interactive process, *B* (■), starts at *T=10*
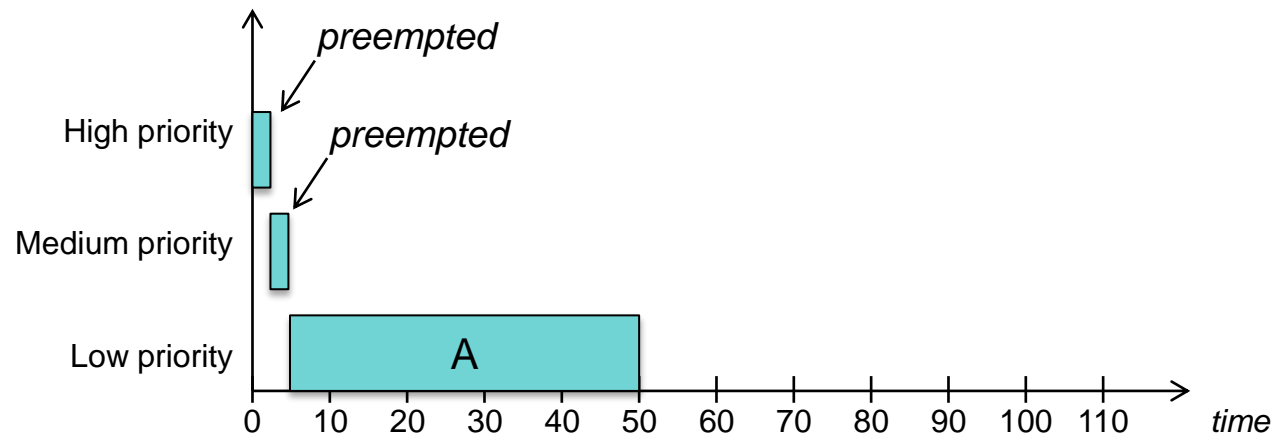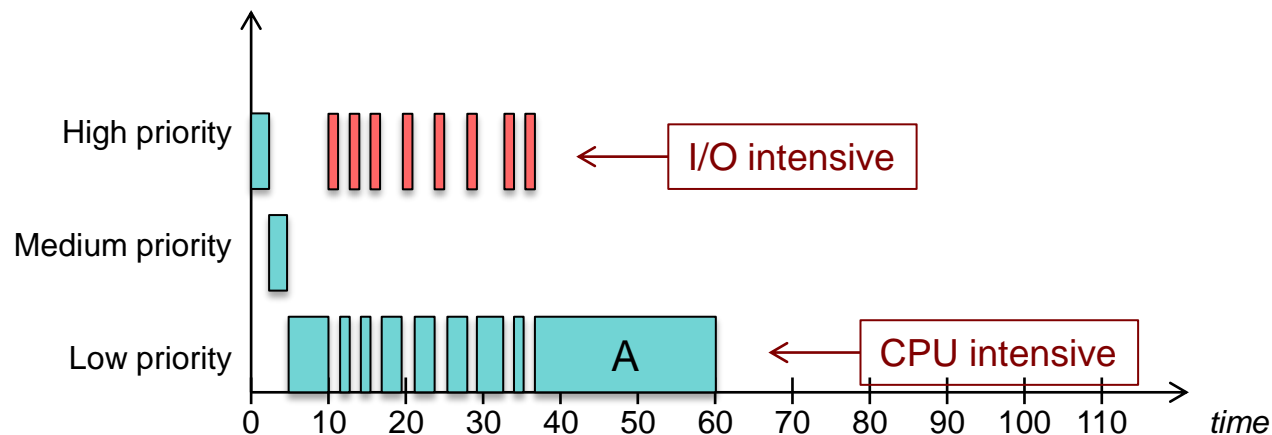- It never uses up its quantum
  - *B* gets priority but spends a lot of its time in the blocked state
  - A (■) gets to run only when *B* is blocked



High priority — I/O intensive

Medium priority

Low priority — A — CPU intensive

time: 0 10 20 30 40 50 60 70 80 90 100 110

# Starvation & aging

- Two problems
  - Starvation:
    If there are a lot of interactive processes, the CPU-bound processes will never get to run
  - Interactive process ending up at a low priority:
    If a process was CPU intensive (e.g., initializing a game) but then became interactive, it is forever doomed to a low priority


- Solve these process aging
  - Increase the priority of a so it will be scheduled to run
    - Simplest approach: periodically, set *all* processes to the highest priority
  - If it remains CPU-intensive, its priority will quickly fall again

# Windows Scheduler

- Two classes:
  - Variable class: priorities 0-15
  - Real-time class: priorities 16-31

- Each priority level has a queue
  - Pick the highest priority thread that is ready to run

- Relative priority
  - Threads have relative levels within their class
  - When a quantum expires, the thread's priority is lowered but never below the base
  - When a thread wakes from wait, the priority is increased
    - Higher increase if waiting for keyboard input
  - Priority is increased for foreground window processes

# Windows Priorities

| | Real-time | High | Above Normal | Normal | Below Normal | Idle |
|---|---|---|---|---|---|---|
| **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 |

# Linux Schedulers – History

- Linux 1.2: Round Robin scheduler (fast & simple)

- Linux 2.2: Scheduling classes (multilevel queue)
  - Classes: Real-time, non-real-time, non-preemptible
  - Basic support for symmetric multiprocessing

# Linux 2.4: O(N) Scheduler

- Multilevel queue with two scheduling algorithms:

- (1) Real-time with absolute priorities (but kernel is not preemptible)
  - FIFO & Round-robin options

- (2) Time-sharing: Credit-based algorithm
  - Each task has some # of credits associated with it
  - On each timer interrupt:
    - Each timer interrupt: running task loses 1 credit
    - If credits for a task == 0, the task is suspended
    - If all tasks have 0 credits:
      - Re-credit: Every task gets credits = credits/2 + priority
    - Choose next task to run: pick the one with the most credits

- Not good for systems with many tasks
  - Re-crediting requires going through every task: O(N)

- Not good for multiprocessor systems
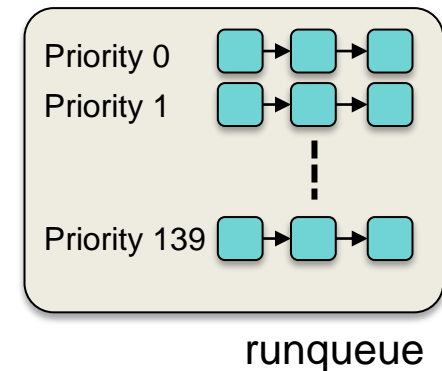  - One queue (in a mutex): contention & no processor affinity

# Linux 2.6: O(1) scheduler goals

Addressed three problems

– Scalability: O(1) instead of O(n) to not suffer under load

– Support processor affinity

– Support preemption in the kernel
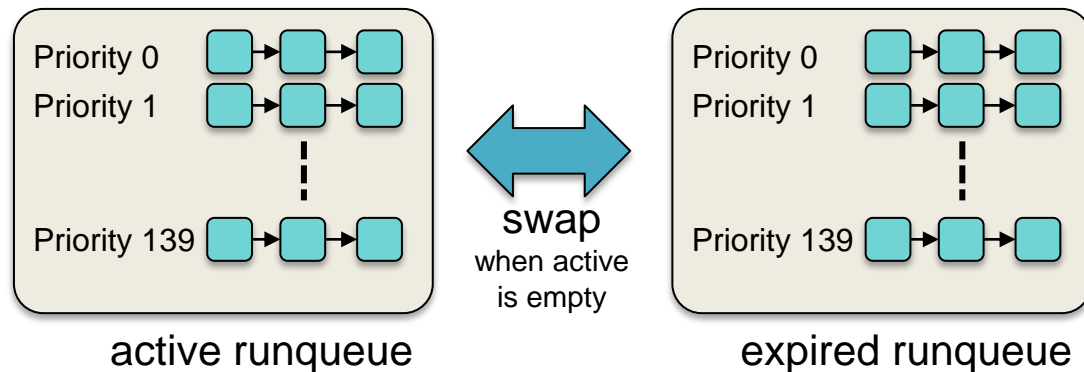  • High-priority (real-time) tasks can interrupt a task running in kernel mode

# Linux 2.6: O(1) scheduler

- O(1) instead of O(N): no increased overhead with more tasks

- One run queue per CPU

- Always schedule the highest priority task
  - Multiple tasks at same priority scheduled round-robin

- Multilevel queue
  - 140 queues (priority levels)
    - 0-99 = real-time
    - 100-139 = timesharing – dynamic priorities
  - Each priority level has its own time slice
    - Higher priority = LONGER time slice

Priority 0
Priority 1
Priority 139

runqueue

# Linux 2.6: O(1) scheduler

- Two sets of queues: active & expired
  - Epoch = time when all runnable tasks begin with a fresh time slice
  - When a task uses up its time slice, it is moved to the expired list
  - When there are no runnable tasks in the active list
    - Active & expired lists are swapped: end of epoch & start of a new one
    - Simulates aging



active runqueue

swap
when active
is empty

expired runqueue

# Linux 2.6: O(1) scheduler

- Real-time tasks: static priorities
  - Choice of round-robin or FIFO

- Non real-time tasks: dynamic priorities → reward interactive tasks
  - I/O-bound processes get priority increased by up to 5 levels
  - CPU-bound processes get priority decreased up to 5 levels
  - "Interactivity credits": +credit for sleeping, -credit for running

- SMP load balancing
  - Every 200ms, check if CPU loads are unbalanced
  - If so, move tasks from a loaded CPU to a less-loaded one
  - If a CPU's runqueue is empty, move from another CPU's runqueue

- Downside of O(1) scheduler
  - A lot of code with complex heuristics
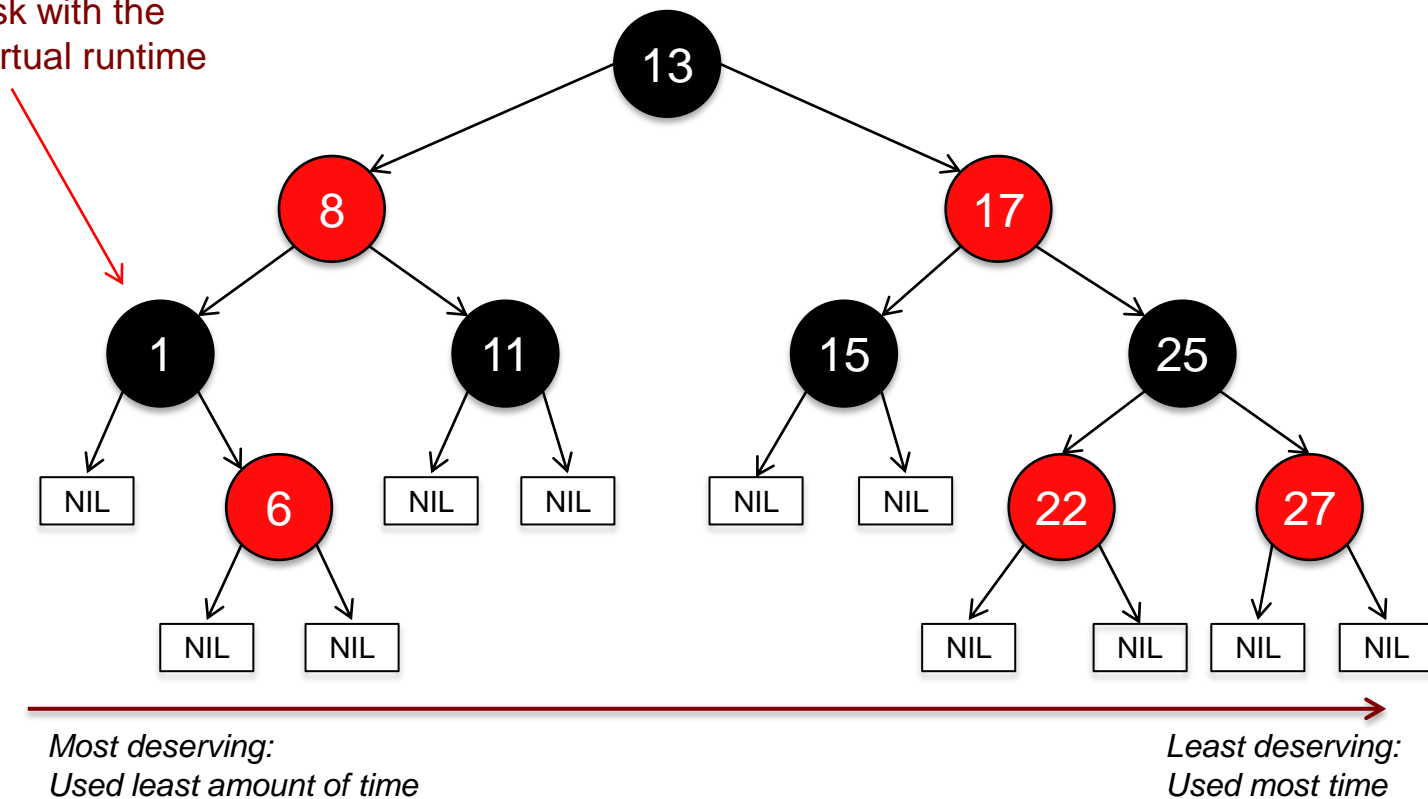
# Linux Completely Fair Scheduler

- Latest scheduler (introduced in 2.6.23)

- Goal: give a "fair" amount of CPU time to tasks

- Keep track of time given to a task: "virtual runtime"

- Basic heuristic: tasks get a fair % of the processor
  - But interactive processors are unlikely to use their share
  - When an interactive task wakes up, the scheduler sees that it used less than its fair share. To try to be fair, it preempts a compute-bound task

- Priorities – affect the rate of "virtual runtime"
  - High priority task's *vruntime* grows slower than the *vruntime* of a low priority task

# Linux Completely Fair Scheduler

No run queues: virtual runtime sorted red-black tree used instead
- – Self-balancing binary tree: search, insert, & delete in O(log n)

Left-most node always
has the task with the
smallest virtual runtime

```
                              13

              8                            17

         1         11              15              25

    NIL      6   NIL   NIL     NIL     NIL     22          27

           NIL   NIL                     NIL   NIL   NIL   NIL
```

*Most deserving:*
*Used least amount of time*

*Least deserving:*
*Used most time*

From: http://en.wikipedia.org/wiki/File:Red-black_tree_example.svg

# The End