

## 2.9. CPU SCHEDULING

### 2.9.1 Basic Concepts

Almost all programs have some alternating cycle of CPU number crunching and waiting for I/O of some kind. (Even a simple fetch from memory takes a long time relative to CPU speeds.)

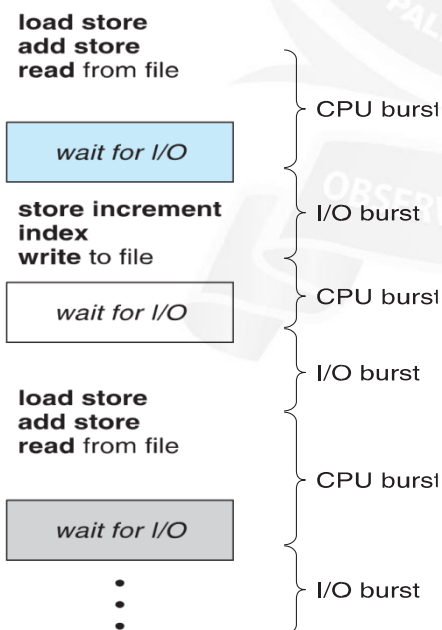
In a simple system running a single process, the time spent waiting for I/O is wasted, and those CPU cycles are lost forever.

A scheduling system allows one process to use the CPU while another is waiting for I/O, thereby making full use of otherwise lost CPU cycles.

The challenge is to make the overall system as "efficient" and "fair" as possible, subject to varying and often dynamic conditions, and where "efficient" and "fair" are somewhat subjective terms, often subject to shifting priority policies.

#### CPU-I/O Burst Cycle

- Almost all processes alternate between two states in a continuing **cycle**, as shown in Figure below :
- A CPU burst of performing calculations, and
- An I/O burst, waiting for data transfer in or out of the system.



## CPU Scheduler

- Whenever the CPU becomes idle, it is the job of the CPU Scheduler ( a.k.a. the short-term scheduler ) to select another process from the ready queue to run next.
- The storage structure for the ready queue and the algorithm used to select the next process are not necessarily a FIFO queue. There are several alternatives to choose from, as well as numerous adjustable parameters for each algorithm, which is the basic subject of this entire chapter.

## Preemptive Scheduling

CPU scheduling decisions take place under one of four conditions:

1. When a process switches from the running state to the waiting state, such as for an I/O request or invocation of the wait( ) system call.
  2. When a process switches from the running state to the ready state, for example in response to an interrupt.
  3. When a process switches from the waiting state to the ready state, say at completion of I/O or a return from wait( ).
  4. When a process terminates.
- For conditions 1 and 4 there is no choice - A new process must be selected.
  - For conditions 2 and 3 there is a choice - To either continue running the current process, or select a different one.
  - If scheduling takes place only under conditions 1 and 4, the system is said to be **non-preemptive**, or **cooperative**. Under these conditions, once a process starts running it keeps running, until it either voluntarily blocks or until it finishes. Otherwise the system is said to be **preemptive**.
  - Windows used non-preemptive scheduling up to Windows 3.x, and started using preemptive scheduling with Win95. Macs used non-preemptive prior to OSX, and pre-emptive since then. Note that pre-emptive scheduling is only possible on hardware that supports a timer interrupt.

Note that pre-emptive scheduling can cause problems when two processes share data, because one process may get interrupted in the middle of updating shared data structures. Chapter 5 examined this issue in greater detail.

- Preemption can also be a problem if the kernel is busy implementing a system call ( e.g. updating critical kernel data structures ) when the preemption occurs. Most modern

UNIX deal with this problem by making the process wait until the system call has either completed or blocked before allowing the preemption. Unfortunately this solution is problematic for real-time systems, as real-time response can no longer be guaranteed.

- Some critical sections of code protect themselves from concurrency problems by disabling interrupts before entering the critical section and re-enabling interrupts on exiting the section. Needless to say, this should only be done in rare situations, and only on very short pieces of code that will finish quickly, ( usually just a few machine instructions. )

### Dispatcher

The **dispatcher** is the module that gives control of the CPU to the process selected by the scheduler. This function involves:

- Switching context.
- Switching to user mode.
- Jumping to the proper location in the newly loaded program.
- The dispatcher needs to be as fast as possible, as it is run on every context switch. The time consumed by the dispatcher is known as **dispatch latency**.

### 2.9.2 Scheduling Criteria

There are several different criteria to consider when trying to select the "best" scheduling algorithm for a particular situation and environment, including:

#### CPU utilization

- Ideally the CPU would be busy 100% of the time, so as to waste 0 CPU cycles. On a real system CPU usage should range from 40% ( lightly loaded ) to 90% ( heavily loaded. )

#### Throughput

- Number of processes completed per unit time. May range from 10 / second to 1 / hour depending on the specific processes.

#### Turnaround time

- Time required for a particular process to complete, from submission time to completion. ( Wall clock time. )

#### Waiting time

- How much time processes spend in the ready queue waiting their turn to get on the CPU.

(**Load average** - The average number of processes sitting in the ready queue waiting their turn to get into the CPU. Reported in 1-minute, 5-minute, and 15-minute averages by "uptime" and "who". )

**Response time**

- The time taken in an interactive program from the issuance of a command to the **commence** of a response to that command.

In general one wants to optimize the average value of a criteria (Maximize CPU utilization and throughput, and minimize all the others. ) However sometimes one wants to do something different, such as to minimize the maximum response time.

Sometimes it is most desirable to minimize the **variance** of a criteria than the actual value.

I.e. users are more accepting of a consistent predictable system than an inconsistent one, even if it is a little bit slower.

**2.9.3 Scheduling Algorithms**

The following subsections will explain several common scheduling strategies, looking at only a single CPU burst each for a small number of processes. Obviously real systems have to deal with a lot more simultaneous processes executing their CPU-I/O burst cycles.

**1. First-Come First-Serve Scheduling, FCFS**

- FCFS is very simple - Just a FIFO queue, like customers waiting in line at the bank or the post office or at a copying machine.
- Unfortunately, however, FCFS can yield some very long average wait times, particularly if the first process to get there takes a long time. For example, consider the following three processes:

Process	Burst Time
P1	24
P2	3
P3	3

- In the Gantt chart below, process P1 arrives first. The average waiting time for the three processes is  $( 0 + 24 + 27 ) / 3 = 17.0$  ms.





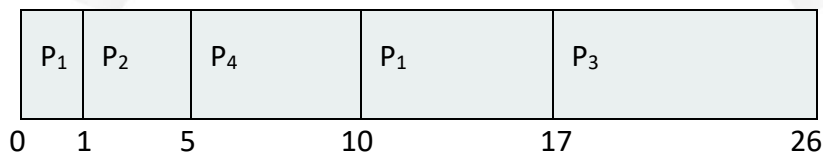
- In the case above the average wait time is  $(0 + 3 + 9 + 16) / 4 = 7.0$  ms, ( as opposed to 10.25 ms for FCFS for the same processes. )
- SJF can be proven to be the fastest scheduling algorithm, but it suffers from one important problem: How do you know how long the next CPU burst is going to be?
- For long-term batch jobs this can be done based upon the limits that users set for their jobs when they submit them, which encourages them to set low limits, but risks their having to re-submit the job if they set the limit too low. However that does not work for short-term CPU scheduling on an interactive system.

SJF can be either preemptive or non-preemptive. Preemption occurs when a new process arrives in the ready queue that has a predicted burst time shorter than the time remaining in the process whose burst is currently on the CPU. Preemptive SJF is sometimes referred to as

**3. Shortest remaining time first scheduling.**

- For example, the following Gantt chart is based upon the following data:

Process	Arrival Time	Burst Time
P1	0	8
P2	1	4
P3	2	9
p4	3	5



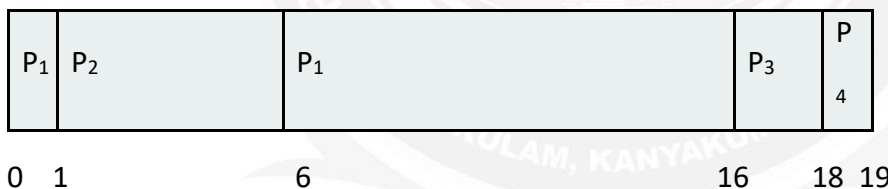
The average wait time in this case is  $(( 5 - 3 ) + ( 10 - 1 ) + ( 17 - 2 )) / 4 = 26 / 4 = 6.5$  ms.( As opposed to 7.75 ms for non-preemptive SJF or 8.75 for FCFS. )

Calculate Waiting time, average waiting time, turn around time, average turn around time

### 3. Priority Scheduling

- Priority scheduling is a more general case of SJF, in which each job is assigned a priority and the job with the highest priority gets scheduled first. ( SJF uses the inverse of the next expected burst time as its priority - The smaller the expected burst, the higher the priority. )
- Note that in practice, priorities are implemented using integers within a fixed range, but there is no agreed-upon convention as to whether "high" priorities use large numbers or small numbers. This book uses low number for high priorities, with 0 being the highest possible priority.
- For example, the following Gantt chart is based upon these process burst times and priorities, and yields an average waiting time of 8.2 ms:

Process	Burst Time	Priority
P1	10	3
P2	1	1
P3	2	4
P4	1	5
P5	5	2



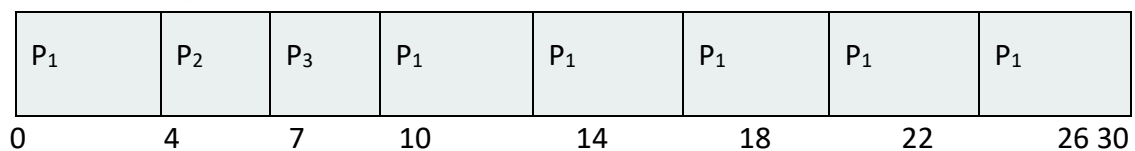
- Priorities can be assigned either internally or externally. Internal priorities are assigned by the OS using criteria such as average burst time, ratio of CPU to I/O activity, system resource use, and other factors available to the kernel. External priorities are assigned by users, based on the importance of the job, fees paid, politics, etc. Priority scheduling can be either preemptive or non-preemptive.
- Priority scheduling can suffer from a major problem known as **indefinite blocking**, or **starvation**, in which a low-priority task can wait forever because there are always some other jobs around that have higher priority.
- If this problem is allowed to occur, then processes will either run eventually when the system load lightens (at say 2:00 a.m. ), or will eventually get lost when the system is shut down or crashes. (There are rumors of jobs that have been stuck for years. )

- One common solution to this problem is **aging**, in which priorities of jobs increase the longer they wait. Under this scheme a low-priority job will eventually get its priority raised high enough that it gets run.
- Calculate Waiting time, average waiting time, turn around time, average turn around time

#### 4. Round Robin Scheduling

- Round robin scheduling is similar to FCFS scheduling, except that CPU bursts are assigned with limits called **time quantum**.
- When a process is given the CPU, a timer is set for whatever value has been set for a time quantum.
- If the process finishes its burst before the time quantum timer expires, then it is swapped out of the CPU just like the normal FCFS algorithm.
- If the timer goes off first, then the process is swapped out of the CPU and moved to the back end of the ready queue.
- The ready queue is maintained as a circular queue, so when all processes have had a turn, then the scheduler gives the first process another turn, and so on.
- RR scheduling can give the effect of all processors sharing the CPU equally, although the average wait time can be longer than with other scheduling algorithms. In the following example the average wait time is 5.66 ms.

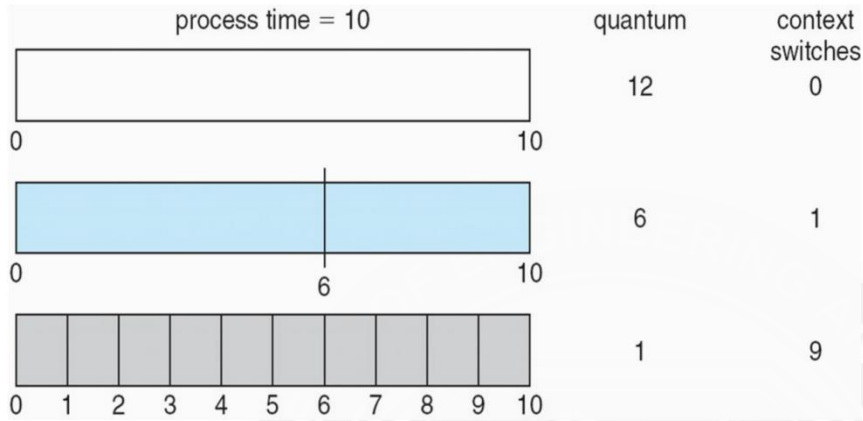
Process	Burst Time
P1	24
P2	3
P3	3



- The performance of RR is sensitive to the time quantum selected. If the quantum is large enough, then RR reduces to the FCFS algorithm; If it is very small, then each process gets 1/nth of the processor time and share the CPU equally.

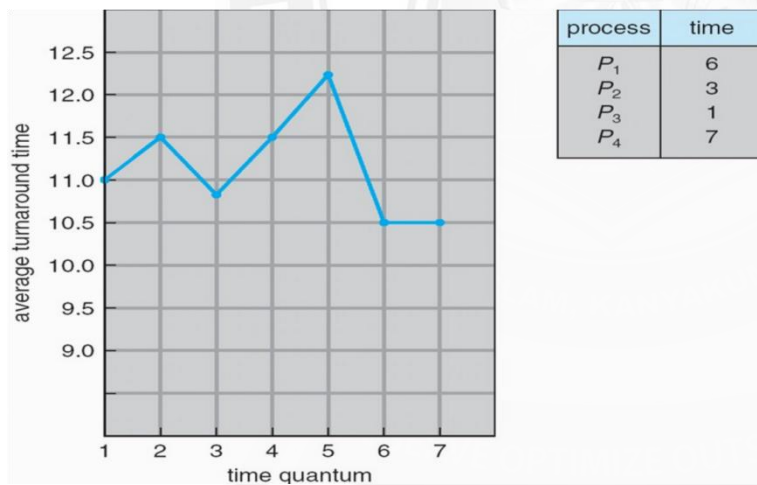


- **BUT**, a real system invokes overhead for every context switch, and the smaller the time quantum the more context switches there are. Most modern systems use time quantum between 10 and 100 milliseconds, and context switch times on the order of 10 microseconds, so the overhead is small relative to the time quantum.



**The way in which a smaller time quantum increases context switches.**

- Turn around time also varies with quantum time, in a non-apparent manner. Consider, for example the processes shown in Figure 6.5:



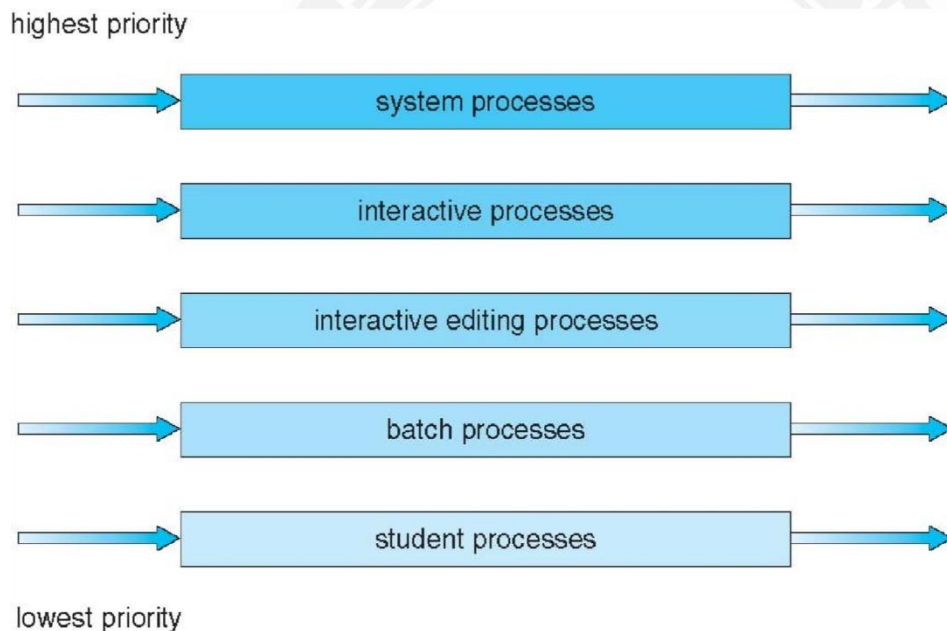
**The way in which turnaround time varies with the time quantum.**

In general, turnaround time is minimized if most processes finish their next cpu burst within one time quantum. For example, with three processes of 10 ms bursts each, the average turnaround time for 1 ms quantum is 29, and for 10 ms quantum it reduces to 20. However, if it is made too large, then RR just degenerates to FCFS. A rule of thumb is that 80% of CPU bursts should be smaller than the time quantum.

Calculate Waiting time, average waiting time, turn around time, average turn around time

## 5. Multilevel Queue Scheduling

- When processes can be readily categorized, then multiple separate queues can be established, each implementing whatever scheduling algorithm is most appropriate for that type of job, and/or with different parametric adjustments.
- Scheduling must also be done between queues, that is scheduling one queue to get time relative to other queues. Two common options are strict priority (no job in a lower priority queue runs until all higher priority queues are empty) and round-robin ( each queue gets a time slice in turn, possibly of different sizes. )
- Note that under this algorithm jobs cannot switch from queue to queue – Once they are assigned a queue, that is their queue until they finish.



## 6. Multilevel Feedback-Queue Scheduling

Multilevel feedback queue scheduling is similar to the ordinary multilevel queue scheduling described above, except jobs may be moved from one queue to another for a variety of reasons:

If the characteristics of a job change between CPU-intensive and I/O intensive, then it may be appropriate to switch a job from one queue to another.

Aging can also be incorporated, so that a job that has waited for a long time can get bumped up into a higher priority queue for a while.

Multilevel feedback queue scheduling is the most flexible, because it can be tuned for any situation. But it is also the most complex to implement because of all the adjustable parameters. Some of the parameters which define one of these systems include:

1. The number of queues.
2. The scheduling algorithm for each queue.
3. The methods used to upgrade or demote processes from one queue to another. (Which may be different. )
4. The method used to determine which queue a process enters initially.

