



# FIRST-COME FIRST-SERVED (FCFS)

The process that requests the CPU first (ie. Comes to the Ready Queue first), is allocated the CPU first.

## FCFS is implemented as follows:

When a new process arrives to the Ready Queue, it is appended to the tail.

When the CPU is free, it is allocated to the process at the head of the queue.

The selected process is removed from the Ready Queue and allocated to the CPU.

The algorithm continues until there is no process in the Ready Queue.

# FCFS has the following properties:

Non-preemptive.

The resulting waiting time depends on the order of the processes in the Ready Q.

The main advantage of FCFS lies in its simplicity.

The disadvantages of FCFS may be summarized in the following:

The resulting waiting time is often long.

Unsuitable for time-sharing systems since it is not guaranteed that each user gets its share fairly (since it is non-preemptive)

A convoy effect may occur in case there is a long CPU-burst process with other relatively short CPU-burst ones. This is explained in more details in the next slide.

# FIRST-COME FIRST-SERVED (FCFS) - THE CONVOY EFFECT

The convoy effect occurs when there is an unstable mixture of CPU- and I/O- bound processes.

## Consider the following example:

Assume there are 10 I/O-bound processes & I CPU-bound process in the Ready Q.

The I/O-bound processes, when assigned the CPU, will leave the CPU quickly.

The I/O-bound processes will be in the Wait status at the I/O device.

During this time, the CPU-bound process is assigned to the CPU.

While the CPU-bound process is running, the I/O-bound processes complete the I/O operation, and are re-appended to the Ready Queue waiting for the CPU.

Clearly the I/O device is idle during this time, and the competition is high on the CPU.

The convoy effect means an uneven distribution of processes that compete for the CPU from one side, and the I/O device from the other side.

Therefore, at one time it is found that most processes are competing for the I/O device whereas the CPU is idle.

Suddenly, at a subsequent time most processes that were competing for the I/O device complete the I/O operation and are re-appended to the Ready Queue.

Now, most processes compete for the CPU whereas the I/O device is idle.

This phenomenon is known as the convoy effect.

# FCFS - EXAMPLES

Given:		Result:		
Process Number Burst Time (in ms)		Waiting Time		
PΙ	24	0	Average Waiting Time	
P2	3	24	= (0 + 24 + 27)/3 = 17  ms	
Р3	3	27		

#### Required:

Calculate the average waiting time using the Gantt Chart.

#### Solution:

Given:Result:Process NumberBurst Time (in ms)Waiting TimeP230Average Waiting TimeP33= (0 + 3 + 6)/3 = 3 msP1246

#### Required:

Calculate the average waiting time using the Gantt Chart.

#### Solution:

P2 P3 P1

0 3 6

# **SHORTEST-JOB FIRST (SJF)**

The CPU is assigned to the process that has the smallest CPU burst.

### SIF is implemented as follows:

When the CPU is idle, the process with the least <u>next</u> CPU burst is assigned to the processor.

If two processes have equal next CPU burst, then the FCFS algorithm is applied.

Note that in SJB, the order of the processes in the Ready Queue is irrelevant.

Moreover, the arrival time of the processes to the Ready Queue is also irrelevant.

## SJF has the following properties:

Non-preemptive.

Often used in long-term schedulers.

# The advantages of SJF may be summarized as follows:

Provides smaller average waiting time as compared to the FCFS.

In general, it is proved to be optimal.

The disadvantages of SJF may be summarized as follows:

It is difficult to know the length of the next burst of a process.

Therefore, it is predicted based on historical data of the same process.

# SHORTEST-JOB FIRST (SJF) - EXAMPLE

Given:	Given:		Result:	
Process Number	Next Burst Time (in ms)	Waiting Time		
PI	6	3	Average Waiting Time	
P2	8	16	= (3 + 16 + 9 + 0)/4	
P3	7	9	= 7 ms	
P4	3	0		

# Required:

Calculate the average waiting time using the Gantt Chart.

Solution:				
P4	PI	P3		P2
0	3	9	16	24

# PREDICTION OF THE NEXT CPU BURST (I)

The value of the next CPU burst is calculated from the following formula:

$$\zeta_{n+1} = \alpha t_n + (1 - \alpha) \zeta_n$$

n is the burst number

 $\downarrow \zeta_{n+1}$  is the value of the next burst

 $t_n$  stores the most recent history (ie. most recent burst value)

 $\zeta_n$  stores the past history (ie. burst value of the past history)

is a constant between 0 and 1. It controls the relative weight of the recent and past history.

The constant  $\alpha$  designates the following:

α

If  $\alpha > 0.5$   $\Rightarrow$  more weight is given to the recent history.

This is used when the history is assumed old and irrelevant.

If  $\alpha$  < 0.5  $\Rightarrow$  more weight is given to the past history.

This is used when the current conditions are assumed to be transient.

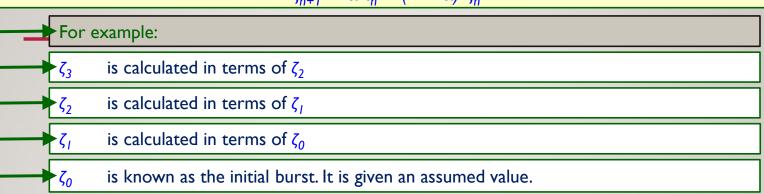
If  $\alpha = 0.5$   $\rightarrow$  equal weight is given to both the recent and past history.

This is the most common case.

# PREDICTION OF THE NEXT CPU BURST (2)

The previously mentioned formula calculates the current burst in terms of the previous one:

$$\zeta_{n+1} = \alpha t_n + (1 - \alpha) \zeta_n$$

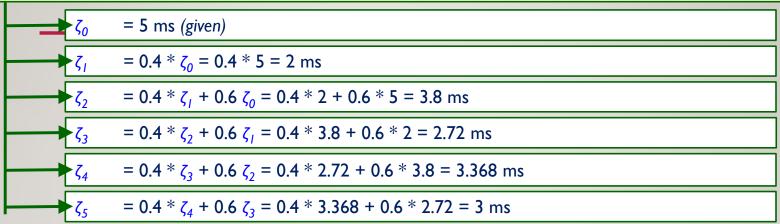


### PREDICTION OF THE NEXT CPU BURST - EXAMPLE

Calculate the 5<sup>th</sup> burst of a process, given that its initial burst is equal to 5 ms, and the recent history is given a weight of 0.4.

Solution: The formula is

$$\zeta_{n+1} = \alpha t_n + (1 - \alpha) \zeta_n$$



# **SHORTEST-REMAINING TIME FIRST (SRTF)**

SRTF is the preemptive version of the SJF.

### SRTF is implemented as follows:

As soon as a new process Pi is appended to the Ready Queue, its next burst time is compared to the remaining time of the currently running process Pj.

If the burst time of Pi is smaller than that of Pj, then Pj is preempted and returned to the Ready Queue.

The CPU is then assigned to Pi.

Therefore, it is important to record the arrival times of all processes.

The arrival time of a process is the time at which it is appended to the Ready Q.

# SRTF has the following properties:

Preemptive.

Often used in long-term schedulers.

The advantage of SRTF is that it is proved to be generally optimal.

The disadvantages of SRTF may be summarized as follows:

Frequent context switching – due to preemption – causes an overhead that cannot be ignored when evaluating the scheduling algorithm.

Like SJF, the value of the next CPU burst should be calculated empirically.

#### **SRTF** - EXAMPLE Given: Process Number Arrival Time Next Burst Time (in ms) PΙ **P2 P3 P4** Remaining Time Solution: Time PΙ **P2 P3 P4 Minimum** 8 NA NA NA 0 NA NA **P2** 9 NA **P2 P2** 5 **P2** 5 \* 9 **P4 P4** 6 9 **P4** \* 8 **P4** 9 9 **P4** 10 9 \* PΙ and so on. The final Gantt chart is shown below:

PI P2 P4 PI P3

# **SRTF** – EXAMPLE (CONT'D)

To calculate the average waiting time, consider the arrival time and the Gantt chart.

#### Given:

P4

Process N	<b>Arrival Time</b>	
PI	0	
P2	1	
Р3	2	

3

PI	P2	P4	PI	P3	
0	1	5	10	17	26

The waiting time of each process is calculated as the difference between the starting time (time at which the process starts to run) and the arrival time (as given).

## **Waiting Time:**

$$PI = 10 - 1 = 9 \text{ ms}$$

$$P2 = I - I = 0 \text{ ms}$$

$$P3 = 17 - 2 = 15 \text{ ms}$$

$$P4 = 5 - 3 = 2 \text{ ms}$$

Average waiting time= (9 + 0 + 15 + 2)/4 = 6.5 ms

If a process is preempted once (or more) then its waiting time is the sum of times the process is waiting in the Ready Queue.

Comparing this result with the non-preemptive algorithm S|F → the

SIF gives an average waiting time of 7.75 ms.