

ECE 353: Systems Software
Lecture 9

Advanced Scheduling

1.0.0

Jon Eyolfson
January 26, 2023



This work is licensed under a [Creative Commons Attribution-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-sa/4.0/)

We Could Add Priorities

We may favor some processes over others

Assign each process a priority

Run higher priority processes first, round-robin processes of equal priority

Can be preemptive or non-preemptive

Priorities Can Be Assigned an Integer

We can pick a lower, or higher number, to mean high priority

In Linux -20 is the highest priority, 19 is the lowest

We may lead processes to *starvation* if there's a lot of higher priority processes

One solution is to have the OS dynamically change the priority

Older processes that haven't been executed in a long time increase priority

Priority Inversion is a New Issue

We can accidentally change the priority of a low priority process to a high one

This is caused by dependencies, e.g. a high priority depends a low priority

One solution is *priority inheritance*

Inherit the highest priority of the waiting processes

Chain together multiple inheritances if needed

Revert back to the original priority after dependency

A Foreground Process Can Recieve User Input, Background Can Not

Unix background process when: process group ID differs from its terminal group ID
You do not need to know this specific definition

The idea is to separate processes that users interact with:

Foreground processes are interactable and need good response time

Background processes may not need good response time, just throughput

We Can Use Multiple Queues for Other Purposes

We could create different queues for foreground and background processes:

- Foreground uses RR

- Background uses FCFS

Now we have to schedule between queues!

- RR between the queues

- Use a priority for each queue

Scheduling Can Get Complicated

There's no "right answer", only trade-offs

We haven't talked about multiprocessor scheduling yet

We'll assume symmetric multiprocessing (SMP)

- All CPUs are connected to the same physical memory

- The CPUs have their own private cache (at least the lowest levels)

One Approach is to Use the Same Scheduling for All CPUs

There's still only one scheduler

It just keeps adding processes while there's available CPUs

Advantages

Good CPU utilization

Fair to all processes

Disadvantages

Not scalable (everything blocks on global scheduler)

Poor cache locality

This was the approach in Linux 2.4

We Can Create Per-CPU Schedulers

When there's a new process, assign it to a CPU

One strategy is to assign it to the CPU with the lowest number of processes

Advantages

Easy to implement

Scalable (there's no blocking on a resource)

Good cache locality

Disadvantages

Load imbalance

Some CPUs may have less processes, or less intensive ones

We Can Compromise between Global and Per-CPU

Keep a global scheduler that can rebalance per-CPU queues

If a CPU is idle, take a process from another CPU (work stealing)

You may want more control over which processes can switch

Some may be more sensitive to caches

Use *processor affinity*

The preference of a process to be scheduled on the same core

This is a simplified version of the O(1) scheduler in Linux 2.6

Another Strategy is “Gang” Scheduling

Multiple processes may need to be scheduled simultaneously

The scheduler on each CPU cannot be completely independent

“Gang Scheduling” (Coscheduling)

Allows you to run a set of processes simultaneously (acting as a unit)

This requires a global context-switch across all CPUs

Real-Time Scheduling is Yet Another Problem

Real-time means there are time constraints, either for a deadline or rate
e.g. audio, autopilot

A hard real-time system

Required to guarantee a task completes within a certain amount of time

A soft real-time system

Critical processes have a higher priority and the deadline is met in practice

Linux is an example of soft real-time

Linux Also Implements FCFS and RR Scheduling

You can search the source tree: FCFS (SCHED_FIFO) and RR (SCHED_RR)

Use a multilevel queue scheduler for processes with the same priority

Also let the OS dynamically adjust the priority

Soft real-time processes:

Always schedule the highest priority processes first

Normal processes:

Adjust the priority based on aging

Real-Time Processes Are Always Prioritized

The soft real-time scheduling policy will either be `SCHED_FIFO` or `SCHED_RR`

There are 100 static priority levels (0—99)

Normal scheduling policies apply to the other processes (`SCHED_NORMAL`)

By default the priority is 0

Priority ranges from $[-20, 19]$

Processes can change their own priorities with system calls:

`nice`, `sched_setscheduler`

Linux Scheduler Evolution

2.4—2.6, a $O(N)$ global queue

Simple, but poor performance with multiprocessors and many processes

2.6—2.6.22, a per-CPU run queue, $O(1)$ scheduler

Complex to get right, interactivity had issues

No guarantee of fairness

2.6.23—Present, the completely fair scheduler (CFS)

Fair, and allows for good interactivity

The $O(1)$ Scheduler Has Issues with Modern Processes

Foreground and background processes are a good division
Easier with a terminal, less so with GUI processes

Now the kernel has to detect interactive processes with heuristics
Processes that sleep a lot may be more interactive
This is ad hoc, and could be unfair

How would we introduce fairness for different priority processes?
Use different size time slices
The higher the priority, the larger the time slice
There are also situations where this ad hoc solution could be unfair

Ideal Fair Scheduling

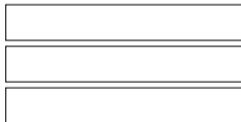
Assume you have an infinitely small time slice

If you have n processes, each runs at $\frac{1}{n}$ rate

1 Process



3 Processes



CPU usage is divided equally among every process

Example IFS Scheduling

Consider the following processes:

Process	Arrival Time	Burst Time
P ₁	0	8
P ₂	0	4
P ₃	0	16
P ₄	0	4

Assume that each vertical slice can execute 4 time units.

Each box represents the time units spend executing

	0			16		24		32
P ₁	1	2	3	4	6	8		
P ₂	1	2	3	4				
P ₃	1	2	3	4	6	8	12	16
P ₄	1	2	3	4				

IFS is the Fairest but Impractical Policy

This policy is fair, every process gets an equal amount of CPU time
Boosts interactivity, has the ideal response time

However, this would perform way too many context switches

You have to constantly scan all processes, which is $O(N)$

Completely Fair Scheduler (CFS)

For each runnable process, assign it a “virtual runtime”

At each scheduling point where the process runs for time t

Increase the virtual runtime by $t \times \text{weight}$ (based on priority)

The virtual runtime monotonically increases

Scheduler selects the process based on the lowest virtual runtime

Compute its dynamic time slice based on the IFS

Allow the process to run, when the time slice ends repeat the process

CFS is Implemented with Red-Black Trees

A red-black tree is a self-balancing binary search tree

Keyed by virtual runtime

$O(\lg N)$ insert, delete, update

$O(1)$ find minimum

The implementation uses a red-black tree with nanosecond granularity

Doesn't need to guess the interactivity of a process

CFS tends to favour I/O bound processes by default

Small CPU bursts translate to a low virtual runtime

It will get a larger time slice, in order to catch up to the ideal

Scheduling Gets Even More Complex

There are more solutions, and more issues:

- Introducing priority also introduces priority inversion
- Some processes need good interactivity, others not so much
- Multiprocessors may require per-CPU queues
- Real-time requires predictability
- Completely Fair Scheduler (CFS) tries to model the ideal fairness