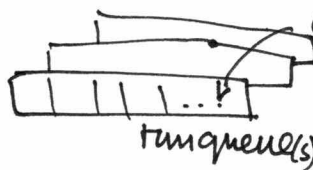


Scheduler (cont.)

Schedule ()

- has a run queue
- selects a task to run
- context switch to that new task

CORE scheduler functionality



enqueue to run queue
(ex. sth. becomes runnable because woken up, new process, moved from another run queue)

ADDITIONAL functionality (outside of schedule())

Example: /kernel/fork.c

(3.10) kernel P = copy_process(...)
wake-up-new-task(p) → set_task_cpu(p, ...)
activate_task(p) → ... enqueue-task(p, p, ...) ...

calls enqueue-task on the p → sched-class

class-specific versions of enqueue-task:

- idle-task has no implementation
- fair has enqueue-task-fair

calls enqueue_entity on the cfs runqueue (cfs_rq) of the sched-entity (se) related to the task we are enqueueing (p → se)

struct cfs_rq has cf on sched/sched.h has "CFR related fields in a runqueue"

! I.E. there is a runqueue that is specific to the scheduling class that this enqueue-task-fair is going to put this task on (i.e. on a class-specific runq).

NOTE: In Linux, the generic runq is composed of class-specific runqs and the class-specific functions stick what they want to enqueue and dequeue to those specific runqs.

RQ:

- CFS-RQ
- RT-RQ

In sched/sched.h:

```
struct cfs_rq;
struct rt_rq;
```

In sched/sched.h: "the main, per-CPU runqueue data structure"

struct rq {

```
struct cfs_rq;
struct rt_rq;
```

```
struct task_struct *idle;
```

};

! **NOTE:** to add your own scheduler, add its class-specific runq here

NOTE: These (cfs, rt, and idle) are all the schedulers implemented in 3.10 kernel!

► **Sched-class** struct in sched/sched.h is the generic scheduling class structure, it has pointers to the following functions:

- enqueue-task = put something on the rung
- dequeue-task = take sth. off of the rung
- pick-next-task = the key schedule func. to select what runs next
- task-tick = called (among others) from scheduler-tick, which gets called by the timer code, w/ HZ frequency w/ interrupts disabled. Checks to see if another task should run. More or less, the interrupt handler for the timer interrupt.
- gets the rung for the current CPU (where the timer interrupt went off),
- gets the current job running on that CPU ($r_q \rightarrow curr$)
- calls the scheduling class' task-tick for the sched-class of the currently running task
- ex: task-tick-idle: does nothing
- ex: task-tick-rt: ... if ($--p \rightarrow rt.time_slice$) return; ...

as long as the time-slice is

positive, just keep going

otherwise "requeue to the end of queue" and set the reschedule flag on the current task: ... set-tsk-need-resched(p); ...

NOTE: The timer interrupt goes off whenever the timer goes off, the timer goes off at some fixed interval (for most architectures). The OS will check, when returning from the timer interrupt, if it needs to reschedule. Most modern hardware works this way. Some hardware allows us to schedule the timer regularly, if we know we will be running the same task.

• Expensive to call the timer interrupt? Say we call each millisecond on a 1GHz CPU and it costs us 1000 instructions (out of the 1 million we execute each millisecond). Not really expensive.

► Scheduling Algorithms:

① **FIFO**: enqueue rung → dequeue → CPU and run till its done. MINUS: while(1) will hose up the CPU. no scheduling decision on timer interrupt

BUT: if current job blocks, it still goes on a waitq and the next one from the rung is scheduled. What happens when the former job wakes up?

- put it at the end - simple, keeps the FIFO principle, typical implementation
- preempt the task that is running
- insert at the front at the queue

② **Round Robin**

goes around, giving each task an equal shot at running (the easiest fair algo)

enqueue → rung → dequeue → CPU. But only run for some time quantum. Say 10ms. Then, if not done, return to back of rung.

• Example: Linux RT scheduler task_tick_rt:

RT scheduler supports two different types of scheduling:

• **SCHED-RR** = Round Robin

• **SCHED-FIFO** = FIFO

task_tick_rt has ... if (p->policy != SCHED-RR) return;
I.e. it runs FIFO tasks till completion!
Otherwise (SCHED-RR), decrease timeslice and, if zero, return to back of queue and reschedule (via reschedule flag).

• The Linux RT scheduler is a little more complex than the above idea of having a queue with either FIFO or Round Robin for all tasks in it.

• The rt class in Linux has priorities (numbers) associated with each task: You run the task w/ the highest priority.

• The caveat w/ priority scheduling is how to break ties... This is either FIFO or Round Robin as described above.

• Ways to implement a priority scheduler (NOTE: under each implementation, a low priority task will not run if higher one is ready)

① One queue w/ tasks w/ diff priorities - scheduler searches the queue to find the job w/ the highest priority and run it (search may be lengthy).

② Multiple queues, one for each priority level - scheduler finds first non-empty queue. NOTE: also, while(1) at high priority will hog up the CPU.

Linux uses this and w/in each queue the jobs are marked SCHED-RR or SCHED-FIFO. Priority scheduling is ~~not~~ not the default in Linux → you need root privileges to change task priorities.

NOTE: If you CTRL-C in a shell to stop an errand high-priority process, you may not be able to since the shell process needs to actually run and it may never do so while the higher priority process runs (the CTRL-C interrupt will still happen and get executed but the shell process won't have a chance to send the associate signal to the errand high-priority process). You can certainly reboot the computer to stop the high-priority process.

NOTE: fork() copies the scheduling parameters (scheduler, priority, etc.) from parent, so these are inherited.

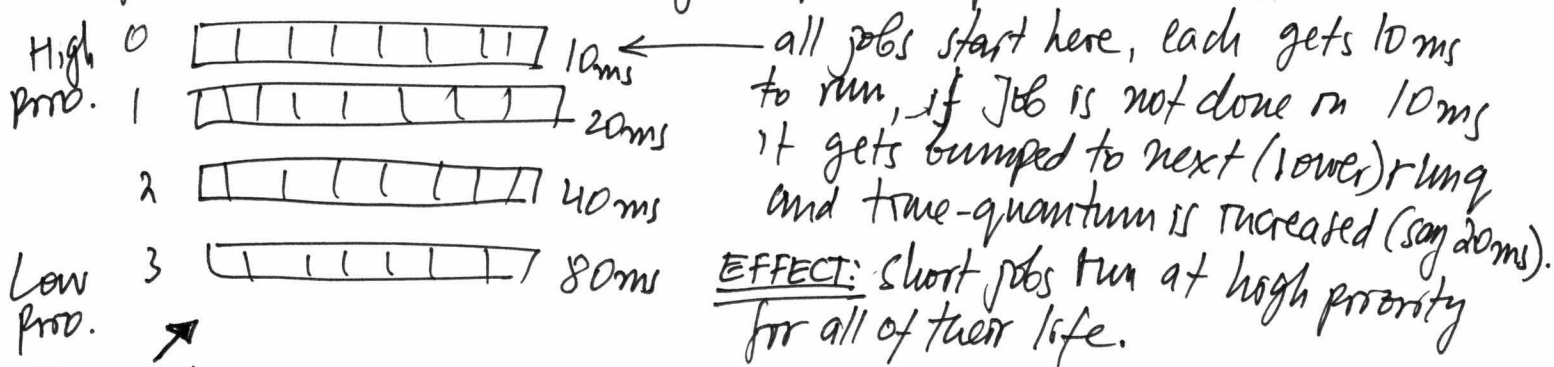
By default, Linux uses the (completely) Fair scheduler (CFS) → the first process that gets created is CFS scheduled and that gets inherited.

To run in a custom (different) scheduler, you need to find that initial task and, before it creates the other processes, change it or them to use that custom scheduler.

- When you implement a scheduler and need to test it w/ some jobs, make sure that the jobs are actually using this scheduler class!
 - Even when you inherit the scheduling parameters of your parent, that does not mean you have to do everything your parent does: ← You adjust some!!!
 - ▶ Say parent has 10ms quantum to run, do we use same for child?
 - Many forks from same parent may usurp the system if each child has 10ms too.
 - You can change so fork() gives $\frac{1}{2}$ time quantum to the child (ex. 5ms left to parent and 5ms to child) - overall quantum stays constant
 - An in-between approach where we reduce the time quantum w/ each fork()
- ③ Shortest Job First: If you run the shortest job first, the wait time is minimized and the job throughput is maximized. (OVERALL!)

• PROBLEM: How do you know how long the job will run?

• An approximation is used w/ many runs in a prioritized list:



• This is also a common way to assign priorities AKA Multi-Level (Queue) Scheduling and is commonly used (ex. interactive commands are short & high priority)

But, this has fallen out of favor because there are some problems w/ it:

- Long-running tasks may still need to respond quickly (ex. web browser, audio player, etc.) but will get buried to a low priority queue.
- You can improve the above a bit by splitting those long-running processes into continuous runs, for example, between sleeping - so, each time your process sleeps, put it back in the High prio que.
- But such base on sleep/run behavior is also not optimal.