

Scheduling

Instructor: Linyi Li
Slides adapted from Dr. B. Fraser

Topics

- 1) Computers seem **magically** able to do **more than one thing at once**.
 - a) How do they have **multiple programs** running "at once"?
 - b) How can **multiple users** log into a computer at once?
- 2) **How does the kernel decide who's turn it is to use the CPU?**

The story so far...

- “*In the beginning*”... CPUs had a **single core** and **one program** running.
- Then “*back in the day*...” computers had a **single core** but **many users**.
 - Each user might have a terminal and want to run programs
 - **How do they share the same CPU?**
- “*These days*...” CPUs have **many cores**, but..
 - **many more processes than cores.**
- **Things that need to run:**
 - ..

```
4proc 1filter
Tree:
[-]-1 systemd (/lib/systemd/systemd --system --descri
  | 161617 fwupd
  | 675 vmtoolsd
  | 161501 anacron
  | [-]-647 avahi-daemon (avahi-daemon:)
  |   | 668 avahi-daemon (avahi-daemon:)
  |   | 676 NetworkManager
  |   | 1004 rtkit-daemon
  |   | 651 dbus-daemon
  |   | 759 ModemManager
  |   | 380 systemd-journal (systemd-journald)
  |   | 677 wpa_supplicant
  |   | 158386 cups-browsed
  |   | 458 systemd-timesyn (systemd-timesyncd)
  |   | 158385 cupsd
  |   | 415 systemd-udev
  |   | 24155 systemd-network (systemd-networkd)
  |   | 663 systemd-logind
  |   | 656 polkitd (/usr/lib/polkit-1/polkitd --no-def
  |   | 121892 rpcbind
  |   | 665 udiskd (/usr/libexec/udisks2/udiskd)
  |   | 1769 geoclue
  |   | 648 cron (/usr/sbin/cron -f)
  | [-]-1514 systemd
  |   | [-]-1658 gnome-shell
  |   |   | [-]-4812 firefox-esr
  |   |   |   | 25734 Isolated Web Co (firefox-esr)
  |   |   |   | 25693 Isolated Web Co (firefox-esr)
  |   |   |   | 4932 Privileged Cont (firefox-esr)
  |   |   |   | 72553 Isolated Web Co (firefox-esr)
  |   |   |   | 72589 Isolated Web Co (firefox-esr)
  |   |   |   | 25697 Isolated Web Co (firefox-esr)
  |   |   |   | 5063 Isolated Web Co (firefox-esr)
  |   |   |   | 157734 Web Content (firefox-esr)
  |   |   |   | 5069 Isolated Web Co (firefox-esr)
```

More Depth

- We will cover scheduling a little to understand the problem.
 - CMPT301 teaches it in depth
- Can read more in OSTEP
 - <https://pages.cs.wisc.edu/~remzi/OSTEP/>
Has in-depth discussions beyond scope of this course.
 - Chapter 7 Scheduling: Introduction
<https://pages.cs.wisc.edu/~remzi/OSTEP/cpu-sched.pdf>
 - Chapter 8 Scheduling: The Multi-Level Feedback Queue
<https://pages.cs.wisc.edu/~remzi/OSTEP/cpu-sched-mlfq.pdf>
 - Chapter 9.7 The Linux Completely Fair Scheduler (CFS)
<https://pages.cs.wisc.edu/~remzi/OSTEP/cpu-sched-lottery.pdf>

CPU Scheduling

CPU Scheduling

- CPU Scheduling

- ..

- Or, sharing multiple cores by multiple processes (beyond the scope of this course)

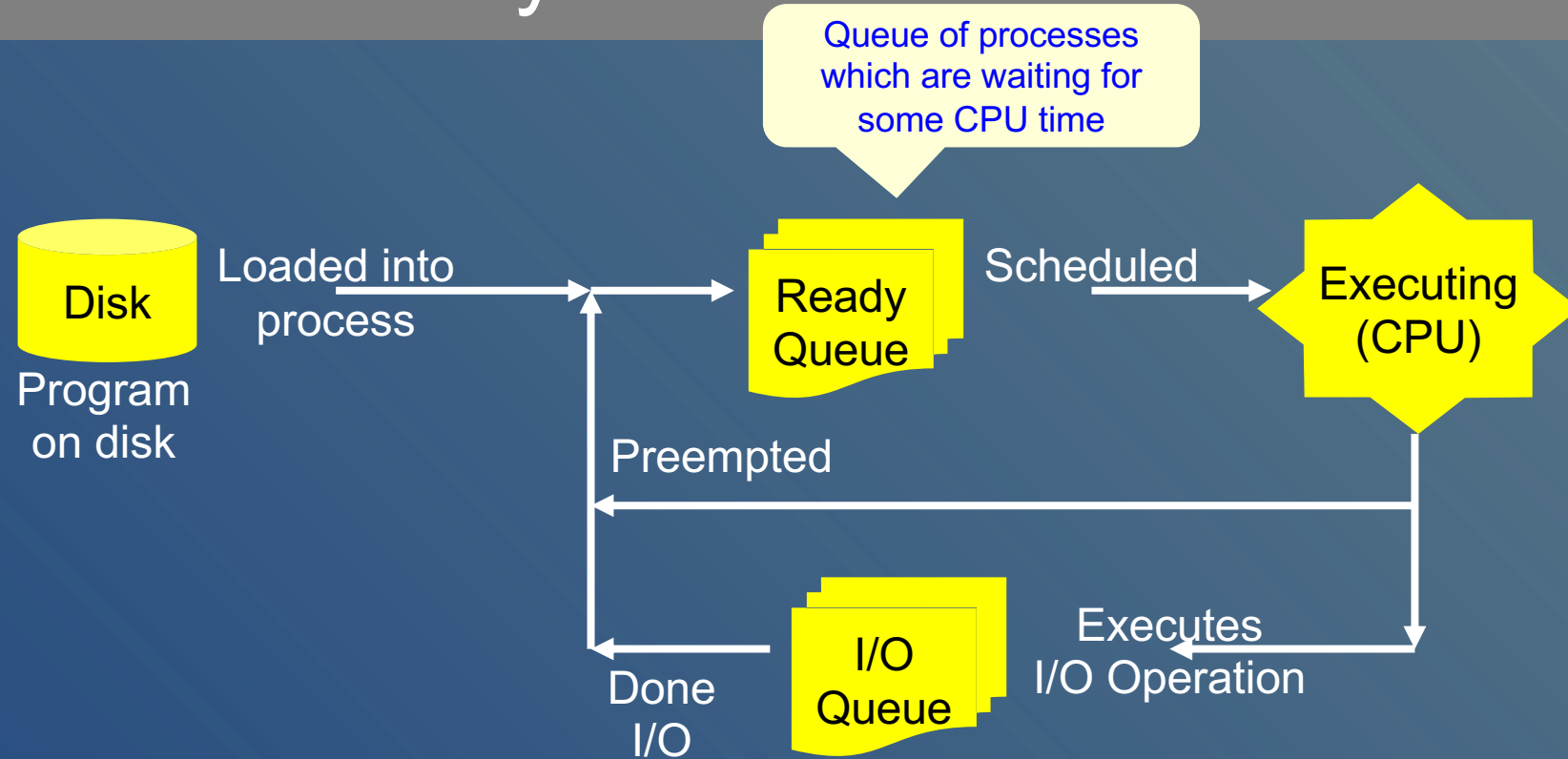
- Context switch

- ..

- This has **overhead** (work) to do this switch, so don't do it too frequently.

- Stopped process can later be **resumed exactly where it left off** once it has another turn on the CPU.

Process Lifecycle



- Scheduling

..

- Scheduler

Component of the kernel that picks the next process to run.

Types of Scheduling Algorithm

- Non-preemptive Scheduling

- A process gives up the core when it:

- ..

- ..

- e.g., I/O, child wait, sleep

- Preemptive Scheduling

- The kernel stops a *process* at any time.

- Preemptable Kernel

- (almost) all *syscalls into the kernel itself* can be preempted!

- Linux Real-time kernel (*PREEMPT_RT*) merged to mainline kernel Sept 2024!

Scheduling Criteria



- We want to maximize:

- CPU utilization : keep the CPU as busy as possible
- Throughput : # of processes that complete their execution per time unit



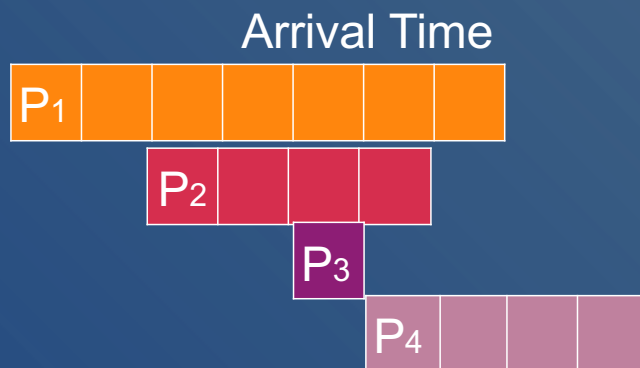
- We want to minimize:

- Turnaround time:
amount of time to execute a particular process
(time from submission to termination)
- Waiting time : amount of time a process has been waiting in ready queue
- Response time : amount of time it takes from when a request is submitted until the first response is produced

Scheduling Algorithms

Scheduling Simplifications

- Each process needs the CPU for a certain amount of time.
 - We'll **assume we know how much time it needs** at the start, but could be estimated.
- Often processes are long lived, but only need the CPU in short bursts of CPU time.
 - Here we'll just look at **one such burst of time**, but in reality a process often has many such bursts.



First Come, First Served (FCFS)

- Simplest algorithm

- ..

- ..

- (once running, a process keeps running)

- Waiting time

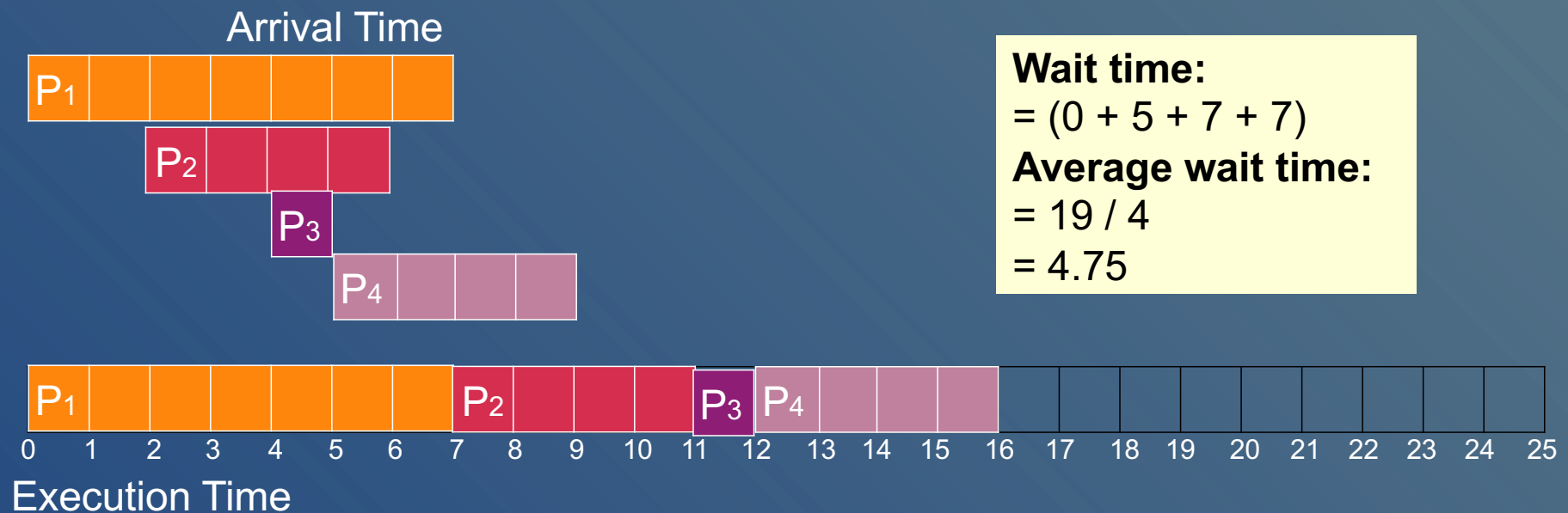
- ..

- Used to assess how good an algorithm is.

There are other metrics are based on scheduling criteria above, but waiting time is easy to calculate, so we'll use it for comparing scheduling algorithms.

First Come, First Serve Example

	P1	P2	P3	P4
Execution Time	7	4	1	4
Arrival Time	0	2	4	5



Wait time:

$$= (0 + 5 + 7 + 7)$$

Average wait time:

$$= 19 / 4$$

$$= 4.75$$

ABCD: First Come, First Served (FCFS)

- What is the **total wait time** for the following processes using **FCFS**? (if same arrival time, order by index)

	P1	P2	P3	P4
Execution Time	40	20	8	10
Arrival Time	0	0	0	0



$$(a) 40 + 20 + 8 = 68$$

$$(b) 40 + 20 + 8 + 10 = 78$$

$$(c) 40 + 60 + 68 = 168$$

$$(d) 40 + 60 + 68 + 78 = 246$$

ABCD: First Come, First Served (FCFS)

- What is the **total wait time** for the following processes using **FCFS**? (if same arrival time, order by index)

	P1	P2	P3	P4
Execution Time	10	20	8	40
Arrival Time	0	0	0	0



$$(a) 10 + 30 + 38 = 78$$

$$(b) 10 + 30 + 38 + 78 = 156$$

$$(c) 10 + 20 + 8 = 38$$

$$(d) 10 + 20 + 8 + 40 = 78$$

- What is the problem with FCFS?
 - A long process can sabotage all other processes.

Shortest Job First (SJF)

- Let's try something where **a long process doesn't sabotage** all other processes.

- **Shortest Job First Scheduling Algorithm**

– ..

– **Non-preemptive:**

Once running, a job runs to completion

– * Theoretically optimal waiting time, if all processes arrive at the beginning

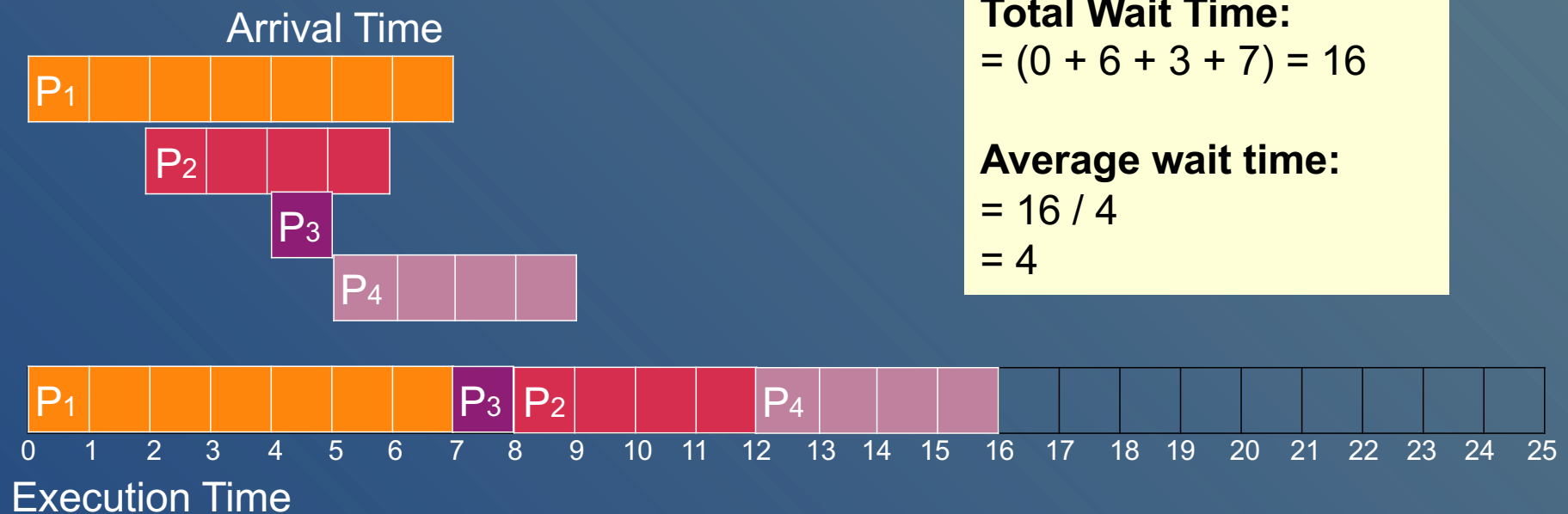
- Prove by reordering

Assume for the sake of discussion, we know how long each process takes.

Not always possible but can be estimated.

Shortest Job First Example

	P1	P2	P3	P4
Execution Time	7	4	1	4
Arrival Time	0	2	4	5



SJF is non-preemptive

Shortest

Total Wait Time:
 $= (0 + 6 + 3 + 7) = 16$

Average wait time:
 $= 16 / 4$
 $= 4$

Shortest Remaining Time First (SRTF)

- Schedule the process with..

- This is preemptive:

-

When a new job arrives,

..

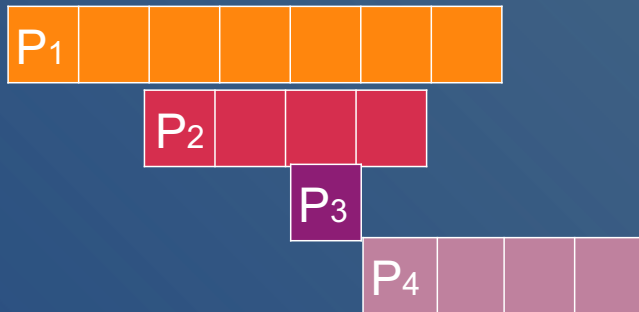
Shortest Remaining Time First Example

	P1	P2	P3	P4
Execution Time	7	4	1	4
Arrival Time	0	2	4	5

Wait Times

P1	P2	P3	P4
0+9	0+1	0	2
= 12			

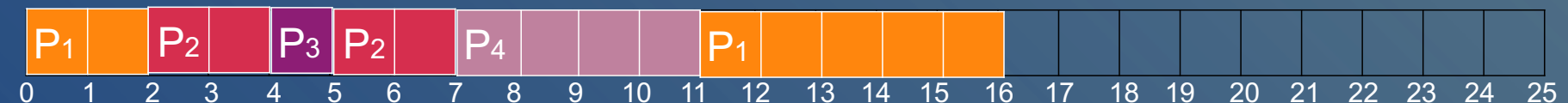
Arrival Time



Average Wait Time

$$= 12 / 4$$

$$= 3$$



Execution Time

SRTF is preemptive

Always pick shortest remaining time

More on SRTF and Wait Time Counting

- When preempted,
 - Wait counter starts;
 - So, for a process

Wait Times			
P1	P2	P3	P4
0+9	0+1	0	2
= 12			

Wait time = end timestamp – execution time – arrival timestamp

- * SRTF (Shortest Remaining Time First) minimizes waiting time
 - Why?
 - Recall SJF (Shortest Job First, the non-preemptive version) conditionally minimizes waiting time
 - SRTF is anytime SJF
 - Prove by contradiction

Round Robin (RR)

- Forget about knowing how long things take:
just give everyone..
- Preemptive
 - Quantum:
 - ..
 - Each x units of time (quantum) the scheduler will:
 - Move currently running process
to the back of the ready queue
 - Take first process in ready queue and runs it
 - Go to next process if quantum used up or current process finishes
 - Newly arrived processes go at back of ready queue

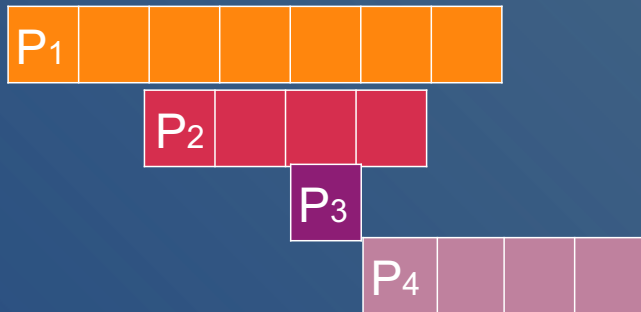
Round Robin Example (Quantum = 3ms)

	P1	P2	P3	P4
Execution Time	7	4	1	4
Arrival Time	0	2	4	5

Wait Times

P1	P2	P3	P4
$0+3+5$	$1+7$	5	$5+2$
$= 28$			

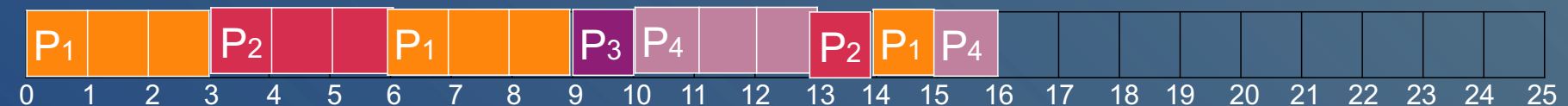
Arrival Time



Average Wait Time

$$= 28 / 4$$

$$= 7$$



Execution Time

Change every
3ms

ABCD: Round Robin

- If the **quantum is very long**, then **round robin is effectively the same as:**



- (a) First come first serve
- (b) Shortest Job First
- (c) Shortest remaining time first

- If the **quantum is very short**, what can **go wrong?**

- (a) Processes do not make progress because they keep being reset when preempted
- (b) Processes do not make progress because they keep being killed when preempted



- (c) Context switch overhead is too high
- (d) The ready queue is likely to be empty

Priority Scheduling

- Priority Scheduling

..

- This can be either preemptive or non-preemptive.

When a new process is added to the ready queue, do we allow a context switch?

- Motivation: real-time tasks with deadlines

- Some systems require hard or soft deadlines for their computational tasks
- e.g., an airplane controller must respond to an outside event (e.g., an incoming bird) within a fixed (short?) time period.

Real-Time Deadlines

- Hard real-time systems:

- ..

- because the consequence of missing a deadline can be catastrophic.

- Soft real-time systems:

- Approximate deadlines that can be missed but should not be by much.

- Real-time tasks usually have higher priorities, i.e., they are more important to run.

Priority Scheduling (cont)

- Task priority is typically expressed as a number (where a smaller number has a higher priority).
- Problem: Starvation
 - ..
 - e.g, If high priority processes keep arriving, low priority processes may never run

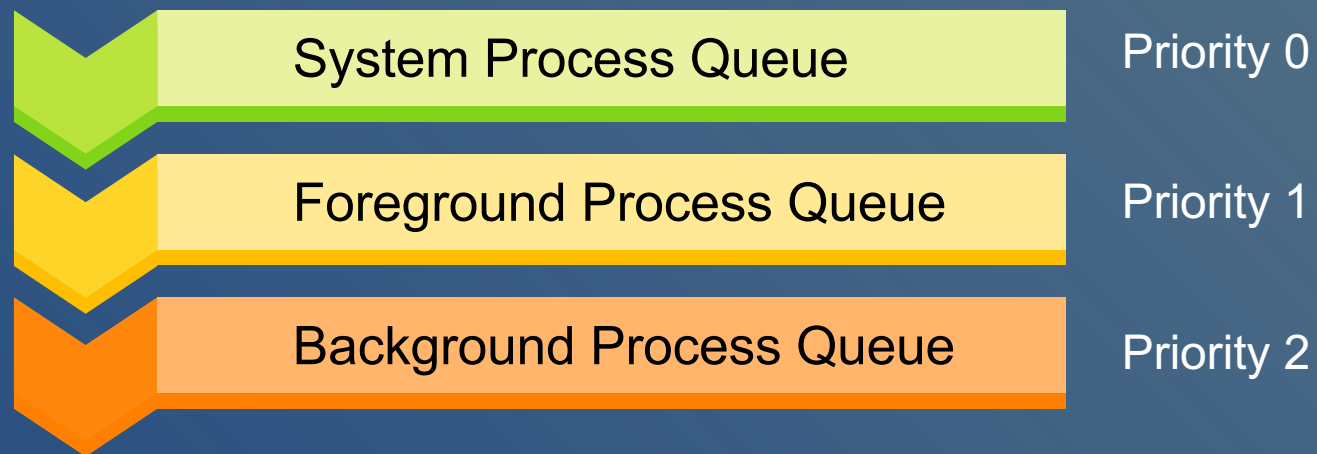
Multilevel Queue Scheduling

- Multilevel Queue Scheduling

- ..

- Each category gets its own

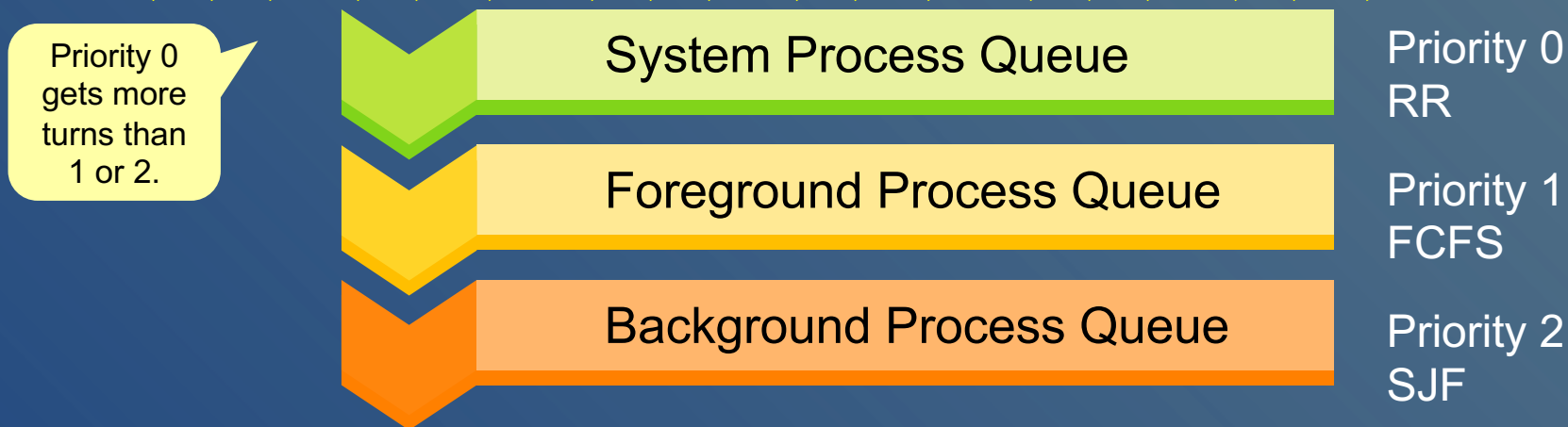
- ..



Multilevel Queue Scheduling

- Each queue gets CPU time based on priority
 - One idea: **Weighted Round Robin**
It's RR, but give more turns to higher-priority queues.
 - e.g., Schedule turns for each priority:

0, 1, 2, 0, 1, 0, 0, 1, 2, 0, 1, 0, 0, 1, 2, 0, 1, 0,



- During each queue's turn:
 - ..
- **Avoids starvation:** Each queue gets a chance to run.

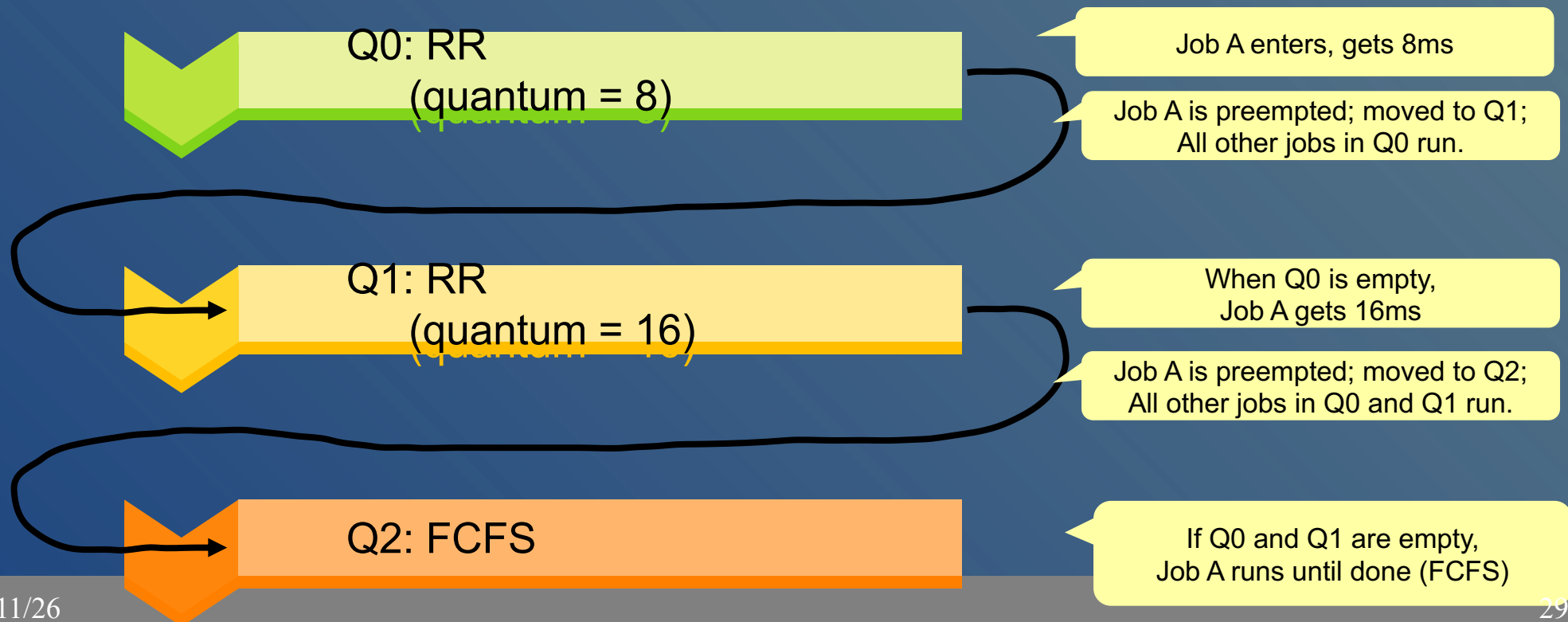
Multilevel Feedback Queue Scheduling

- Multilevel Feedback Queue

- Use multiple queues.

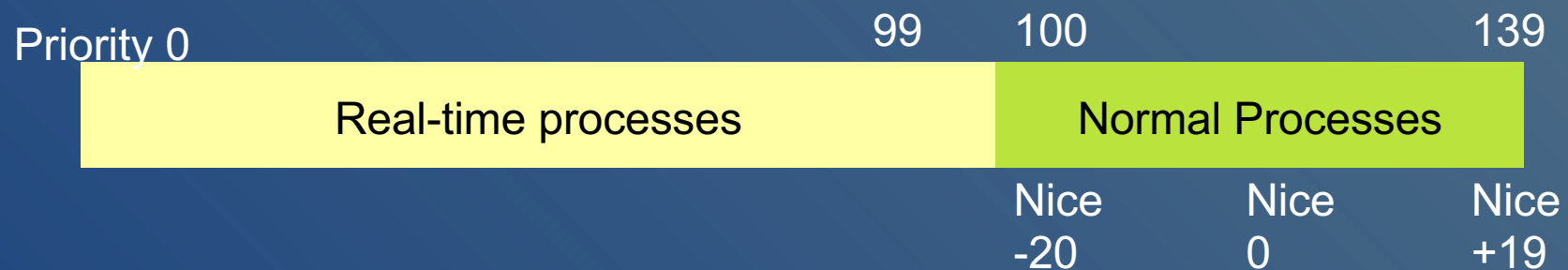
- ..

- Like Multilevel Queue, but processes lose priority via **aging**:
Lower priority by moving to lower queue if process runs too long.



Linux is Nice

- Linux categorizes processes into two classes
 - Real-time processes (priority values 0 to 99)
 - Normal processes (priority values 100 to 139)
- We can use nice value to..
 - change/assign a priority for a normal process.
 - Nice values range from -20 to +19 (lower nice == higher priority).
 - The default nice value is 0.
 - The nice -20 = priority 100, etc.



Linux Completely Fair Scheduler (CFS)

- Longer running processes get a lower priority.
 - ..
 - Older processes lose priority (aging)
- CFS tries to make sure that
 - ..
 - CFS uses **virtual run time** instead of physical (actual) run time
 - **Virtual run time** = physical run time + decay-formula
 - Higher decay with lower priority
 - **i.e.**, the decay formula produces a bigger value for a lower priority
 - *Stored internally in a red-black tree based on virtual run time*

Process Types

Interactive vs. batch

–Interactive

- Mainly user driven; regular desktop applications

–Batch

- Program runs from start to end; no interaction needed.
e.g., Compiling a program, Data analytics

I/O bound vs. CPU bound

–I/O bound

- More I/O than computation e.g., format change, such as CSV to XML

–CPU bound

- More computation than I/O e.g., compression, cryptography, etc.

–(More: Memory bound, Communication bound, ...)

Summary

- Scheduler picks what job to run next.
- Algorithms:
 - First come, first serve
 - Shortest job first
 - Shortest remaining time first
 - Round Robin
 - Multilevel queue
 - Multilevel feedback queue
 - Completely Fair Schedule
- Drawing process scheduling diagrams
 - Compute Wait time, average wait time