

# Lecture 27 — Scheduling, Idling, Priorities, and Queues

Jeff Zarnett

`jzarnett@uwaterloo.ca`

Department of Electrical and Computer Engineering  
University of Waterloo

November 14, 2015

# Scheduling Algorithms, Continued

Carrying on from last time, we will examine some more scheduling algorithms.

We will introduce a new measure: **normalized turnaround time**.

This is the ratio of the turnaround time (the time waiting plus the amount of time taken to execute) to the service time (the time it takes to execute).

We can tolerate longer processes waiting a comparatively longer period of time.

Goal: minimize not only the normalized turnaround time for each process, but to minimize the average over all processes.

Calculate  $R = \frac{w+s}{s}$  where  $w$  is the waiting time and  $s$  is the service time.

The service time is, as always, a guess.

When it is time to select the next process to run, choose the process with the highest  $R$  value.

A new process will have a value of 1.0 at the beginning, because it has spent no time waiting (yet). Thus it is not that likely to get selected.

Jobs with a small  $s$ , i.e., short jobs, are likely to get scheduled quickly.

The HRRN approach introduces something important: age of the process.

A process that has spent a long time waiting rises in priority until it gets a turn.

So processes will not starve, because even a process that is expected to have a very long  $s$  will eventually have a high enough  $R$  due to the growth of  $w$ .

We still need to estimate  $s$ , which may or may not be simple guessing.

# Multilevel Queue (Feedback)

For the most part, until now, we have treated processes more or less equally (except when we have taken the highest priority process).

While it might seem very fair, it may not be ideal for a situation where processes behave differently.

A desktop or laptop has many processes running (foreground/background).

We can apply different scheduling algorithms to different types of process.

# Multilevel Queue (Feedback)

The multilevel queue takes the ready queue and breaks it up into several.

A process can be in one (but only one) of the queues.

It is assigned to the queue based on some attribute of the process (priority, memory needs, foreground/background, et cetera).

The foreground queue, for example, could be scheduled by Round Robin, and the background by First-Come, First-Served.

When there are multiple queues, we also need a way of choosing which of the queues to take from next.

We might say some queues have absolute priority over others, or we might have time slicing amongst the queues.

This could be balanced evenly (rotate through each) or give more time slices to some queues at the expense of others.



An example of this was the Compatible TimeSharing System on the IBM 7094.

Give CPU-Bound processes longer blocks of time to execute so they would not have to spend so much time swapping.

In the highest priority class, a process got 1 time slice; in the next one down, a process got 2 time slices; the third class meant 4 time slices, and so on.

If a process ran up against the limit of a time slice it was moved down a class.

So it got a lower priority, but when it did get selected to run, it was able to run with a lower chance of being interrupted.

Like a few schemes, we have seen so far, this is a ratchet.

A process can move down in the priority list, but there does not appear to be a way for it to move up.

A process that needed a lot of CPU early on was going to be punished “forever”.

If the user pressed the Enter key, it might be a sign the process was likely to become interactive (and therefore should move up in priority).

Some genius user (there's always one), figured out that by pressing the enter button repeatedly, his long running processes would finish faster.

This was a bit unfair; his processes got priority over the others.

Things really broke down when this individual decided to be nice:  
... He told all his friends.

Suddenly everyone was doing it and the benefit of the system was lost.

This scheduling algorithm may also be referred to as **feedback**.

We do not have any information in advance about how long processes will be.

Assign priority based on the amount of CPU time assigned so far.

A process that has used a lot of CPU so far gets lower priority.

Promise the users something and then fulfill that promise.

If there are  $n$  users, each gets an equal share ( $1/n$ ) of the CPU time.

Or with  $m$  processes, each process gets  $1/m$  of the CPU time.

The system must keep track of how much CPU time each process has received.

It then considers the how this value compares to the ideal.

If a process has a value of 0.5, it had only half the CPU it “should” have received.

If it has a value of 2.0, it has had double.

So the scheduling algorithm is then to run the process with the lowest score, trying to keep all values as close to 1.0 as we can

The lottery is a system to give predictable results with a simple implementation.

The premise is that every process gets some number of “lottery tickets” for each resource (e.g., CPU).

When a decision has to be made, a lottery ticket is selected at random.

The process that holds that ticket gets that resource.

This system provides clarity; if a process has priority  $p$ , what does that mean?

If a process has a fraction  $f$  of the total tickets, then we can expect that process to get about  $f\%$  of the resource.

When a process is created or terminates this may increase or decrease the number of tickets, or result in their redistribution.



More important processes get more tickets & have a higher chance of winning.

If there are 100 tickets outstanding, if a process has 25 of them, it has a 25% chance of winning any given draw.

To increase a process's chance of winning, give it more tickets.  
To decrease it, give it fewer.

Unlike in the real lottery, though, there is always a winner.

Co-operating processes may be permitted to exchange tickets.

A client that sends a request to a server might then give its tickets to the server.

This increases the chance the server gets the resources to run next and respond.

This is a lot less overhead than guaranteed scheduling.

We do not keep track of how much of the resource a process has received.

Assuming that the lottery system is sufficiently random, over time the resource allocation will tend towards the proportions of the tickets each process holds.

If process *A* has 20% of the tickets, *B* has 30%, and *C* has 50%, then the CPU will be given to the processes in approximately a 20:30:50 ratio, as expected.

Sometimes our scheduling algorithm cannot produce a new process to run next because there is, quite frankly, nothing to do.

The actual implementation of the idle thread may vary across different systems.

In some cases it is just repeatedly invoking the scheduler.  
In others it does a bit of useless addition.  
Or it might just be a whole bunch of NOP instructions.

The CPU can be told to halt/switch to a low power state.

Whatever it actually “does”, the idle thread never has any dependencies on anything else and is always ready to run.

Since the idle task does not necessarily do much, why have it?

It prevents having special cases in the scheduler, first of all.

It also provides accounting information about how much of the time the CPU is not doing anything.

In fact, a lot of the time on the desktop or laptop, task manager will tell you that “System Idle Process” is taking up a large percentage of the CPU.

You will recognize that this just means the CPU is not doing anything;  
It does not mean that some system process is using up all your CPU's time.

Saving power by shutting down (parts of) the processor seems like a nice savings of energy (and potentially increases battery life).

But: time when the CPU is doing nothing might potentially be put to use.

There are usually some accounting and housekeeping tasks that the CPU can be doing when it has nothing else.

For example, the OS could collect statistical data, or defragment the hard drive.

Sometimes we get into a situation called a **priority inversion**.

A high priority process is waiting for a low priority process.

$P_1$  is high priority and is blocked on a semaphore, while  $P_2$  is in a critical section.

As  $P_2$  is low priority, it might be a long time before  $P_2$  is selected again to run and can finish and exit the critical section.

$P_1$  cannot run, because it is blocked, and it could be blocked for a long time.

In the meantime, other processes with lower priority than  $P_1$  (but higher than  $P_2$ ) carry on execution.

Having  $P_1$  waiting for the lower priority processes is rather undesirable.



The solution is **priority inheritance**.

The right thing to do is to bump up the priority of  $P_2$ , temporarily, to be equal to that of  $P_1$ , so that  $P_1$  can be unblocked as quickly as possible.

To generalize, a lower priority process should inherit the higher priority if a higher priority process is waiting for a resource the lower priority process holds.

So  $P_2$  will get selected, will execute and exit the critical section.

Its priority then falls, meaning  $P_1$  will be selected and may continue.

A famous case of priority inversion took place on the Mars Pathfinder rover.

The solution was to enable priority inheritance.

Let's see analysis of the effectiveness of some of the scheduling algorithms.

We will assume that processes arrive according to a Poisson distribution (randomly), and have exponential service times.

Note that any scheduling relationship that chooses the next item without caring how long it will run (the expected service time) obeys this relationship:

$$\frac{T_r}{T_s} = \frac{1}{1 - \rho}$$

where:  $T_r$  is the turnaround time (waiting plus execution);  
 $T_s$  is the service time (average time in running state); and  
 $\rho$  is the processor utilization.

This includes priority-scheduling, where the priority is assigned by some means other than based on their (predicted or known) execution times.

If we do make a distinction based on the expected service times, then we get some meaningful results.

In our scenario, we will have two different priority classes (“high” and “low”, how exciting) with different service times.

Preemption, in this context, means that a low priority process will be interrupted as soon as a higher priority process is ready.

Assumptions: 1. Poisson arrival rate.

2. Priority 1 items are serviced before priority 2 items.
3. First-come-first-served dispatching for items of equal priority.
4. No item is interrupted while being served.
5. No items leave the queue (lost calls delayed).

## (a) General formulas

$$\lambda = \lambda_1 + \lambda_2$$

$$\rho_1 = \lambda_1 T_{s1}; \quad \rho_2 = \lambda_2 T_{s2}$$

$$\rho = \rho_1 + \rho_2$$

$$T_s = \frac{\lambda_1}{\lambda} T_{s1} + \frac{\lambda_2}{\lambda} T_{s2}$$

$$T_r = \frac{\lambda_1}{\lambda} T_{r1} + \frac{\lambda_2}{\lambda} T_{r2}$$

## (b) No interrupts; exponential service times

$$T_{r1} = T_{s1} + \frac{\rho_1 T_{s1} + \rho_2 T_{s2}}{1 - \rho_1}$$

$$T_{r2} = T_{s2} + \frac{T_{r1} - T_{s1}}{1 - \rho}$$

## (c) Preemptive-resume queueing discipline; exponential service times

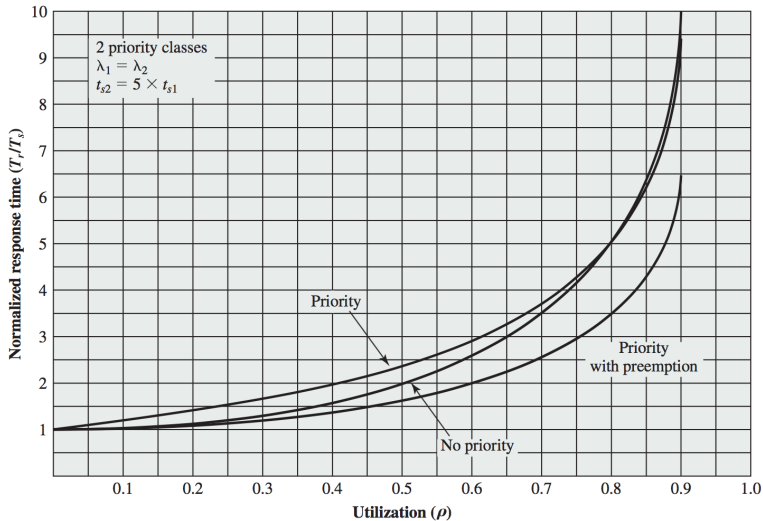
$$T_{r1} = T_{s1} + \frac{\rho_1 T_{s1}}{1 - \rho_1}$$

$$T_{r2} = T_{s2} + \frac{1}{1 - \rho_1} \left( \rho_1 T_{s2} + \frac{\rho T_s}{1 - \rho} \right)$$

Assume we have an equal number of arrivals between high and low priorities.

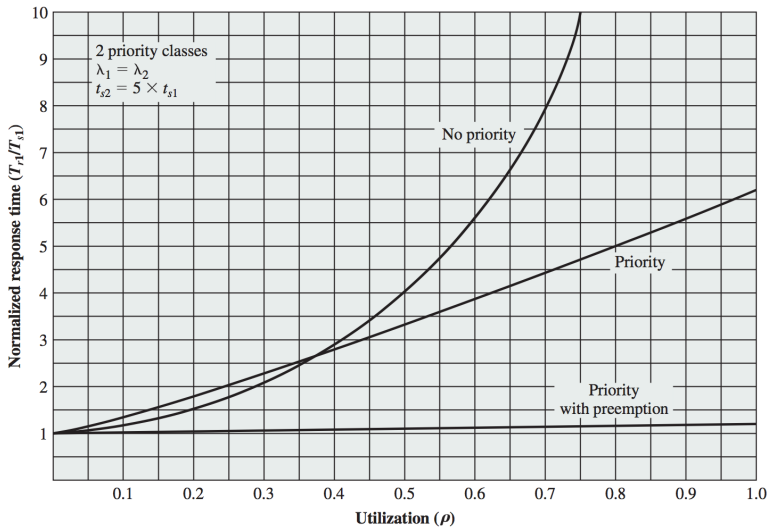
Low priority tasks take about five times as long as high priority tasks.

# Overall Normalized Response Time

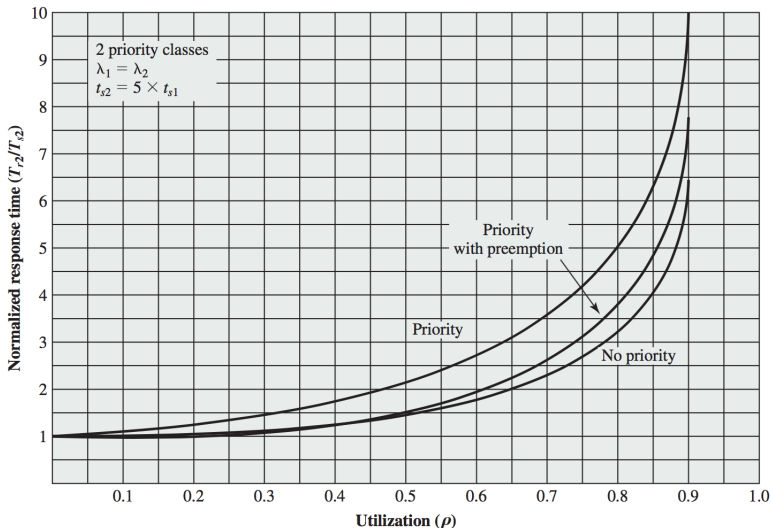




# Shorter Process Normalized Response Time



# Longer Process Normalized Response Time

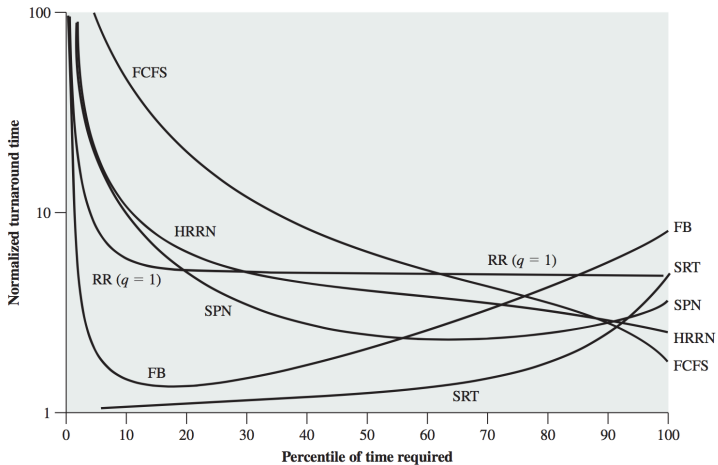


A simulation with 50 000 processes with an arrival rate of  $\lambda = 0.8$  and  $T_s = 1$ .

Processor utilization  $\rho$  is also 0.8.

Each process is grouped into service time percentiles of 500 processes.

# Simulated Turnaround Time



# Simulated Wait Time

