



## Topic Notes: CPU Scheduling

The CPU is the first of several scarce resources that we will consider. A primary function of an operating system is to determine which processes (and, in turn, users) get to utilize the CPU(s).

CPU scheduling can be done at three different levels:

1. *Long-term Scheduling* – also known as batch scheduling. Decide which jobs/processes are allowed into the system.
2. *Short-term Scheduling* – or interactive scheduling. Decide from a collection of ready processes which gets the CPU next.
3. *Medium-term Scheduling* – or memory scheduling. Decide if/when a process should be “swapped out” or back in based on memory available.

We will discuss primarily short-term scheduling here, though we will discuss all of the algorithms that may be split between batch and interactive scheduling.

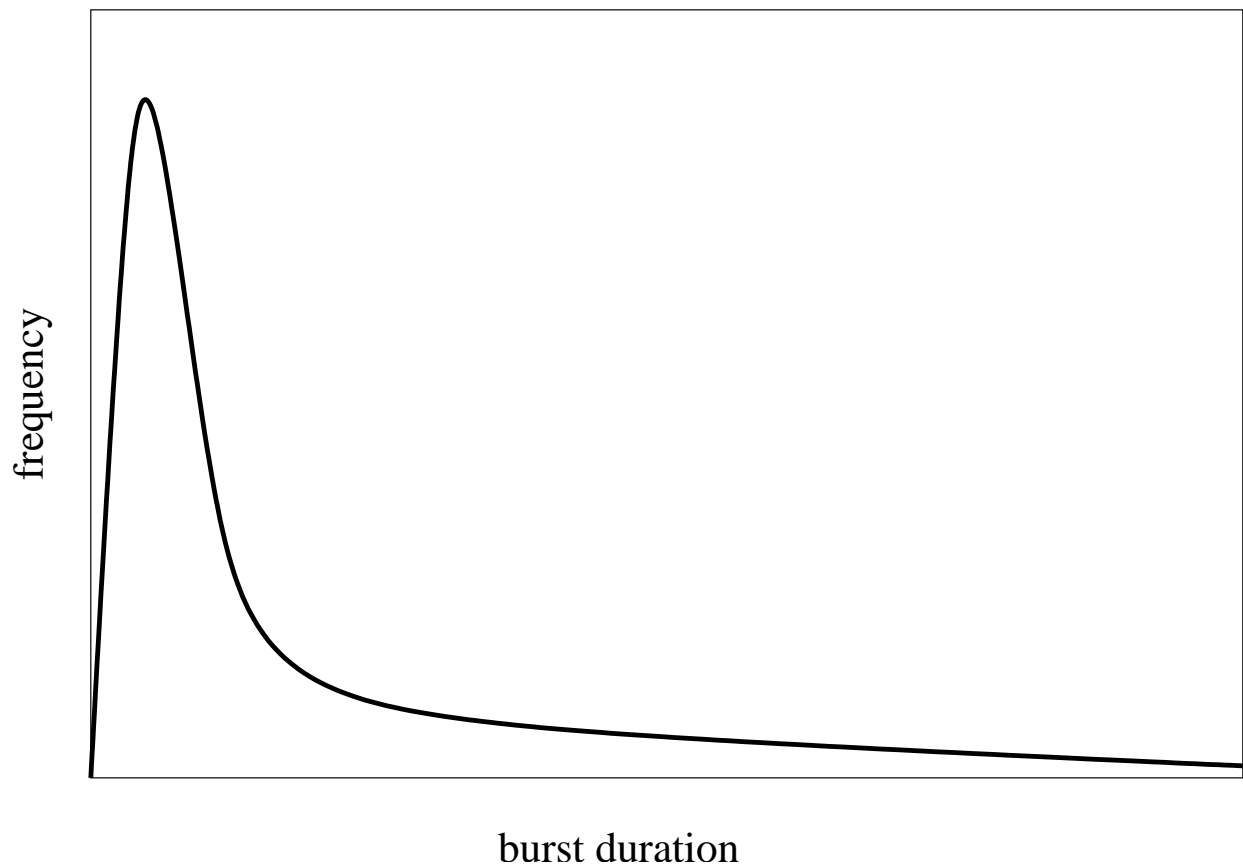
For the purposes of this discussion of general scheduling algorithms, I will use the term “process” but it could also apply to threads in many cases.

A typical process alternates between the need for CPU and the need for I/O service throughout its lifespan. This is called the *CPU-I/O burst cycle*.

It is this fact that makes multiprogramming essential. When a process needs I/O, it’s good to have another process ready to move in and take advantage of the available CPU resource.

The amount of time that it can make use of the CPU is known as its *CPU burst time*.

A typical distribution of CPU burst times looks like this:



Processes may be categorized as:

- *CPU-bound* – process does not need much I/O service, almost always want the CPU
- *I/O-bound* – short CPU burst times, needs lots of I/O service
- *Interactive* – short CPU burst times, lots of time waiting for user input (keyboard, mouse)

The type of processes in the system will affect the performance of scheduling algorithms.

A short-term CPU scheduling decision is needed when a process:

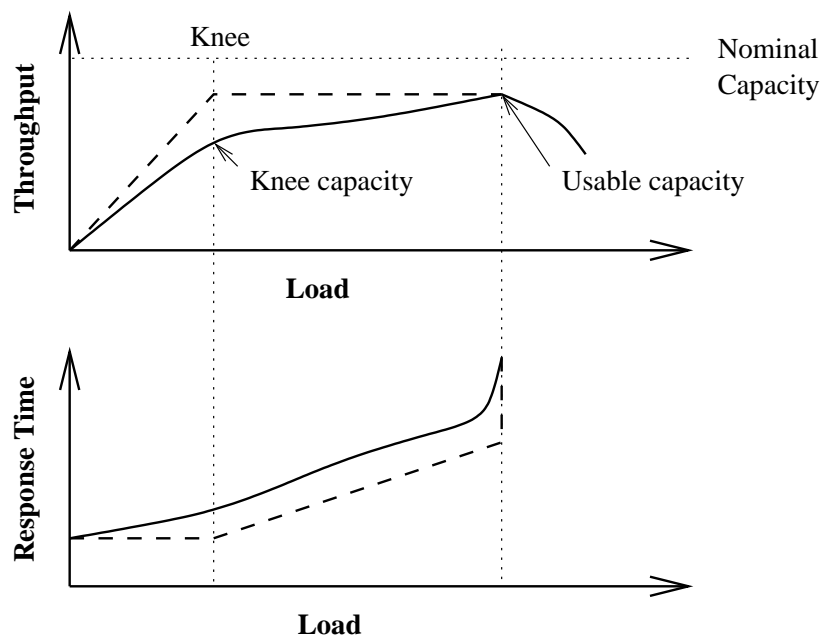
1. switches from a running to a waiting state (non-preemptive)
2. switches from a running to a ready state (preemptive)
3. switches from a waiting to a ready state (preemptive)
4. terminates (non-preemptive)

---

## Goals of a CPU scheduler

- maximize *CPU Utilization* – keep the CPU busy doing useful work (owner)
- maximize *Throughput* – rate of process completion (owner)
- minimize *Turnaround Time* – amount of time to execute a particular process (user)
- minimize *Waiting Time* – amount of time that a process is waiting in the ready queue (user)
- minimize *Response Time* – amount of time it takes from a process' arrival until its first turn on the CPU (user)

## Response Time's Relationship to Throughput



## CPU Scheduling Algorithms

### First-Come, First-Served (FCFS) Scheduling

As its name suggests, the first process to arrive gets to go first. It is a non-preemptive FIFO system.

Example:

Consider 4 processes  $P_1$ ,  $P_2$ ,  $P_3$ , and  $P_4$  that have burst times of 18, 3, 5, and 4, respectively.

If the processes arrive in the order  $P_1$ ,  $P_2$ ,  $P_3$ ,  $P_4$ , FCFS will service them in that order. We visualize this with a *Gantt Chart*:

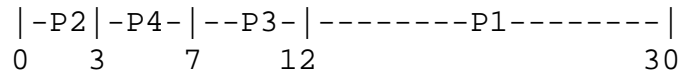
```

|-----P1-----| -P2- | -P3- | -P4- |
0                18  21   26   30

```

Waiting times:  $P_1 = 0; P_2 = 18; P_3 = 21; P_4 = 26$ , Avg:  $\frac{0+18+21+26}{4} = 16.25$ .

But what if the processes arrive  $P_2, P_4, P_3, P_1$ ?



Waiting times:  $P_1 = 12; P_2 = 0; P_3 = 7; P_4 = 3$ , Avg:  $\frac{12+0+7+3}{4} = 5.5$ .

FCFS characteristics:

- Penalizes short jobs
- Rewards long jobs
- Large variance in throughput
- Sensitive to arrival order
- Is starvation free – every job gets its turn
- Easy to implement

### Shortest-Job-First (SJF) Scheduling

a.k.a. Shortest Process Next (SPN)

Choose the process with the smallest next CPU burst.

Non-preemptive – process does not leave until its burst is complete

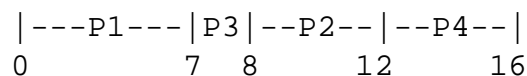
Preemptive – if a new process arrives with CPU burst length less than the remaining time of the currently executing process, we preempt. This is also known as *Shortest Remaining Time First (SRTF)*.

Example:

Consider these processes

Process	Arrival Time	Burst Time
$P_1$	0	7
$P_2$	2	4
$P_3$	4	1
$P_4$	5	4

Non-preemptive:



Average waiting time:  $\frac{0+6+3+7}{4} = 4$ .

Preemptive (SRTF):

| -P1- | -P2- | P3 | -P2- | --P4-- | --P1--- |  
 0     2     4 5     7     11     16

Average waiting time:  $\frac{9+1+0+2}{4} = 3$ .

SJF characteristics:

- It is *optimal* in minimizing waiting time
- Penalizes long jobs
- Rewards short jobs
- Somewhat sensitive to arrival order
- Gives optimal nonpreemptive throughput
- Permits starvation
- Difficult to predict burst/service times. We can try to estimate it based on previous bursts, or by averaging (text discusses this). However, this leads to an approximate SJF, which would not necessarily provide the optimal schedule.

## Priority Scheduling

Associate a priority with each process and always choose the one with the highest priority.

SJF/SRTF are examples of this, with the priority assigned as the burst times.

Biggest problem: starvation of low-priority jobs.

Can deal with this by *aging* – increasing the priority of processes that are not getting a chance.

## Round Robin (RR) Scheduling

Each process gets a small unit of time on the CPU (the *time quantum*), typically 10-100 ms. After this time, the job is preempted and added to the end of the ready queue.

For  $n$  processes and quantum  $q$ , each process gets  $\frac{1}{n}$  of the CPU time. No process waits more than  $(n - 1)q$  for its next turn on the CPU.

RR characteristics:

- Preemptive (at quantum  $q$ )

- Less sensitive to arrival order
- Quantum should not be too small relative to context switch time
- At overly large  $q$ , approximates FCFS
- Low overhead
- Starvation impossible

Example,  $q = 20$ :

Process	Burst Time
$P_1$	53
$P_2$	17
$P_3$	68
$P_4$	24

| -P1- | -P2- | -P3- | -P4- | -P1- | -P3- | -P4- | -P1- | -P3- | -P3- |  
 0    20    37    57    77    97    117   121   134   154   162

What if we changed  $q$  to 5? Lots more context switching - more overhead. You can see this in the queueing system lab.

Real values for a quantum tend to be in the 20-200 ms range, usually about 100 ms. This, amazingly, has been consistent for decades as hardware and operating systems have evolved.

Perhaps this is more a function of the delay that a human is willing to accept than something related to hardware speeds.

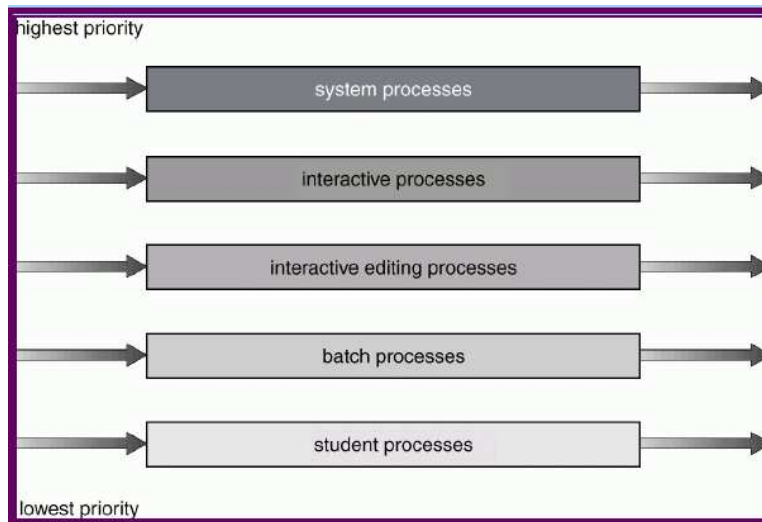
### *Multilevel queues*

We can partition the ready queue into a number of queues.

Perhaps with 2 queues, one can have an RR discipline for interactive jobs and the other can have a FCFS discipline for batch jobs.

Need to schedule among the queues as well, then.

Perhaps, a priority system – take jobs from high priority queues:



## Multilevel Feedback Queues

This is a method used in real systems.

Processes may move among queues; aging may be implemented this way.

Parameters:

- number of queues
- when to move a process to a different queue
- where new processes will enter
- scheduling among queues

A simple example:

Three queues,  $Q_0$  has  $q = 50ms$ ,  $Q_1$  has  $q = 100ms$ ,  $Q_2$  is FCFS.

Processes enter at  $Q_0$ , if not finished after 50 ms, move to  $Q_1$ , if not finished after an additional 100 ms, move to  $Q_2$ .

$Q_2$  is served only when  $Q_1$  and  $Q_0$  are empty.  $Q_1$  is served only when  $Q_0$  is empty.

Alternately, if there are a large number of queues, the time given to jobs in each state can be allocated among the queues to continue to make progress on low-priority jobs while still paying most attention to the high-priority ones.

How does this benefit interactive jobs? How does this benefit CPU-bound jobs? I/O-bound jobs?

“Niceness” – see Unix “nice” command.

Traditional SysV Unix systems use this kind of scheduling:

```

scheduler:
while (no process is picked to execute) {
  for (every process on run/ready queue) {
    pick highest priority process that is loaded in memory;
  }
  if (no process is eligible to execute) {
    idle until awakened by interrupt;
  }
}
remove chosen process from run/ready queue;
switch context to that process, resume its execution;

```

Ties in priority are broken by scheduling the process that has waited the longest (FIFO)

Each process has a priority field, function of its recent CPU usage.

Processes that have used the CPU a lot get lower priority.

Highest priorities are for low-level system processes. Then a class of kernel processes, then a class of user processes ( $n$  levels of priority within).

Processes that sleep in kernel mode are given high priority so they can continue immediately when they are awakened (I/O completes, for example).

A user-level process that gets preempted gets a “recent CPU usage” field that is set to the amount of time it just spent on the CPU.

Priority is calculated as:

$$\text{priority} = (\text{``recent CPU usage''} / 2) + (\text{base priority})$$

Once a second the system’s clock handler recomputes priorities of all ready processes by halving the recent CPU usage field.

This quickly “ages” processes which had a lot of usage but were then denied the CPU as other processes run.

This makes “interactive” processes get a higher priority, as they need little CPU. When they want the CPU, they will get it.

Add in the idea of “nice” to be able to modify process priorities:

$$\text{priority} = (\text{``recent CPU usage''} / \text{const}) + (\text{base priority}) + (\text{nice value})$$

## Fair-share scheduling

What if the goal is to divide the CPU fairly among a group of users, regardless of the number of processes they start?



With a typical RR or feedback queue system, a user can launch lots of jobs.

It can be done exactly – keep track of the proportion of time each process/user has had access to the CPU and what proportion each was supposed to have. The processes most below their fair share are selected to run.

This was done in Unix SysV by Henry in 1984.

*Lottery scheduling* is one way to do this.

Each process gets a number of lottery tickets proportional to its fair share of the CPU.

The scheduler holds a lottery at each scheduling decision point and the process with the winning number get a prize – a trip on the CPU for up to one quantum!

This can be used to implement priorities – higher priority processes just get more tickets and hence a better chance to win each time.

If tickets are given to users to dole out among their processes, it can produce fair schedules even if some user tries to put many more processes on the system.

A new high-priority job can have a good chance for a good response time if it is given lots of tickets.

See the Petrou paper for details.

---

## Comparing Algorithms

How can we compare the algorithms and approaches?

- Deterministic modeling – consider a set of processes and see what happens with various algorithms. Our Gantt charts were an example.
- Queuing models – take a mathematical approach. Given mathematical descriptions of the kind of processes and the underlying system, determine expected performance. This would likely be a major topic in a graduate level operating systems course.
- Implementation – try it out on a real system. May be impractical, but it's the best test, of course.
- Simulation – devise a model and simulate it. See Lab 2.

---

## Multiple Processors

How do we do CPU scheduling when we have more than one CPU?

For simplicity (and reflecting most realities) we assume symmetric multiprocessing: all CPUs can run any task.

Options:

- separate queues for each CPU
- one set of queues shared among all CPUs

Having separate queues for each CPU makes it easy to choose the next process to run but we need to figure out in which CPU's queue to place a job and if/when to move jobs among the CPUs.

Having a shared single queue or set of queues means that any idle CPU can run any ready job.

But think about what this means if multiple CPUs need to choose a new process at exactly the same time.

This is possible: they are truly concurrent – 2 independent processors.

Could they both choose the same job and either duplicate work (bad) or maybe corrupt kernel structures?

If the CPUs have to wait and be sure only one is choosing a job next, that could be slow and would not scale well.

There is a potential advantage to schedule a process on the same CPU on which it was last executed. If we stay on the same CPU, there's a chance for cache reuse. Otherwise, we'll have misses for sure as the process ramps up on the new CPU.

This is called *affinity*.

But we don't necessarily want a process to be pinned to a CPU forever – processes come and go and this will not give a long-term load balance.

---

## Examples

Linux 2.6 and FreeBSD 5.x have introduced relatively new schedulers. More on these later.

The book says some things about a few systems:

- Solaris:
  - 4 major classes for scheduling:
    - highest priority to real-time tasks
    - next highest to system/kernel service threads
    - interactive
    - time-sharing

They separate interactive and time-sharing to try to give GUIs and similar things a very fast response, even on a system with a good number of time-sharing type processes competing for the CPU.

Within a class, there are priorities (0-59).

Priorities have their own quanta specified: higher priority get a shorter quantum. (200-20)

Also note that we can learn something about the Solaris approach by checking out the man page for “ts\_dptbl” on bullpen.

Note how configurable this is.

This is a good example of the separation of mechanism and policy.

Check out the policy on bullpen with the dispadmin command.

Try running some processes using the dostuff program on rivera.

- Windows XP:

Strict priority scheme. Highest-priority task always runs next.

- Linux:

Old scheduler (up to Linux 2.4) was traditional multilevel feedback queue with priorities.

Problems with efficiency under heavy loads and with SMP scalability.

New scheduler (Linux 2.6) is much more efficient – constant time.

Read about it in the text and in the supplemental reading. It’s easy reading and worth 10-15 minutes of your time.

It is fully preemptible – when a high priority task arrives, a lower priority task that is executing will be preempted for quick response times. Processes executing kernel code can be preempted - this was not the case in the previous Linux kernel.

For SMPs, one “runqueue” for each processor. Claim: scalability to 64-way SMP.

Processes move among runqueues when a CPU is idle, or periodically when not.

Dynamic priorities are computed and I/O bound jobs are given longer time slices.

The McKusick and Neville-Neil book talks in detail about FreeBSD scheduling.

<http://portal.acm.org/citation.cfm?id=1035594.1035622&coll=portal&dl=ACM&idx=1035594&part=periodical&WantType=periodical&title=Queue&CFID=39231776&CFTOKEN=54087012>

FreeBSD up to 5.1 uses the 4.4BSD scheduler:

- prioritized run queues
- always choose the highest priority
- multiple processes at a priority are executed in round robin fashion.
- multilevel feedback queues
- immediate switch to newly-arrived higher priority job if the current job is in user mode only (interrupt is generated)

- see methods: `resetpriority()` (note that the code matches the formula in the article), `setrunnable()`, `wakeup()`, `roundrobin()`, and `schedcpu()`, many defined in `/sys/kern/kern_synch.c` on a FreeBSD 4.x system.

How to compute the priorities? See the FreeBSD scheduling article.

Based on two values associated with a thread:

1. `p_estcpu` – estimate of recent CPU utilization of the thread
2. `p_nice` – “nice” value between -20 (high priority) and 20 (low priority), by default is 0

This is similar to the Unix SysV approach discussed earlier.

It does take into account the system “load average” when deciding on the decay rate of the CPU usage field.

Like SysV, it recomputes priorities once per second.

Blocked tasks do not need their priorities recomputed until they return to the system, so their recent CPU usage is computed as a function of system load and sleep time when they are returned to the ready state.

Even so, consider a heavily-loaded system. Once a second, lots of processes need to have their priorities recomputed and may be moved among the queues.

FreeBSD 5.2 and beyond use the ULE scheduler.

Like the Linux O(1) scheduler, it addresses SMP and is not dependent on the number of threads.

Why the name? It’s implemented in `sched_ule.c`

We can check it out on any FreeBSD 5 or higher machine, such as the cluster head node, which is running 6.2.

It includes per-processor queues to allow for affinity scheduling.

Processes migrate to another CPU only when there is an idle processor to occupy.

ULE also addresses symmetric multithreading (SMT) – aka hyperthreading.

System sees multiple CPUs but it’s really one core CPU that is multiplexing in hardware.

From the point of view of scheduling, these are multiple CPUs, but they are a little different in that there should be little penalty for migrating among these “virtual” processors.

Each processor has three queues:

- idle queue – all “idle” (low priority) threads – run only when there’s no one else who wants the CPU
- current queue – set of jobs ready to run

- next queue – another set of jobs ready to run but only after the current queue empties

After all jobs from “current” are gone, the current and next queues are swapped.

Interactive jobs are inserted into current, for good response.

It decides which jobs are interactive based on the ratio of sleep time to run time.

The ACM Queue article is short and pretty easy reading as well.

---

## **Final Thoughts**

Modern schedulers are often concerned with SMP.

It is desirable to have things configurable.

It is desirable to be able to select a task quickly even when there are a lot of jobs in the system.

Choosing the time slice in a priority system: we seem to see some systems that give high-priority jobs a longer quantum, others give high-priority jobs a shorter quantum. Why?