

CPU Scheduling - Short term Scheduling

(SOS) lecture 05

- Cpu scheduling used in multiprogramming systems to utilize the CPU usage in positive manner.

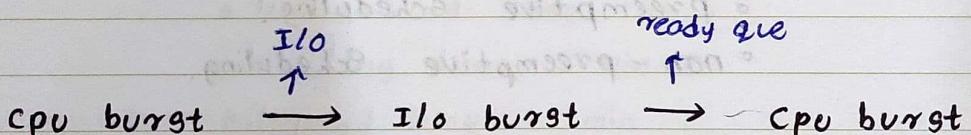
- have 2 types of burst in Scheduling.

• Cpu burst

• I/O burst

- CPU burst - The time took the process to begin executing in the processor.

- I/O burst - When the CPU is waiting for I/O for further execution.



- Start with CPU burst and end with also CPU burst.

- have 2 types of bound measures.

- I/O bound

- CPU bound.

CPU

I/O bound - time took to complete a task determined with speed of CPU. have few but long CPU burst

• I/O bound - time took to complete I/O operation. have short CPU burst (dealing with CPU with short time)

CPU Scheduler. (Short term Scheduler)

- This will select an process from ready que and allocate to CPU.

- CPU Scheduler will select a new process only in,

- Process Switcher from running state to Waiting state (I/O Wait) (Parent wait for child end)
- Switches from running to ready (Interrupt)
- Waiting to ready (after I/O completion).
- Terminator (end of a process).

① and ④ → must select a new process.

② and ③ → can select any.

- CPU Scheduling has 2 types.

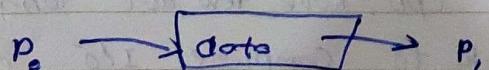
- Preemptive Scheduling.

- Non-preemptive Scheduling.

- Preemptive - has a timer. each process given a time to execute in CPU. (nendays) - 2, 3

- Non-preemptive - CPU only released by the time of termination or to move to waiting que. no timer. (1, 4)

- Preemptive Scheduling can cause race conditions. many processes access the same data at given time



Dispatcher.

- This component give the control of the CPU to the Scheduler that select the current process, C short term Scheduler or CPU Scheduler.
 - The time took to Stop a process & Start another called dispatch latency.

◦ Dispatcher involves,

- When switching context
- running user application
- When restarting a program.

Using below points we can measure the algorithm.

- CPU utilization - max
- Throughput - max
- Turnaround time - min
- Waiting time - min
- Response time - min

refer 6.10 for more.

There are 3 types of CPU scheduling algorithms

- First come, first served (FCFS)
- Shortest job first (SJF)
- Round Robin (RR)

① first come first served (FCFS)

- Which process Coming first will be Served
- if the burst time was short average waiting time will be reduced.

Conway effect - Short processes behind long process
to this solution SJF algorithm came.

② Shortest Job first (SJF)

- SJF gives minimum of waiting time
- In here first have to find the CPU burst.
- In preemptive called Shortest-remaining-time-first.

CPU burst can be predicted using exponential averaging formula.

$$\tau_{n+1} = \alpha t_n + (1-\alpha) \tau_n$$

- if there is no recent history $\alpha = 1$

- if there it will be the value.

- if we don't care $\alpha = 0$.

③ Round Robin algorithm (RR)

o each process gets a small unit of CPU time. (time quantum) (10 - 100 milliseconds).

o in round robin performance depends on a time quantum.

o if que is large \rightarrow FCFQ

o if q is small \rightarrow have to do context switch

o We want quantum to be large.

o good for interactive System.

④ Priority Scheduling.

o based on the priority CPU will allocated and each process given a time number

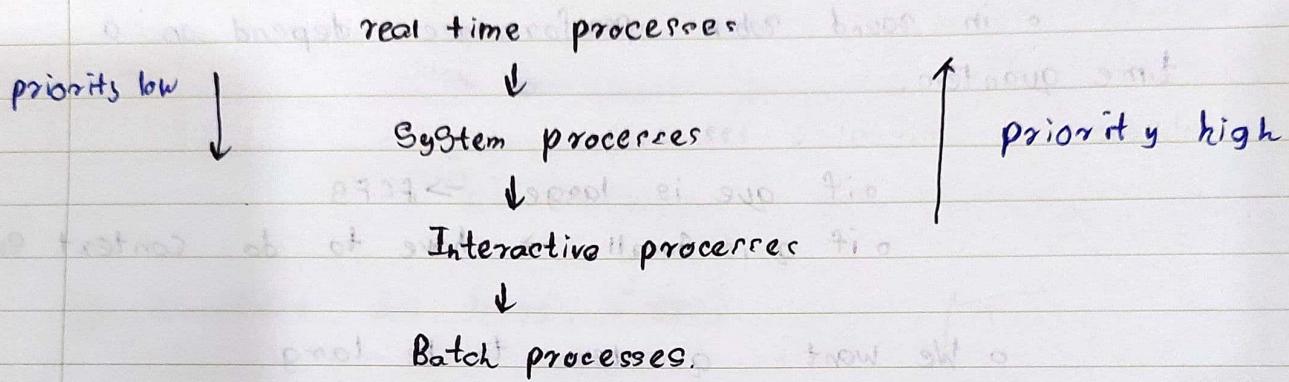
o Smallest priority \rightarrow highest integer and vice versa.

o SJF is priority Scheduling method.

Starvation - low priority processes may not execute
Solutions (aging) - increase the priority of the process.

⑤ Multilevel que.

- We can have Separate que for each priority. Prioritization is based on process type.



⑥ Multilevel feedback que.

- here processes can move between que.

by this method priority que can be moved to higher priority que.

Thread Scheduling.

there are 2 types of threads

• user-level threads

• kernel-level threads.

* if the threads support then threads will be scheduled not the processes.

o Thread library will schedule user-level threads to light weight processes (Lwp). also known as process contention scope (pcs)

↳ System contention scope (sc) will decide which kernel thread schedules into CPU.

~~Priority Based Scheduling~~

Multiprocessor Scheduling

- o This is used when a system has more than a single processor.
- o Multiprocessors can be:
 - o Multicore CPU's
 - o Multi-threaded core's
 - o NUMA system
 - o Heterogeneous multiprocessing

Symmetric multiprocessing (SMP)

Processor schedules itself. has a common ready que. also has a private que of threads

refer 5.38 for more, 5.39

Multithreaded Multicore System.

- o each core has a thread, if one thread has not enough space it will switch to another thread.
- o Chip multi-threading (CMT) - assigns each core multiple threads.

Multicore - this system has 2 levels of Scheduling.

- o OS decides which software thread to run on CPU.
- o each core decides which hardware to run on physical core.

Load balancing.

- o how the work load is distributed among processors. 2 ways to do it.
- o Push migration
- o Pull migration.

Processor affinity.

- o When thread running on one processor the cache of that processor information of processor stored by that threads.

o Load balancing effect of affinity. 2 types of affinity.

- o soft affinity
- o hard affinity

5.43

Real time Cpu Scheduling.

- have 2 types of real time system.
 - soft real time 5.45
 - hard real time.

Event latency.

- the time took to event occurse to being serviced. 2 types of latency.
 - interrupt latency 5.46
 - Dispatch latency.

Algorithm Evaluation.

- Deterministic modeling.
 - We have predefined work load and test it on each algorithm. (performance).
- Queueing model.
 - Compute the throughput, utilization, waiting estimate the cpu burst and arrival time distribution.
- Little's formula.
 - Valid for any algorithm and for any arrival distribution.
- Simulation.
 - to get more accurate result we can use this.
- implementation