

● Scheduler (cont.)

► Scheduling Algorithms (cont.)

- Most of these, we discussed in context of one CPU, w/ multiple, we can still use one q. per CPU, so some discussion.
- ① ● FIFO
 - ③ ● Shortest Job First / Shortest Burst Job First (min waiting time) (shortest burst time, related to)
 - ② ● Priority (w/2)
 - ② ● Multi-Level Queues / Multi-Level Feedback Queues (longer jobs get to lower queues) (when jobs sleep, they go back to higher queues) (w/3)
 - ② ● Round-Robin - good for fairness

► FAIRNESS in more detail:

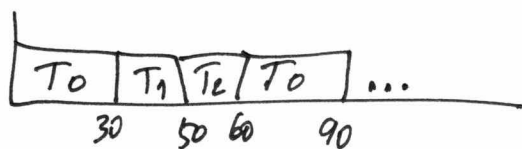
④ ● Weighted Round Robin

Like Round-Robin, but associated w/ each task is a weight that adjusts its allowed run time.

$w_i = \text{weight}$

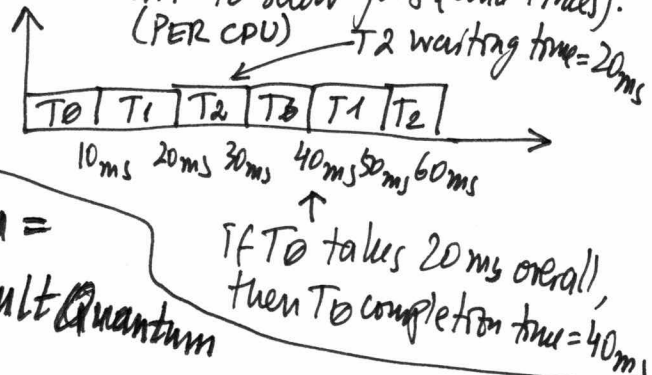
where the weight is a ratio scale (w/ meaningful zero and ratios).

EX: $T_0 = 3, T_1 = 2, T_2 = 1$



SOME TERMS:

- waiting time = how long you're runnable but not run (that time)
- arrival time = when you arrive at runq
- completion time = when task completes or time until task completes (Atat turnaround time)
- GANTT chart to show jobs (and times): (PER CPU)



Time quantum = $w_i \times \text{default Quantum}$

DISADVANTAGE: Suppose for jobs above the weights were 3000, 2000, and 1000. According to this algo, T_0 will run for 3000ms, T_1 for 2000 etc. That would not feel right (T_0 will run for 30 seconds before T_1 runs!).

So, we would like to keep the proportionality but w/ lower granularity

⑤ ● Lottery Scheduling ⇒ Line up all jobs & give them tickets, the number of which is proportional to the weight of the job.

EX: In above, distribute 3000 tickets: 3000 for T_0 , 2000 for T_1 , 1000 for T_2

⇒ Pick a winning ticket & run associated job for one time quantum (say 10ms)

⇒ Run Lottery again, etc. NOTE: ticket gets put back, i.e. chances are the same each round

- Over time, the ratio at which T_0, T_1, T_2 run is 3:2:1 and jobs are run in 10 ms increments.
- Alternatively, you could try to normalize the weights, so $w \times \text{default quantum}$ is something feasibly small BUT finding a common factor such that all jobs' weights are divisible by it (1000 in prev. example) may not be trivial.
- The lottery approach gets around the above issue by using statistics.
- Problem with the lottery is time is:

 $O(n)$

Which is worse than FIFO or Round Robin which are $O(1)$.

- Another problem is in the short term we may not get the 3:2:1 ratio.

A solution:

⑥ Weighted Fair Queuing (WFQ):

Each job has a weight, but we also introduce a virtual time associated w/ it.

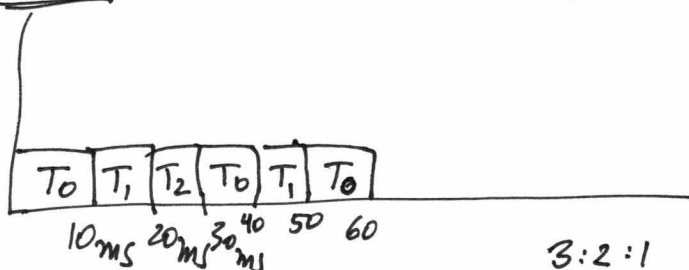
$$w_i \quad VT_i = \frac{1}{w_i} * E(t) \quad \leftarrow \text{the amount of time job } i \text{ actually runs.}$$

Then run the task with the smallest virtual time for a time quantum. (say 10ms)

Example:

	T_0	T_1	T_2
w	3000	2000	1000

Gantt:



$\frac{1}{600} = GVT_i$	VT_{i1}	0	0	0 \leftarrow pick any (say T_0)
$\frac{2}{600} = GVT_1$	VT_2	$\frac{1}{300}$	0	0 \leftarrow pick T_1 or T_2 (say T_1)
$\frac{3}{600} = GVT_2$	VT_3	$\frac{1}{300}$	$\frac{1}{200}$	0 \leftarrow pick T_2
$\frac{4}{600} = GVT_3$	VT_4	$\frac{1}{300}$	$\frac{1}{200}$	$\frac{1}{100} \leftarrow$ pick T_0
\vdots	VT_5	$\frac{2}{300}$	$\frac{1}{200}$	$\frac{1}{100} \leftarrow$ pick T_1
	VT_6	$\frac{2}{300}$	$\frac{1}{100}$	$\frac{1}{100} \leftarrow$ pick T_0
	VT_7	$\frac{1}{100}$	$\frac{1}{100}$	$\frac{1}{100} \leftarrow$ Things will repeat from here on

NOTE: That the ratios are perfect and in increments of the time quantum!

Scale to avoid using floating point on the ratios, i.e. use fixed point.

NOTE: Since we search for the smallest ^{operations} VT , performance is not $O(1)$ but $O(\log n)$.
 Overflow is not a problem: suppose we ~~intentionally~~ work with 64 bits and scale by 2^{20} (to avoid float point) - we have 2^{44} bits left and it will take some time to overflow (a lot of time!).

- PROBLEM** is when a new job shows up... If you set its VT to 0, it will run till it catches up w/ the rest of the jobs. You could set it to the lowest VT , or keep track of global $VT = \sum(w_i) \cdot \text{time} \dots$

- Keep track of Global VT as:

$$GVT = \frac{1}{\text{Sum}(W_i)} \cdot \text{Time}$$

Sum of all weights

- We could assign new jobs the GVT at the time they are added to the run queue.

← then you can compare each job's VT to the GVT and run the job that is ~~most~~ smaller than the GVT by the most.
 $VT_i \approx GVT$ means job i ran more or less equally w/ other jobs, considering weights.

- Linux's CFS uses a variant of the Weighted Fair Queuing (WFQ) algo - the virtual time is called "virtual runtime."

CLARIFICATION: Weighted Fair Queuing (WFQ) does not actually use the VT, but the virtual Finishing Time (VFT) of tasks:

$$VFT_i = VT_i + \frac{1}{W_i} * \text{Time Quantum}$$

- VFT is thus, the VT that you would have if you were to run.

Repeat example w/ VFT:

NOTE: this also breaks ties...

NOTE: Proportions here are 3:2:1

again →

↑
most

	T_0	T_1	T_2	
W	3000	2000	1000	
VFT ₁	1/300	1/200	1/100	← pick T_0
VFT ₂	2/300	1/200	1/100	← pick T_1
VFT ₃	2/300	2/200	1/100	← pick T_0
VFT ₄	3/300	2/200	1/100	← pick T_1
VFT ₅	3/300	3/200	1/100	← pick T_0 (any would be ok)
VFT ₆	4/300	3/200	1/100	← pick T_0 (T_0/T_2 tie)
VFT₇	4/300	3/200	1/100	← pick T_1
VFT ₇	4/300	3/200	2/100	

Every job has a deadline, run the job w/ the earliest deadline first.

⑦ Deadline Scheduling:

Property: If you run the jobs in deadline order and a schedule exists that can achieve all deadlines, the deadline scheduling will use that schedule! I.e. it is optimal. caveat: If no such schedule exists, this algo is not very efficient.

- So far, we mostly considered uniprocessor scheduling. Multiprocessor is similar, since we can repeat the algo on each CPU w/ its own queue. BUT:
- How do we assign a job to a CPU to begin with?

► Load Balancing: Multiple CPUs, each w/ own queue, where do we put a new job

- Least # of tasks CPU - count # tasks on each queue & put job

Does not account for job weights, or the running time of jobs, or

- Least sum of weights CPU - ^{resolved here} use CPU w/ least sum of weights ^{priority} ^{resolved here}
- Least priority CPU - use CPU w/ lowest max priority on queue

NOTE: Load Balancing may need to happen not only when a new job shows up, but also when jobs exit and/or block, also if a CPU is running the idle task and other CPUs run jobs (idle balancing), also periodically.

- Idle Balancing & Load Balancing involve moving a job from one CPU to another. From a synchronization viewpoint this involves locking multiple run queues (and the associated danger of deadlocks).

- Periodic Load Balancing may be done with a kernel daemon, for example

❗ Linux Source Example in `/kernel/sched/fair.c` look for "Load Balance" (3.10 kernel)

- More examples from Linux Source:

⊗ the scheduling class is the primary thing by which you implement a scheduler - the underlying core scheduler calls functions from the scheduling class.

⊗ functions that get called are:

- pick-next-task
- enqueue-task / dequeue-task

• many other hook functions to the scheduler class that the core scheduler promises to call at a particular point in time so that you can make a decision at that point in time how the scheduler will work.

EXAMPLE: `core.c` → `schedule()` calls `pick-next-task()` BUT BEFORE THAT it calls `put_prev-task(rq, prev)` which notifies the scheduler class just before the pick.

COMSW4118-15-5

10/24/2017

- For example, this way just before you pick a new task you can update the virtual time.
- pre-schedule() is another hook call back that happens before pick-prev-task and idle-balance. You can implement a pre-schedule function there.
- post-schedule() is called after you are done picking the next task.