

Scheduling

What is process Scheduling?
& describe different types of Scheduler?

⇒ It is the task of the OS to Schedule the processes of different states like Ready, Running, Waiting... According to Scheduling algorithm with the help of different Schedulers.

⇒ Process Scheduling allocates the CPU time to each process for execution.

Why do we need?

○ For CPU utilization in multi-programming & multi-tasking os we need Scheduling.

Def: The act of determining which program is in the ready state and should be moved to the running state is known as process scheduling.

⇒ The number of processes is more than resources we have only one processes, but processes are more than one.

⇒ These process are not sensible enough the schedule themselves one after another so the OS take this responsibility.

⇒ OS Schedule these process using the scheduling algorithms.

The main purpose of process scheduling is to keep the CPU busy all the time and to deliver minimum response time for all programs.

Scheduling Queues:

⇒ A place where the processes waits to get the resources are known as job queue.

⇒ There are three kinds of queues:

- o Job Queue

- o Ready Queue

- o Waiting Queue

Job queue:

- o When a process enters in the system, it is placed in a queue. That queue is called job queue.

- o These queue of the process resides in the secondary memory.

Ready Queue:

- o The processes that are placed in the main memory and waiting for the CPU time are called ready queue.

Waiting queue:

- o If the process is waiting for some kind of I/O operations it will be kept in waiting queue & wait here.

Process migration:

o A process begins in the job queue & moves b/w the ready queue and waiting queue and eventually finishes the task. This process of moving b/w the queue is called process migration.

What is Scheduler?

⇒ The part of OS that "allocates the computer resources" to the processes is known as scheduler.

⇒ It uses "Scheduling algorithms" to decide which process it must allot to the CPU.

Types of Schedulers

⇒ There are three types of schedulers.

- o Long-term Scheduler

- o Short-term Scheduler

- o medium term Scheduler

1) Long-term Scheduler-

- o Responsible for managing the processes in a job queue which is in the secondary memory.
- o It checks which process wants to go first priorities that process in the list
- o Also controls the no. of processes in the memory is called the degree of multi-programming.

Steps it performs:

- o Get a new job, check its priority, put it in the queue, give it to the ready queue.
- o Also known as job scheduler.
- o It has less scope in time-sharing OS.
- o It executes relatively infrequently.

Use:

- o It is mainly present in batch OS like IBM 1360 mainframe computers.
- o It is not present in 'Windows' OS.

→ It brings "new process to ready state".

Short-term Scheduler:-

- o Responsible for managing the processes in the ready queue which is in the main

main memory.

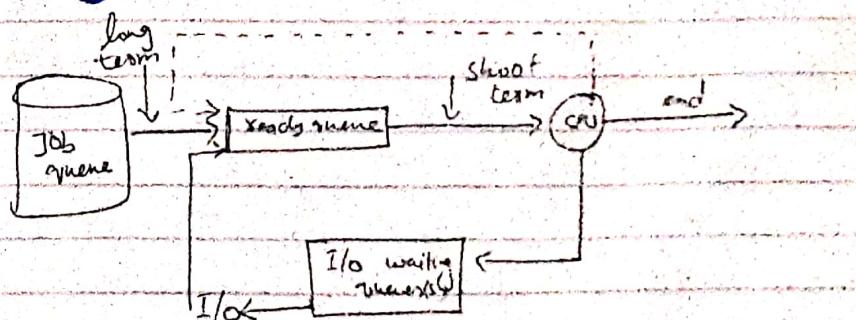
- o Processor also gets the processes from this queue because it is near the processor. It is also called CPU scheduler.

- o In Simple, select a process from ready queue to running state of CPU.
- o In some OS it is called dispatcher which is responsible for

- o Switching to user mode
- o Context switching.

- o It is also invoked when an event occurs which may lead to interruption of the current process.
- o Highly frequently executes.

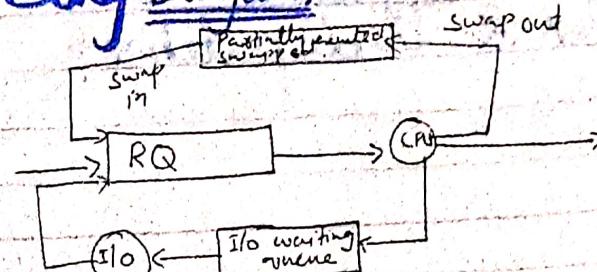
Diagram:



iii) Medium term Schedules:-

- => Responsible to keep the things in the smooth flow.
- => it smooth by keeping the ways of the processes clear.
 - o who blocks the ways other process?
 - o A process may come into the main memory and may do nothing other than keeping the memory free.
- => medium term schedules take this process out of the main memory and place it in the secondary memory.
- => Sometime a process of high priority want to come in the main memory, but the memory is not enough because it is reserved by some low priority process. medium term play its role swap that low priority process out and make space for high priority.
- Perform swapping in & swapping out. Responsible for suspending & resuming the process.

Long Diagram



PP-SQ (2014-2020)

- o Explain long term schedules.
- o Difference between short-term & medium term schedules.

(Long)

- # Describe the differences b/w long term, short-term and medium term schedules.

What is difference b/w Schedules & Dispatches?

Dispatches

- Def: A dispatcher refers to a module that provides the control of the CPU to that process that gets selected by the short-term - schedules.

Dependencies:-

- The overall working of a dispatcher depends entirely on the scheduler. It means that the dispatcher has to wait until & until a scheduler performs the selection of processes.

Types:-

- has no major type. it is a type of code segment.

Schedules

- it is basically a system or software that helps in selecting a process out of various available processes.

- The overall working of a scheduler occurs independently. It happens pretty immediately as & when required.

- Basically, there are three types of schedules.

Time taken:

- Dispatcher latency refers to the time taken by a dispatcher to stop one process and to start another process.

Algorithm

- No specific algorithm for the implementation of a dispatcher.

Functions:

- Switching of Context
- Switching to user mode
- Jumping into the user program's process location for restarting that given program

- The overall time taken by a scheduler is very negligible & neglectable.

- A typical scheduler can work on various types of algorithm like LL, SJF, FCFS etc.

- Perform the selection of various processes.

Modes of Scheduling / Types

Non-Preemptive

Def:

- In non-preemptive mode, once it's a process enters into a running state, it continues to execute until it terminates or blocks itself to wait for I/O or by requesting some OS services.

Basics:

- Allocation of resources to process holds till terminated or switch to waiting state.

Interrupt:

- Process cannot be interrupted.

Preemptive

- In preemptive mode currently running process may be interrupted and moved to the ready state by the OS.

- Allocation of resources to process for limited time.

- Process can be interrupted.

Starvation:

- If process of long burst time is running then low priority process has to starve.

Flexibility:

- Not flexible.

Cost:

- Not cost counter switching overhead.

- It has no overhead of switching process.

System use:

- Mainly batch system use.

CPU utilization:

- Low CPU utilization.

Example:

- FCFS
- SJF

352CFIR

- But short burst process has more share.

flexible:

- cost associated.

- It has overhead of switching process.

Time sharing:

- System use.

High:

- Round Robin
- SRFT
- Priority

Terminology used in scheduling

Algorithm:

1) Arrival time:

- Time at which process arrives in ready queue.
- Arrival time can vary from process to process.

remain memory \Rightarrow CPU Process and

2) Completion time:

- Time at which process completes its execution.
- process can be left \Rightarrow CPU

3) Burst time: /Execution time/ Running time -

- Time required by the process for its CPU execution.
- Total time going for 6 processes

4) Response time:

- Time it takes to start responding from arrival time.

CPU responds \Rightarrow when process come for it is the chance b/w the arrival of a CPU and the first time

A Time taken by a process since it enters a ready queue till the process execution is completed.

5) Turn around time:

- How long a process takes to execute it. Limited by the speed of the CPU device.
- Calculated by the following formula
Turn around time = Completion time - Arrival time

$$SAT = CT - AT$$

$$TAT = WT + BT$$

Use time of CPU

6) Waiting time:

- Amount of time that a process spends while waiting in a ready queue until it gets the CPU (for its turn)
- Sum of all times spent in waiting queue.
- it is calculated by the following formula: it is not limited by the speed of CPU

$$WT = TAT - BT$$

$$WT = \frac{CT}{\text{Arrival time}} - \frac{BT}{\text{Arrival time}} = \frac{CT - BT}{\text{Arrival time}}$$

Objective of CPU Scheduling

Maximize

- CPU utilization

- Throughput

Minimize

- Turnaround time

- Waiting time

- Response time

CPU Scheduling Criteria

- Every CPU algorithm has different properties. So criteria used for computing these algorithms.

i) CPU Utilization

- keep CPU as busy as possible.

Range: 0 - 100%.

in practice: 40% - 90%.

ii) Throughput:

- Rate at which processes are completed per unit time.
- Increases for short processes.
- Decreases for processes of huge size.

iii) Average turn around time:-

- Sum of all TAT of Procn
No. of procn.

iv) Average waiting time:

- Sum of all waiting time of Procn
No. of procn.

Average response time:

- Sum of all response time of Procn
No. of procn.

b) Factors of CPU

Share in fairly manner with each process so that process does not go into starvation.

Scheduling algorithm

- A scheduling algorithm is the algorithm which tells us how much CPU time we can allocate to the processes and when.

i) FIFO Come, First served Scheduling

- CPU Scheduling algorithm is the FCFS algorithm.
- The process that requests the CPU first is allocated the CPU first.
- It can be managed with FIFO rule.
- + working from back.

Pros:

- Simple method
- No starvation
- Less overhead
- Better for long processes
- Non-preemptive

Cons

- Avg. wt, response time is poor.
- Convey effect
- Criteria: "Arrival time."

Example:-

Process	AT	BT	FT	TAT	WT
P ₁	2	6	17	17 - 2 = 15	15 - 6 = 9
P ₂	5	3	24	24 - 5 = 19	19 - 3 = 16
P ₃	1	8	11	11 - 1 = 10	10 - 8 = 2
P ₄	0	3	3	3 - 0 = 3	3 - 3 = 0
P ₅	4	4	21	21 - 4 = 17	17 - 4 = 13

Process	AT	BT	FT	TAT	WT
P ₁	2	6	17	17 - 2 = 15	15 - 6 = 9
P ₂	5	3	24	24 - 5 = 19	19 - 3 = 16
P ₃	1	8	11	11 - 1 = 10	10 - 8 = 2
P ₄	0	3	3	3 - 0 = 3	3 - 3 = 0
P ₅	4	4	21	21 - 4 = 17	17 - 4 = 13

$$AWT = \frac{0 + 0 + 8 + 4 + 3}{5} = \frac{9}{5}$$

Note: In non-preemptive scheduling algorithm WT and response time are same.

(iii) Shortest-Job-First (SJF)

\Rightarrow It outcomes the problem of FCFS

$$AWT: \frac{9+16+2+0+13}{5} = \frac{40}{5} = 8$$

Ex 8:

Process	AT	BT	FT	TAT	WT	RT
P ₁	2	2	4	4 - 2 = 2	2 - 2 = 0	0
P ₂	0	1	1	1 - 0 = 1	0	0
P ₃	2	3	7	7 - 2 = 5	2	2
P ₄	3	5	12	12 - 3 = 9	4	4
P ₅	4	9	16	16 - 4 = 12	3	3

\Rightarrow In SJF processes are scheduled according to the burst time of these processes.

\Rightarrow This is the best approach to minimize waiting time.

\Rightarrow If two processes have the same burst time the FCFS should be applied.

\Rightarrow This Scheduling algorithm is optimal if all processes are available at the same time. (unless arrival time is 0 as all)

Wait time:

	P ₁	P ₃	P ₄	P ₅
Wait time:	2	4	7	12

Arrival time is same for all)

Types:-

- o Preemptive (SRTF/SPN).
- o Non-preemptive (FIFO by default)

Non-preemptive:

⇒ Once the CPU is given a process it cannot be preempted until the current CPU burst finishes.

Algorithm:

- o Sort all the processes according to the arrival time.
- o Then select that process that has minimum arrival time and minimum burst time.
- o After completion of the process make a pool of processes that arrives afterwards till the completion of the previous process and select that process among the pool which is having minimum burst time.

Problem with non-preemptive:-

- o if the arrival time for processes are different, which means all the processes are not available in the ready queue at time 0, and some jobs

arrive after some times in such situation, sometimes process with short burst time have to wait for the current process execution to finish because in non-preemptive SJF, on arrival of a process with short duration, the existing process execution is not halted/stopped to execute the short job first.
o This leads to the problem of "starvation", where a shorter process has to wait for a long time until the current longer process gets executed. This happens if shorter jobs keep coming, but this can be solved using the concept of aging.

Preemptive Shortest Job First:-

- o In preemptive shortest job first scheduling, jobs are put into ready queue as they arrive, but as a process with short burst time arrives, the existing process is preempted or removed from execution, and the shorter job is given the CPU.

Priority Scheduling

- Each process has its own priority (highest priority)
- All available processes (highest priority) gets the CPU.
- If the user uses FCFS.

Model: Preemptive & Non-preemptive.

Priority: External & Internal.

- Users the number higher the priority

Pros

- Every process gets response for the highest priority.

Cons

- Indefinite blocking or starvation.
- Solution:

Aging: increasing the priority of low priority process based on the time it's waiting.

mode: preemptive.

Pros:

- Fair amount of CPU is allocated for each job.
- It is not affected by context switch.

Round Robin algorithm:

- Each selected process is assigned a time interval called quantum or time slice.

Cons

- If a process is blocked or terminated before the time slice has elapsed then switching is done and other process is scheduled to run.
- If process needs the CPU least longer than 'time slice' then the process is running at the end of time quantum.

Now that process will preempt move to end of ready queue and CPU will allocated to process which instant of ready queue.

Cons:

- Leads to starvation for processes with less burst time on the hand to repeat the cycle many times.

- o Round Robin scheduling is FCF's scheduling with preemptive mode.

Note:

- \Rightarrow with decreasing value of time quantum
- o No of context switch increase
 - o Response time decrease
 - o Chances of starvation decreases
 - o Thus, smaller value of time quantum is better in terms of response time.
 - o Time quantum should be such that is neither too big nor too small.
PP-quantum

(iii) Ans

efficient resource allocation

optimized throughput

fairness

effective resource management

P ₁	P ₂	P ₃	P ₄	P ₅	P ₆	P ₇	P ₈	P ₉	P ₁₀	P ₁₁	P ₁₂	P ₁₃	P ₁₄	P ₁₅	P ₁₆														
2	4	6	8	9	11	13	14	16	18	20	22	23	25	26	28	30	31	32	33	34	35	36	37	38	39	40	41	42	43

Q: What is processor affinity?

⇒ Processor affinity refers to the assignment of specific processes or threads to run on particular CPU cores or processors with a multicore system.

⇒ it allows for control over CPU resource allocation

⇒ Processor affinity can be set at the process or thread level, enabling processes or threads to consistently execute on any cores.

⇒ Benefits of processor affinity:

- Performance optimization by reducing cache thrashing
- Control over resource allocation

(i) (PP-2017)

What is Starvation?

⇒ Starvation refers to a situation where a process or task is unable to make progress or obtain necessary resources, even if it is ready to execute.

⇒ It occurs when certain processes or tasks are consistently favored over others, leading to a deprivation of resources for the starved processes.

⇒ This can be caused by

- o Priority inversion
- o unfair resource allocation
- o flaws in scheduling algorithm
- o Deadlocks

⇒ Starvation can negatively impact system performance and fairness, resulting in lower throughput and longer response times.

⇒ To mitigate starvation, techniques like

- o Priority aging
- o RR time slicing
- o fair resource allocation strategies are employed

P.

(How to create thread? etc.)

Q.

Past Papers Short Questions

(Ch #5)

- (i) What does it mean to preempt a process? (Repeat)
- (ii) Difference b/w preemptive and non-preemptive Scheduling? (Repeat)
- (iii) What advantages is there in having different quantum sizes at different level of a multi-level queuing system?
- (iv) Which one of the following scheduling algorithm could result in Starvation?
FCFS, SJF, Round Robin.? (Repeat)
- (v) Different types of real time Scheduling? (See)
- (vi) What is the processor Affinity? (Repeat)
- (vii) What are turnaround & response times?
+ waiting times? (Repeat)
- (viii) What is meant by the term dispatch latency?
- (ix) Define Starvation?
- (x) List the various criteria must be considered in making a good Scheduling algorithm?
- (xi) What is RR Scheduling?