

- The state of the thread can be changed any no. of times.

## \* Scheduling Policy:

- 5 diff. sched algorithms
  - i). FIFO ii). Round Robin iii). Others.
- Real-time      -Real time.      -Non Realtime

- Total 160 priorities available - 0 to 127.

- 0-99 - Real Time Priority.

- 100-139 - Non Real-time Priority.

- Round Robin & FIFO are implemented using the same data structure.

- Round Robin with no time quantum is FIFO.

- In Round Robin a new process joins at the back of the queue and in FIFO a process joins at the front of the queue.

- So in FIFO the same process keeps on executing until the process is completely executed.

- A non-real time thread will always have a non-real time & algorithm so no need to change the policy.

- These algorithms are only for processes in the same priority no.

- Default is others.

- For real time threads/processes the priority value need also be set like 0 and 99.
- All these are done using attribute function.

25/09/19 Nice Value -20 → +19

- Nice value + 100 → Gives the priority no. for Non Real time processes (From 100 to 139).

- For real time process priority value need to be set by user but for non-real time process the priority is calculated automatically using nice value.

- Each priority no. maintains a separate linked list of processes.
- Each queue can use diff. priority algo.

- CPU Protection Timer
  - Time slice expiry.
  - In FFO any algo. time slice expiry is a decision point.
  - In FFO, the processes join at the start of the queue after time slice expiry.

- pthead\_attr\_struct \* schedparam - Set the priority value for a process.

- Structure is used.
- Pass the add / pointer of the structure.
- Default value: Priority no 0.

- \* Priority after scheduling param — get the priority value of a process.
- \* Inherit sched — By default new threads have to given their scheduling policy by themselves.
  - The new threads can also inherit the policy from the parent thread.
- \* Scope of a thread — In linear there is only one scope — system scope.
- \* pthread\_scope\_system — All the threads compete with all other threads of the complete system, because each thread is a process.
- \* pthread\_scope\_process — CPU is given a particular process and any thread cor from that process can execute acc. to scheduling policy.
- \* All threads — Threads compete within its own process only because each thread does not make a process.
- \* Set Sched Param —
  - Attributes of any thread can be changed after creation but not with the attribute functions.

- pthread\_set sched param - Changes the scheduling and priority of a particular thread by giving id as a arg.
- It is not necessary to call the function from the thread which needs to be changed, can be called from any function.
- \* Self & Equal
  - pthread self - Returns the id of calling function. thread.
  - pthread\_equal - Check wether two thread identifiers refer to the same thread.
- \* Detach - Change the thread to detach after a thread is created as joinable.
  - A detached thread cannot be joined again after creation.
  - pthread\_detach (pthread + th) -
    - works only if no thread has made this thread joinable.
    - If another thread has joined this thread then detach will not work.
- \* Exit - pthread\_exit (void \*retval).
  - Return value of thread is void \*
  - So, as many variables can be stored in

heap and the add. of heap can be passed.

- Stack gets deleted after thread is exited so don't pass add. of stack.

- Heap is not cleared by returning from a function.  
It needs to be cleared by the user.

• Finalisation Function - If a thread is stuck somewhere then this function clears thread from all those places.

- Joinable thread has to give the it's return value/structure to other thread. Till then it's not removed.

- Detached thread doesn't need to do so.

### Join - (thread join)

- Blocks the thread until some other thread (passed hd. by argument) has finished execution.

• void \* thread return - Add. which stores add. of some other variable.

- The thread which dies passes the add. of its return value to the calling thread.

10/19

### \* Threading Issues:

- The scope of thread is always till its parent process only and not above the level.
- If exec() is executed in a thread then the process which created thread will die (In Linux) - Because the kill signal goes to the process group id, which the thread belongs to.
- Usually all the thread will get killed if any one thread executes exec() because group leader will gets killed.

### \* Create Threads -

- clone() - system call.
  - do\_fork() - Function in kernel.
- File handles are shared b/w all the threads so no. of files opened will not be increased by creating a new thread.

### \* Implicit Threading -

- Thread Pool - Some threads already created and ready to be used to speed up sys.
  - Parallel Lang. - OpenMP, MPI, Cilk
- After the use of thread is completed it is kept back in the pool.

## \* Signal Handling -

- One signal handler for the process and all its threads but diff. signal mask for all threads.
- User defined signals can be masked but not critical signals.
- Signal is given to the process and not the threads.

. TASK KILL - wait, suspend, blocked state with can be sleep stopped with only SIGKILL signal.

- 90% of the time pthread\_exit() will release all the data structures of thread. Only 10% of time the threads are joinable.

## \* Thread Cancellation -

- Asynch - Thread dies immediately.
- Deferred - Thread dies after reaching a (synch) certain cancellation pt. even though the signal was present before.
- Pthread libraries support deferred cancellation.
- pthread\_join() is a cancellation point.
- Deferred is preferred over asynch.

9/10/19

### A CPU Scheduling :-

\* Long Time Scheduler - Admit prog. from user to ready.

- Improves multiprogramming.
- Few or like Linux does not have long time scheduling.

→ More jobs available ready for execution.

- But less resource available for each job because each process takes some resources to be ready.

- very rare.

\* Short term Scheduler - From ready to running (dispatches)

- 4 major decision pts - i). Time slice expiry
- ii). Process Completion
- iii). Process moves to waiting state from running.
- iv). Higher priority job arrived. - switches from waiting/ready to running.

- (ii) & (iii) are voluntary context switch. - Preemptive

- (i) & (iv) are involuntary context switch. - Non Preemptive

- A non-preemptive scheduler will have only (i) & (ii) as decision pts.

- A preemptive scheduler all of the above are decision pts.

- Most frequently run scheduler.

- Scheduler is also a process so it has highest priority of 0.

- \* Middle Term Scheduler - From system to swap/suspend state, or vice versa

- Reduces multiprogramming.
- Increases resource available for each process.

- Usually by adding a new job the CPU utilization increases.
- But if by adding a new job the CPU utilization decreases then a process needs to be suspended and mid-term scheduler is called.

- + Background Process - Usually only CPU process or very less I/O operation.

- . I/O Bound Process - More time in I/O than in CPU for computation.

- Most systems are designed for I/O bound process.
- Responsiveness is needed i.e. process should be given enough time.

- Preemptive scheduling takes more time than non-preemptive scheduling because to check for whenever a new process joins ready state it has to be checked for priority.

- + Dispatch latency - Time taken to switch a process.

11/10/19  
★

- Stopping the set of currently running process and loading the new process.
- Fixes Non productive work done in dispatching latency.
- That's why quantum time should not be too low - atleast 1000 times that of dispatching latency.
- \* Response Time - Time taken b/w the submission of a req./ process until the first response/output come.
  - Not see the completion of the job.
- \* Turn Around Time - Total time taken to complete the job after its arrival/in the submission.
  - When scheduling as shortest process first, the response time decrease.
  - User Oriented - less response time needed.
    - More responsive.
  - System Oriented - Effective and efficient utilization of the processor.
  - Server machines are designed for 24x7 operation.
- \* Waiting Time - waiting in ready queue and not in waiting state.

running process

done in a

not to  
that of dispa-

the submit-  
process until  
one.

Db.

run to  
b after its

is first, the

needed.

cient  
processor.

= Operatia

e and

11/10/19

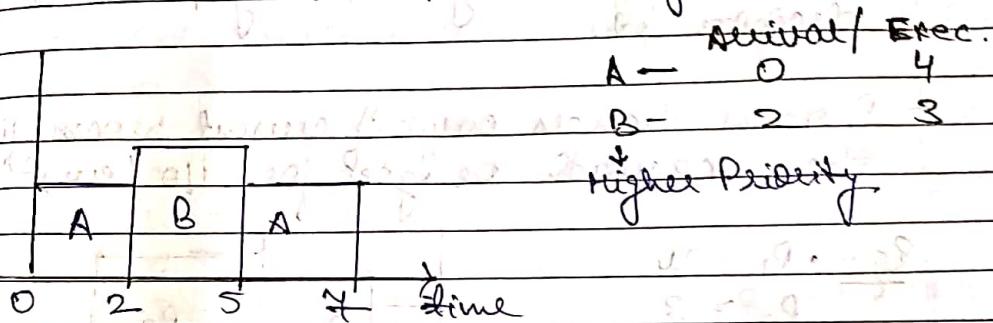
## CPU Scheduling -

### CPU Utilization -

Prod. time - executing prg.

unprod. time - Management work.

Idle Time - Can be used to improve performance  
like power management.



- Turn Around Time (A) =  $7-0 = 7$   
 " " " (B) =  $5-2 = 3$

- Normalised Turn Around Time =  $\frac{\text{Turn Around Time}}{\text{Execution Time}}$

- Norm. Turn Around Time (A) =  $7/4 \approx 1.75$   
 " " " (B) =  $3/3 = 1$

- Best Norm turn around time is 1.  
 Norm turn around time  $\geq 1$ .

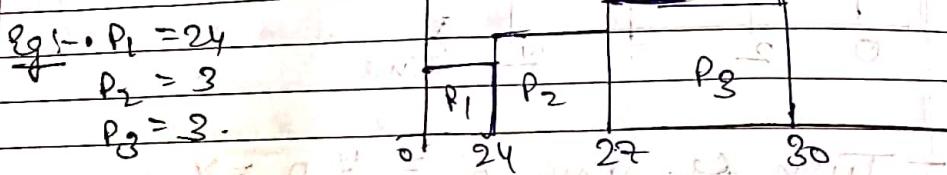
- Waiting time (A) = 3  
 " " (B) = 0

- Response time (A) = 0  
 " " (B) = 0.

Waiting Time = Turn A. Time - Exec. time  
 in ready queue - I/O time

\* First-Come-First-Serve

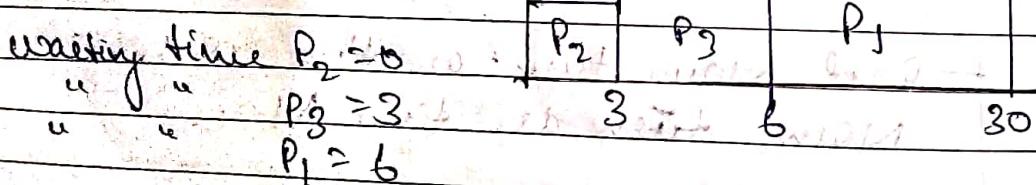
- Simplest algo.
- uses simple queue.
- By default non-preemptive algo. i.e. time slice expiry or higher priority process is not a decision pt.
- Process changes only if current process finishes the execution or goes for I/O / event / signal.



Waiting time  $P_1 = 0$   
 " "  $P_2 = 24$   
 " "  $P_3 = 24$

Avg. waiting time =  $\frac{24+24}{3} = 16$

$P_2 \rightarrow P_3 \rightarrow P_1$



Avg. waiting time =  $\frac{(0+3+6)}{3} = 3$

- Execution time
- I/O time
- Shorter duration
- Priority is not considered at all.
- Conveyer Effect - Small processes have to wait for long time due to longer for processes coming before it.
- I/O bound process suffer because I/O processes are more in no. and shorter duration while as CPU processes are longer but less in no.

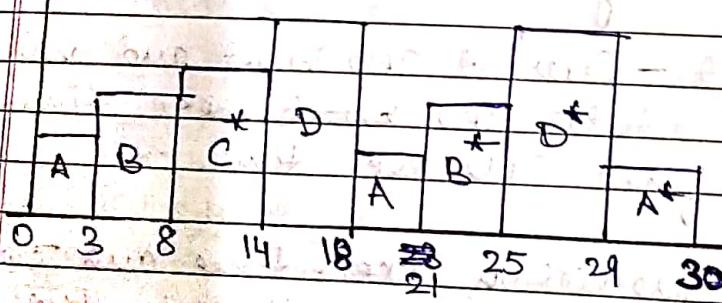
	Average time (avg)	Exec time (total)
A	0	7 = $3C + 5I/O + 3C + 5I/O + 1C$
B	2	9 = $5C + 2I/O + 4C$
C	4	6 = $6C$
D	6	8 = $4C + 1I/O + 4C$

Arrival time      Ready queue

0 Exec time      → 3C 5C 6C 4C 3C 4C 4C 1C

Ready queue      → A B C D | A B D A

Arrival time      → 0 3 2 4 6 8 10 19 26



→ Response Time (A) = 0

$$\text{(")} \quad (B) = 3 - 2 = 1$$

$$\text{(")} \quad (C) = 8 - 4 = 4$$

$$\text{(")} \quad (D) = 14 - 6 = 8$$

14/10/19

- Turnaround time (A) =  $30 - 0 = 30$   
 (B) =  $25 - 2 = 23$   
 (C) =  $14 - 4 = 10$   
 (D) =  $29 - 6 = 23$

- Waittime Time (A) =  $30 - 7 - 10 = 13$   
 (B) =  $23 - 9 - 2 = 12$   
 (C) =  $10 - 6 = 4$   
 (D) =  $23 - 8 - 1 = 14$

- Norm TAT (A) =  $30/2$   
 (B) =  $23/9$   
 (C) =  $10/6$   
 (D) =  $23/8$

#### \* Shortest Job First :-

- Non-Preemptive type - Shortest job first.

Preemptive Type - Shortest remaining time, first.

• CPU Burst - Current CPU burst and not the total CPU exec. time.

- Only the first CPU exec. time.

- Based on prediction, exec. time and I/O time cannot be known before the complete exec. of prog.

- Gives min. avg waiting time.

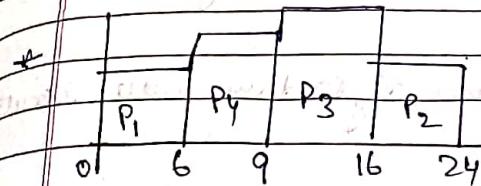
- Picks the process which has shortest next CPU burst.

$$\begin{aligned}0 &= 30 \\-2 &= 22 \\4 &= 10 \\6 &= 23\end{aligned}$$

$$\begin{aligned}3 - 10 &= 13 \\&= 12\end{aligned}$$

14

- Starvation - If jobs with smaller CPU time keeps on coming so the process with higher CPU times does not get a chance or it starves.



Arrival time      Burst time

P <sub>1</sub>	0	6
P <sub>2</sub>	1	8
P <sub>3</sub>	2	7
P <sub>4</sub>	3	3

$$\text{Avg. waiting time} = (0+15+7+3)/4 = 6.25$$

- For preemptive type (shortest job first), the decision pt. is only termination of current process or process goes to I/O.
- Even if a shorter job comes in blue, then also the current process will not stop.

### \* CPU Burst

$$T_{n+1} = \kappa T_n + (1-\kappa) T_n ; \quad 0 \leq \kappa \leq 1.$$

T<sub>n</sub> - Actual length of n<sup>th</sup> CPU burst

T<sub>n</sub> - Predicted value of last " "

T<sub>n+1</sub> - " " " " next " "

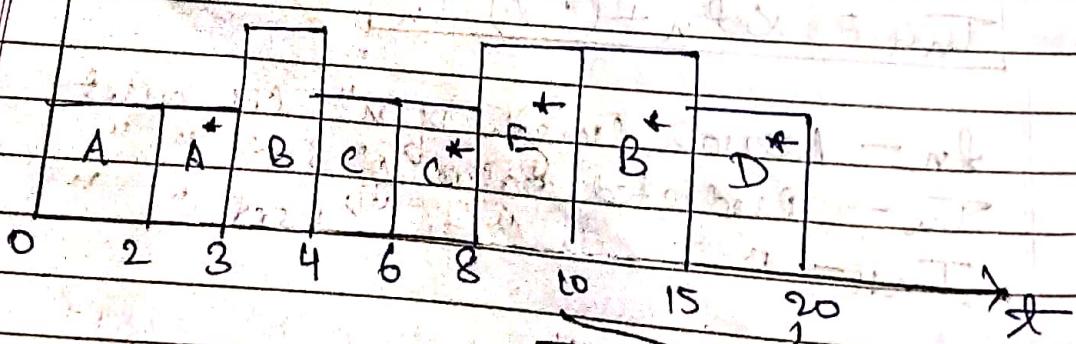
Generally  $\kappa = 1/2$

- If  $K=0$ , next CPU burst prediction depends only on the previous history.
- If  $K=1$ , prediction depends only on the last CPU burst.
- Any priority driven algorithm, will result in starvation.
- FIFO & Round Robin are not priority driven algo., they are queue based algo. — Non-starvation.

### \* Shortest Remaining Time First :-

- Preemptive time sharing algorithm with
- Decision Pt. - Termination, I/O, Arrival of new process.

Process	Arr. Time	Ex. Time	Time ended
A	0	3	3
B	2	6	8
C	4	4	12
D	6	5	11
E	8	2	13



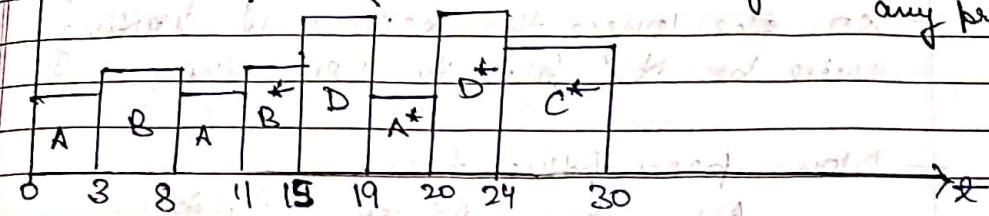
Tie break b/w B & D, so we choose B because it had already executed for some time before.

### Aer. Exe.

eg 1	A	0	7	$- 3C + 5T + 3C + 5T + 1C$
$\Rightarrow$	B	2	9	$- 5C + 2T + 4C$
$\Rightarrow$	C	4	6	$- 6C$
$\Rightarrow$	D	6	8	$- 4C + 1T + 4C$

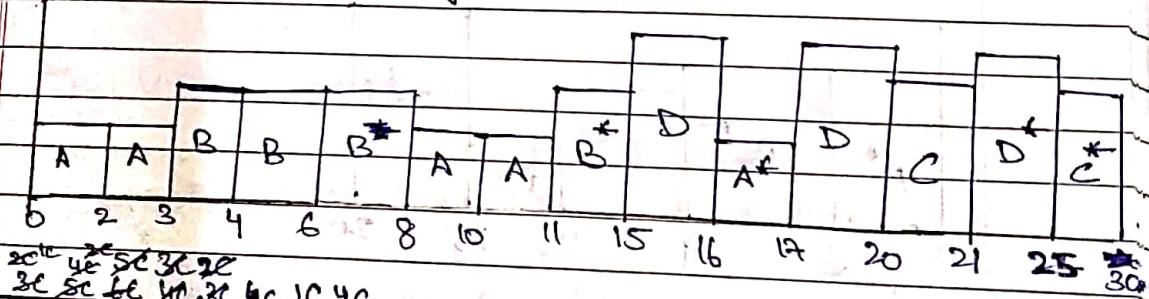
### - Shortest Job First (SJF) -

This can be exchanged without any prob.



Arrival time → 3C 8 6C 4C 3C 4C 10 4C  
Ready queue → A | B | C | D | A | B | C | D |  
Served time → 0 2 4 6 8 10 15 20

### - Shortest Remaining Time First (SRTF).



The no. of decision pts. in SRTF is greater than in SJF but avg. waiting time is less.

At each decision pt. there will be context switching because schedule is called even if the same process is continuing operation.

→ & D, so we it had already come before.

15/10/19 P

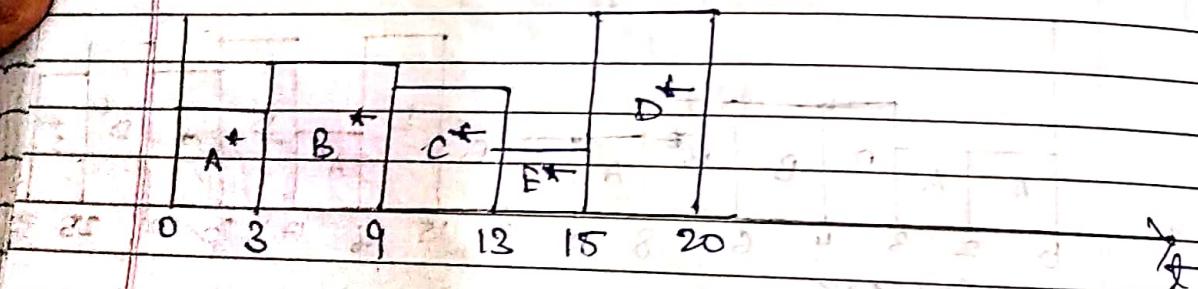
### \* Highest Response Ratio Next (HRRN) :-

$$\text{Response Ratio} = \frac{\text{Time Spent waiting}}{\text{Expected service Time}} + \frac{\text{Expected Remaining Time}}{\text{Expected service Time}}$$

- Ratio increases with waiting of the process.
- so the longer the process is waiting, higher will be its priority next time.

- Non-preemptive type.

	Arr.	Exc.	Response Ratio
A	0	3	$\frac{0+0}{3} = 0 \checkmark$
B	2	6	$\frac{2+3}{6} = 1 \frac{1}{2} \checkmark$
C	4	4	$\frac{4+2}{4} = 3 \frac{1}{2} \checkmark$
D	6	5	$\frac{6+5}{5} = 2 \frac{1}{5}$
E	8	2	$\frac{8+2}{2} = 5 \frac{1}{2} \checkmark$



- E waited less than D because if execution time is less than waiting time has greater impact on priority.
- so this algo. gives priority to process with small CPU burst in blue. To. marks
- when a process goes to I/O and comes back

Page No.: 1 / 1

Page No.: 1 / 1

+ (RRN) :-

+ Expected Review  
Time  
and service Time.

of the process.

Waiting time - higher

me. -  $\frac{1}{T}$

then waiting time is calculated from the time of its 1<sup>st</sup> arrival and not from the start.

### 10/14 Priority Scheduling :-

- starvation in both preemptive and non-preemptive

- Aging - As the waiting time increases, the priority of the process increases.

• Solution to starvation.

• LFU :- E.g. - Count incremented each time the program is used.

- Multiple priority queue - each queue representing 1 priority.

- In Linux - Priority = 140 queue total,

- 0 to 139 - lowest

Highest Priority

• If queue no. n doesn't have any process then it checks in (n+1)<sup>th</sup> queue.

Eg:-

Arr. Ex. Priority

P1 0 3 5 8

P2 1 2 3

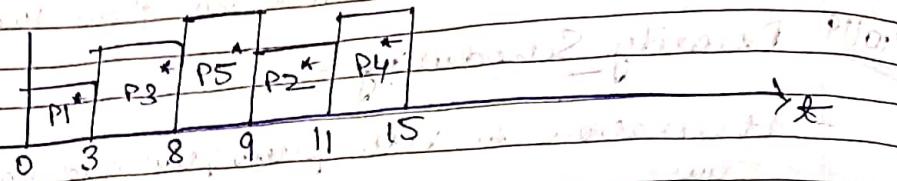
P3 3 5 2

P4 4 4 4

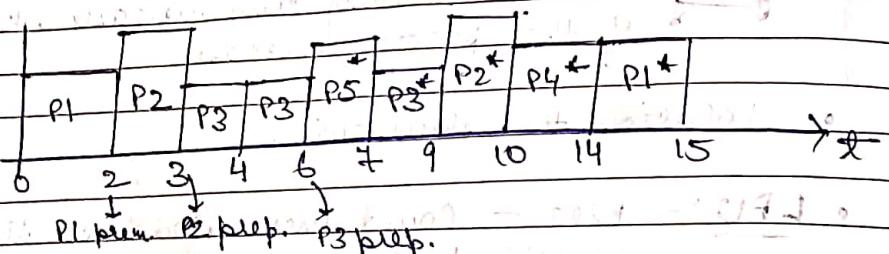
P5 6 1 1

- Lesser the no higher the priority.

- Non-Preemptive - No preemption, only completion



- Preemptive -



Eg: All. Ex.

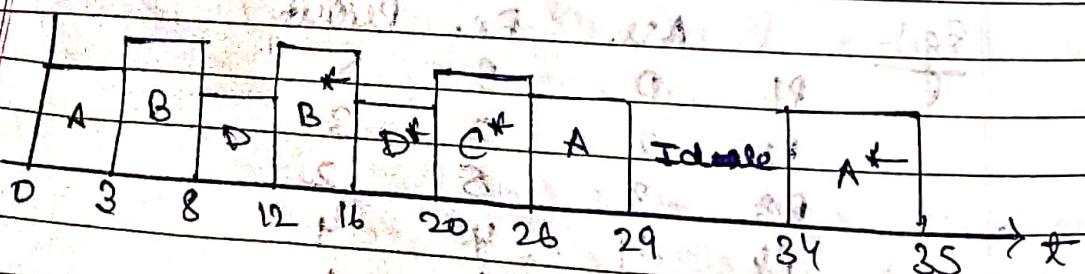
				Periodicity
A	0	7	$= 3C + 5I + 3C + 5I + 1C$	3 *
B	2	9	$= 5C + 2I + 4C$	1
C	4	6	$= 6C$	2
D	6	8	$= 4C + 1I + 4C$	0

Exc. time  $\rightarrow$  3C 5I 6C 4C 3C 4C 4C 1C

Process queue  $\rightarrow$  A | B | C | D | X | B | D | A |

Arrived item  $\rightarrow$  0 2 4 6 8 10, 13, 34

- Non Preemptive.

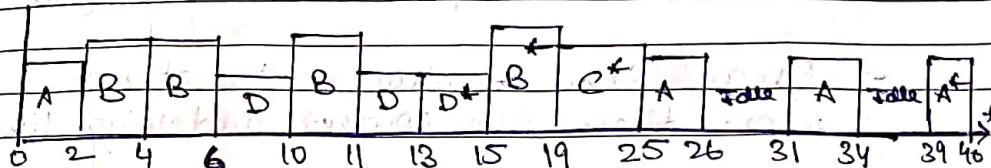


CPU Utilization —  $\frac{30}{35} = \frac{6}{7}$

+ Preemptive

12 3C  
3C 3C 6C 4C 4C 4C 3C 8C  
[A] [B] [C] [D] [E] [F] [G] [H]

0 2 4 6 11 13 31 39



$$\text{CPU Utilization} = \frac{40 - 30}{40} = 0.75$$

- No. of decision pts. in no preemptive is always greater than in non-preemptive.

Priority

- 3  
1  
2  
0

+ Round Robin (RR)

- Usually 10-100 ms.
- No starvation.

- For n processes in ready queue and time quantum q. Mar. response time for any process is then (n-1)q time units.

- At every decision pt (i.e. process completion or time slice expiry or TIO) the scheduler is called to check for a change in process (even if there is only 1 process/job in the queue).

- RR guarantees CPU protection because after every quantum time the CPU calls the scheduler which is a kernel process.

- CPU protection - time - RR uses the time to check for quantum time expiry.

- If quantum time is more, no. of context switch is less so less quantum context switch and more less non-productive work (and vice versa)
- Quant time it should be atleast 1000 times more than the context switching time.

Eg:- Ass. Eqn. SI

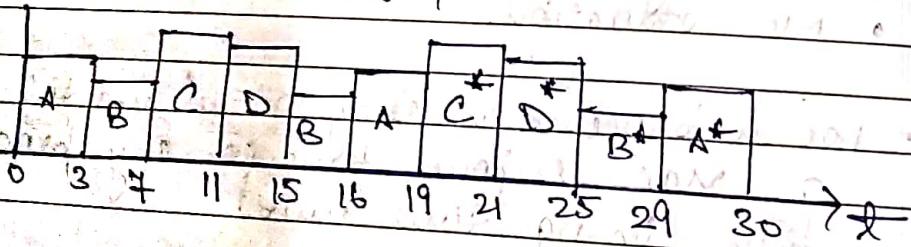
A	0	$T = 3C + 5SI + 3T + 1C$
B	2	$9 = 5C + 2T + 4C$
C	4	$6C$
D	6	$8 = 4C + T + 4C$

Quantum time = 4

~~3C 5C 6C 4C 3C 1C 3C 2C 4C 4C 1C~~

~~|X|B|X|D|X|A|X|C|X|X|X|~~

~~0 2 4 6 8 10 8 11 16 18 24~~



- Preemptive schedule - Not with arrival but with time slice expiry.

- Decision bits are completion/termination, T/W or time slice expiry (no new process arrival).

- Disadvantage - If a process is going out (for

entry to  
the expiry.  
decision pt.  
itself and  
vice versa)  
1000 times  
time.

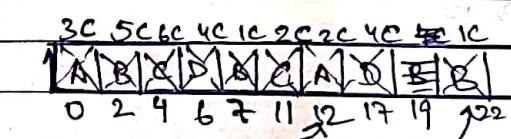
study time  
Page No.: / /  
Date: / /

I/O (or event) before time slice expiry (like A)  
then it loses its remaining time and has to  
join back the queue.

Virtual Round Robin (VRR):

- Auxiliary Queue — To overcome the above problem.
  - when a process (like A) comes back, it joins aux. queue and not main queue with remaining time quantum and not the complete quantum.
  - Aux. queue has higher priority than main queue, so if there is any process in aux. queue then it will be executed first.

Eg: Same Question with aux. queue.



3(4C) 4(3C) - 1(1C)

A	B	X
8	19	29

A	B	C	A	D	B	C	B	A	A	D	B	A
0	3	4	11	12	16	17	19	22	24	28	29	30

- Gives advantage to I/O bound processes or user interactive processes.
- Job/process arrival is not a decision pt.

16/10/19

\* Priority Scheduling with Round Robin :-

- Used in Linux.

- Same priority :- round robin. (Round robin is a preemptive scheduling algorithm)

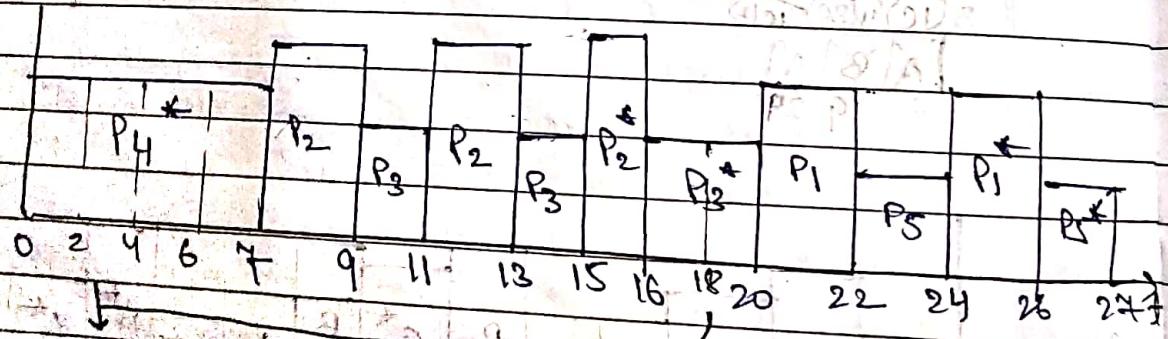
Eg 1	Burst time	Priority
P <sub>1</sub>	4	3
P <sub>2</sub>	5	2
P <sub>3</sub>	8	2
P <sub>4</sub>	7	1
P <sub>5</sub>	3	3

(Quantum time = 2 sec. we will take 2 sec.)

Priority 1 Qu → [P<sub>4</sub>]

Priority 2 → [P<sub>2</sub> P<sub>3</sub> P<sub>2</sub> P<sub>3</sub> P<sub>2</sub>]

Priority 3 → [P<sub>1</sub> P<sub>5</sub> P<sub>1</sub> P<sub>5</sub>]



Time slice exp. are also present here as decision pt.

Eg 2 → P<sub>0</sub> with priority 0 arrives at t = 12.

- If it is preemptive scheduling then arrival is ignored, instead it continues execution.

a decision pt so CPU is given to  $P_1$  because it has highest priority.

In Linux, if a process is preempted in blue (before time slice expiry) then that process (here  $P_1$ ) will join at the front of its priority queue with only the remaining quantum time and not the full time.

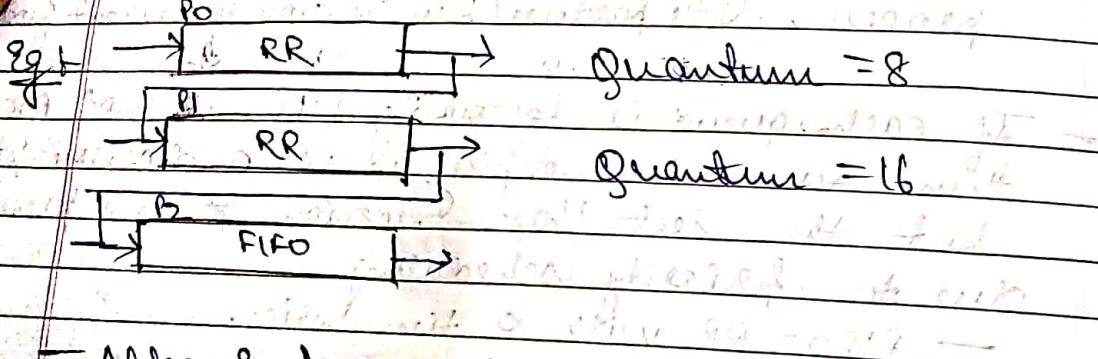
- In this case decision pt is all the 4 cases -
  - Time slice expiry due to RR of each queue.
  - Rest three decision pt. due to priority scheduling scheme (i.e. FIFO, arrival of new process, Idle entry/termination of current process)
- If each queue is following FIFO and not RR, then time slice expiry is not a decision pt., but the rest three decision pt. are present due to priority scheduling.
  - FIFO = RR with  $\infty$  time slice.

### \* Multilevel Queue :-

- More than 1 queue available for diff. types of processes
  - 1 for foreground and other for background.
- Starvation of background process.
  - To overcome this, time sharing b/w the foreground, background processes.
- Sys Process - Scheduler, Priority dispatcher.

## \* Multi Level Feedback Scheduler

- Process can join in any any of the levels.
- Process can be preemted & be demoted to any of levels (i.e. both increasing the priority and decreasing as well)
- Scheduling alg. can be diff. for diff. queue (i.e. one can have RR, others FIFO, etc.)



After 8 time unit, processes demotes from priority 0 to priority 1 queue and then after 16 Quantum time unit it is moved to third priority level.

So for interactive / I/O bound process (they usually have smaller CPU time) they are always in PQ queue because they have CPU have small CPU so they complete execution before 8 units (usually).

### \* Priority Calculation:

$$CPU_j(i) = \frac{CPU_j(i-1)}{2}$$

$$GCPU_k(i) = \frac{GCPU_k(i-1)}{2}$$

$$P_j(i) = Base_j + \frac{CPU_j(i)}{2} + \frac{GCPU_k(i)}{4 \times W_k}$$

$CPU_j$  - Measure of processor utilization by process  $j$  through

$GCPU_j$  - " " " of group  $k$  "

$P_j(i)$  - Priority of process  $j$  at the beg. of interval  $i$  ;  
lower value means higher priority.

$Base_j$  - Base priority of process  $j$ .

$W_k$  - Weighting value assigned to group  $k$ .

- Group Count - CPU may be grouped segregated in various groups

- Any CPU running a particular process increases its CPU process count as well as Group count.

- Group count is common for all CPU of a group

- $0 < W_k \leq 1$  and  $\sum_k W_k = 1$
- Base priority of process remains same.

Eg 1  
+ Process A - Group 1 -  $w_1 = 0.5 ; w_2 = 0.5$   
Process B - Group 2  
" C - Group 2

- If B is executing then group count of C also increases and vice versa.

time	PA	PB	PC
Base CPU Alloc	Base CPU GCPU	Base CPU GCPU	Base CPU GCPU
60 0 0	60 0 0	60 0 0	60 0 0
60 0 0	60 0 0	60 0 0	60 0 0
90 30 30	60 0 0	60 0 0	60 0 0
30 30	60 60	60 60	0 60
74 30 15	90 30 30	75 0 0	30

#### + Fair Share Scheduling :-

- If a group has more processes than some other group, then first group has lower priority than the second group.

18/10/19  
Eg 1 PA - Group 1 -  $w_1 = 0.8 , w_2 = 0.2$  \*  
PB - " 2  
PC - " 3

$$Base_1 = 30$$

$$Base_2 = 40$$

$$Base_3 = 50$$

Time Quantum A = 50

" " B = 40

" " C = 30

0.5;  $w_2 = 0.5$

if count of C  
versa. 0

PA			PB			PC		
Base	CPU	GCPU	Base	CPU	GCPU	Base	CPU	GCPU
30	0	0	40	0	0	50	0	0
✓	1	1	✓	1	1	✓	1	1
50	50	50	0	0	0	0	0	0
49	25	25	40	0	0	50	0	0
✓	1	1	✓	1	1	✓	1	1
25	25	25	40	40	40	0	40	40
39	12	12	45	20	20	45	0	20
✓	1	1	✓	1	1	✓	1	1
62	62	62	20	20	20	0	20	20
54	31	31	57	10	10	62	0	10
✓								

### \* Traditional UNIX Scheduling :-

$$CPU_{j|i} = \frac{CPU_{j|i-1}}{2}$$

$$P_j(i) = Base_j + \frac{couli}{2} + nice_j$$

- Nice value = -20 to 20
  - lower being higher priority.

### \* Linux Scheduling :-

- Interactive Performance - Process with higher I/O should get more chance.
- Good wakeup Perf - New process should be admitted with not much of a delay.

- Fairness - No starvation - Auv. queue concept (Pending queue).
- Exception - Real time process
  - All real time processes are critical processes so they will not wait for fairness and execute as soon as possible.
  - Conventional / user created process follow fairness.
- Real Time Process - Priority - 0 - 99 (100)
- Conventional / user - " " - 100 - 139 (100).

### \* SMP - Symmetric Multiple Processor.

i) Feq. Sym - Same operating freq. of diff. processes.

Eg - Intel atom.

ii) ISA Sym - Same ISA of diff. processor

- A job can be transferred from one core / CPU to other core / CPU if they support same ISA.

• SMP efficiency - If any core / CPU is idle then processes / load can be shifted to that core / CPU from a heavily loaded core / CPU.

• SMP affinity - Processes can set what all cores it can execute - Hard affine.

Aux. queue concept  
(Pending queue).  
process

are critical  
not wait for  
soon as possible.

processes follow

0 - 99 (low)  
100 - 139 (high)

process:

freq. of diff.

diff. processes

in one core/  
multiple use

as well then  
it should  
from a

what all  
- Hard  
affine.

- Soft affine - Process can shift from one core to other.

- most processes are soft affine.

#### \* System Calls -

- nice() - Check nice value.

- renice() - Change nice value.

- sched\_yield() - leave the CPU by itself.

- Nobody is forcing that process.

- Join the ready queue.

Eg:- Producer-consumer problem.

#### \* Linux Schedulers -

- SCHED\_RR → Real-time

SCHED\_FIFO

SCHED\_OTHERS → Non-real/Conventional

\* SCHED\_FIFO - When a process moves to ready queue it will join at the start of the queue.

- So if there are no real-time process with higher priority than the current process is available, then the same process will continue until it finishes (or higher priority job comes). These will be decision pt. but same process will continue in CPU.

- If a higher priority process (Real-time) joins then CPU is given to this new process.
- Joining of non-real time processes does not matter because they are always lower priority than real time process.
- Line 2.4 - Non-real time
- Line 2.6 - Real-time
  - waiting time in admitting new process is decreased
- \* SCHED-RR - when a process moves to ready queue then it will join at the back of the queue.
  - Within same priority number/level, some processes may be FFO and others may be RR.
  - The only diff is the place where they join in ready queue, other things remain same.

### 2110119 Conventional Scheduling - Preemptive Priority

- Two diff types - i) Static Priority  
ii) Dynamic
- Short time sched. is based on dynamic priority
- Dynamic priority is based on static priority
  - Static
  - Nice value is entered by user -  $too = 128$   
 $-20 \text{ to } +19$

(Real-time) join new process!

new process always have lower priority.

in admiring a decision

wishes to ready  
join at the

level, some  
process may be RR,  
here they join  
is remain same.

negative priority

negative priority

on dynamic

static priority  
nice value

use - 100 to 139  
-20 to +19.

- Static priority is always b/w 100 to 139 irrespective of whether it is real time or static concurrent process because it is calculated from nice value.

$$\boxed{\text{Static Pr. No.} = 120 + \text{Nice value}} \\ (100 - 139) \quad (-20 + 19)$$

- Quantum time of a process is based on the static priority no. of that process - higher priority (or less no.) will get more quantum time.

- A new process will inherit the nice value/static priority from parent which can be changed later using nice() sys. call.

$$\rightarrow \text{Static. Pr. No.} = \underbrace{\text{MAX_RT_PRIO}}_{100} + (\text{Nice}) + 20$$

- -ve nice value are for system processes or user inf. processes.

- The nice value are usually user level process.

- For priority 120 and above -

$$\boxed{\text{Quant. time} = (140 - \text{St. Pr.}) * 5 \text{ ms.}}$$

- For priority below 120 -

$$\boxed{\text{Quant. time} = (140 - \text{St. Pr.}) * 20 \text{ ms.}}$$



- Dyn. P.E. should always lie in the range of 100 to 139, i.e. -

$$\text{Dyn. P.E.} = \max(100, \min(\text{st. P.E.} - \text{bonus} + 5, 139))$$

- when the process starts for first time, A bonus is taken as a default value which usually is 0.
- lower dynamic priority no. has higher priority.

- O.S. Bonus  $\leq 10$ .
- If bonus is less than 5 then dyn. priority  $\frac{\text{no.}}{\text{reduced}}$  moves from static priority no. and vice versa
- Avg. sleep time reduces when process is running in CPU and increases when process is in I/O.
- Avg. sleep time does not change when process is in ready queue.
- Interactive process may get some advantage but ~~batch~~ process does not get any advantage.

### \* Interactive Process -

$$D.P. \leq \frac{3}{4} S.P. + 28$$

$$\Rightarrow S.P. - \text{bonus} + 5 \leq \frac{3}{4} S.P. + 28$$

$$\Rightarrow \boxed{\text{Bonus} - 5 \leq \frac{S.P.}{4} - \frac{28}{3}}$$

- Since avg. sleep-time changes after each CPU execution  
so bonus value, dynamic priority and thus queuing time may change for next execution.

Page No.:  
Date: / /

- $\frac{S.P. - 28}{4} = \text{Interactive delta}$

$$\Rightarrow -3 \leq \Delta \leq 6.$$

- A process is called interactive process based on both its avg. sleep time (Bonus value) and its static priority no.

- A process with ~~SP~~ lesser static priority can be int. process with less bonus value.

Eg: • SP. = 100.

$$\text{Bonus} - 5\Delta = \frac{100 - 28}{4}$$

$$\text{Bonus} - 5\Delta = -3 \Rightarrow \text{Bonus } \Delta = 2$$

- ! Avg sleep time should be greater than 200ms.

- SP = 120.

$$\text{Bonus} - 5\Delta = \frac{120 - 28}{4}$$

$$\text{Bonus} - 5\Delta = 30 - 28$$

$$\text{Bonus } \Delta = 2$$

- ! Avg. sleep time  $\Delta = 200 \text{ ms}$

21/10/19  
N  
go

- SP = 139.

$$\text{Bonus} - 5\Delta = \frac{139 - 28}{4}$$

$$\text{Bonus} - 5\Delta = 34.75 - 28$$

$$\text{Bonus } \Delta = 6.75 + 5 = 11.75$$

- ! Not possible.

- ! Process with SP = 139 can never become int.

each CPU execution  
and thus ~~growing~~ <sup>growing</sup>

Page No.:  
Date: / /

Study time  
Page No.:  
Date: / /

## Interactive Process

To avoid starvation in Linux, there are two queue - Active Queue & Expired Queue.

- Any process after completing it's allocated quantum time will go to expired queue at the same priority level acc. to it's next priority value (Because dyn. pr. may/may not change).
- The process with highest priority will be selected from active queue to execute next (Not from expired).
- So a process will get only one chance to execute in 1 cycle.
- Exception - If a process is interactive process then it may join the active queue back instead of expired queue.

No real process will ever expire i.e. it will never go to expired queue and will always finish from active queue. None real process will join active queue.

129	131
99	
0	0

Active      Expired

- After the active queue is becomes empty the expired queue will become active queue and active queue becomes expired queue.
- Only the pointer needs to be exchanged.

- Interactive process will not join the active queue again if -
  - The oldest expired process has already waited for long time.
  - An expired process has higher static priority (lower no) than the interactive process.
- Other than these two conditions the int. process will join back the active queue always
- Every batch process (CPU bound usually) will join the exp. queue after time slice expiry.
- Any new process or I/O/suspended process will always join the active queue.
- Only time slice expired process will go to expired.

### \* Real time - process scheduling - Preemptive priority

- Time quantum is always decided from static priority for both types of processes.
- Real time process dynamic priority is not calculated but static priority directly given (Parameters based on SCHED PARMS).
- Dynamic priority of real time priority does not change.
- If during a running process a new task process with higher priority comes then the current

process has to complete and give CPU to new process. After the execution of new process, the previous process will continue execution only for the remaining time quantum.

All the processes at priority level finishes execution a higher priority lower priority (higher) priority.

• Generation quantum for short-time processes also but there is no mechanism to stop that in real-time process.

• user yield - sys. call to give up CPU by a running process voluntarily.

- The batches will execute with next time after the remaining time quantum only.

• Any process from outside (i.e. new admission or job's arrival/suspended state) will join at the end of the queue.

• The batches which is already in the free CPU (i.e. currently running) will have option to join in next batch or leave it.

• Every iteration of the same priority of processes (inwards batches) happens until Quantum time is elapsed after which the next quantum time is applied.

- Queues - Each queue is a circular linked list.
  - Each link is a dummy node which has the no. of nodes attached to it (No. of process in that priority).
- Bitmap - The lowest priority available level at which a process is available is set 1, other bits are 0.
- OOS scheduler.
- decals\_task\_prio() - Updates the dynm. pr. and any sleep time for the conventional process coming back from I/O.
- Direct Sched. - Current process blocked or finishing execution then scheduler is called.
  - Non preemptive decision pts. 6 are direct calling to sched.
- Lazy Invocation - By setting the ~~THE NEED~~ ~~TA NEED~~ REQUEST flag to 1.
  - When a process currently executing moves from kernel space to user space.
  - Preemptive dec. pt. all lazy invocation.  
(Time slice expiry: done by Scheduler\_tick()  
function, higher priority process came at ready queue: checked done by fly-to-wake-up() function).

study time  
Page No. \_\_\_\_\_  
Date: / /

st circular

which has  
(for no. of

mails level  
be in set 1,

ma. for and  
tional process

cked or  
e is called.  
e direct

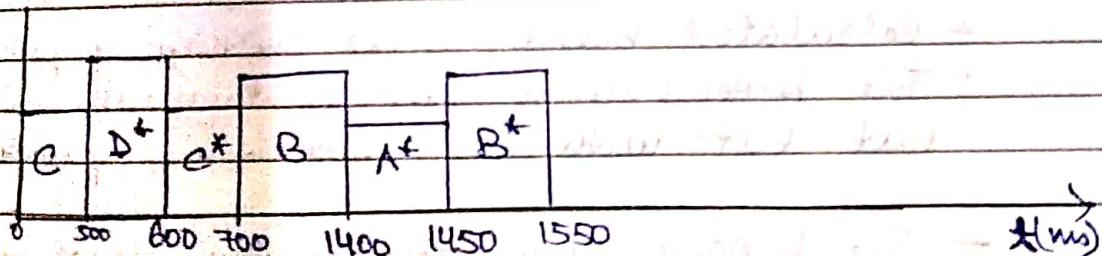
IF NEED  
to flag to 1.  
mutes from

action  
- tick ()  
come at  
try to make

- When a `setScheduler()` system call is issued - A process may get higher priority than the currently running process.

Egt Process	Nice	Ex. time (ms)	Sleep time (ms)	Realtime Priority	Sched Policy
A	8	50	0	-	SCHED_OTHER
B	-15	800	0	-	SCHED_OTHER
C	-5	600	-	30	SCHED_RR
D	0	100	-	20	SCHED_FIFO

Process	S.P.	Q.t (ms)	D.P.
A	128	60	133
B	105	700	110
C	115	2500	30
D	120	100	30



- Assuming A, B, C, D come in ready queue in that order.

B moves to exp.

queue - It can not

become int. process

because curr. sleep time

is 0 and CPU is in

taking place so avg.

sleep time remains 0.

Now the active queue

becomes empty so

expired queue becomes

active and B executes

because it is the only available process there.

22/10/19

## ★ Completely Fair Scheduler:

- Completely Fair Share - for CPU sched.
- Completely Fair Queue - for I/O (disk) sched.

- Red Black tree is used instead of queue.
- Picks the process that has used the least amount of time.
- If process which spends a lot of its time sleeping then its spent time value will be low as it will be selected for execution with low CPU time.
- Order of complexity is  $\log(N)$  for insertion/deletion etc.
- Min allocated time = 1 ms.

- Virtual Runtime - Time provided to a given task.
- Calculated based on the history of process.
- The process with least runtime will be executed next but with the least time quantum.
- The process with least runtime is at the left of tree.
- Complexity of picking the next process to run is O(1) because while adding/deleting a node the left most node will be found in the cache.
- Self Balancing - No path will be twice as long as any other node so max complexity can be  $2\log(N)$ .

- Alloc. arg. is used to initialise a new task.
- Grp. sched. is used to decide virtual run time.

Time slice -

$$\text{timeslice} = \text{sched. period} \times \left( \frac{\text{task's weight}}{\text{total weight of tasks in sum queue}} \right)$$

- Task's wt depends on the task's nice value.
- Min wt = 15 corresponding to nice value 19.
- Max. wt = 88761 corresponding to nice value -20.

- Weight is the granularity of that nice value.

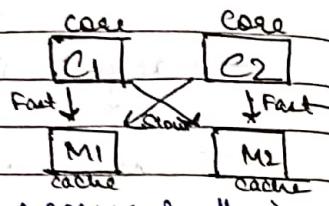
- Time slice becomes smaller as the load increases.

- Granularity of a particular nice value is not a fixed value but a range of value which out of which 1 is given to a process dep. on the type of process.

### \* Multi Processor Scheduling -

- Single Cycle data path - IF, ID, EX, MEMR, MEMW will all happen in 1 clock cycle.
- Cycle width is big.
- Multi Cycle data path - Each process will happen in 1 cycle each.
- Cycle width is small.
- Buff

- Each cycle uses the OFR of previous cycle so buffer latch is needed to store the OFR.
- Pipelining - Multiple int. can use execute parallel at diff. stage (one at IF stage, other at ID, other at EX and so on) at the same time.
- Intel i4 20 stage pipelined.
- Super scalar - Multiple pipeline for ..
- NUMA - Non uniform memory access.
- Symmetric Multiprocessing (SMP) - Multi core



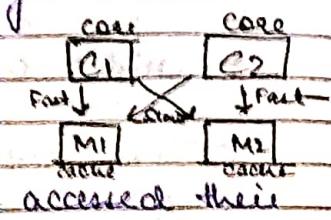
- Two mechanism - i) 1 Global queue for all the cores ii) Diff. queues for diff. cores.

- Global Queue - Not possible.
  - If a process running in Core 1 goes to T0 and comes back, it should come back at the same place because all cache data are stored in cache of core 1. If it comes from core 2 then data has to be moved from cache 1 to cache 2.

off of previous cycle  
at to above that why

can see execute  
stage here at if  
bus at EX and so

line from  
y access.



) - Multicore

MP.

queue for all  
diff. cores.

ne I goes to I/O  
come back at  
cache stage

If it a comes  
be moved from

unless used local queue for core.

Page No. 1 Date 1/1/2024

- cache to cache

- local Queue

- If multiple core reg. a job at the same time  
then it may happen that multiple cores  
may pick the same job. To avoid prevent this  
if one core is accessing the queue to pick up the  
a process then the queue has to locked so that  
other cores cannot use it at that time.

- local Queue - load sharing is reg. explicitly

- In global queue ~~this is not reg. (One action-~~  
~~array).~~

\* Memory Stall Cycle - If a process needs a data  
which is not in its cache

then it has to wait in CPU till data comes  
from L2 cache. If again L2 cache miss then  
it has to fetch from primary memory and  
process has to wait in CPU.

- Till that time CPU will execute "nop" so  
no productive work is done that time.  
(Reg-write = 0, Mem-write = 0).

- So getting data from I/O is better because that  
process gets blocked and other process can  
execute in other CPU.

\* Multi Threading - Each core has  $\geq 1$  hardware thread.

- Hyperthreading in intel.

- Diff. from software threads / os threads.

- If in a core 1 thread needs data from mem.  
then core switches to other ~~core~~ thread for  
execution until the data is fetched for 1st thread

Effects

- The Great Famine was one of the greatest ones in India & there is much evidence towards the same because all the work is in the same area of that year - the famine greater than any other because it was more severe and longer.
- Peasants suffered the most unanticipated by the government with freedom from serfdom while landlords were to benefit in a great way.
- Effect on migration not allowed, certain only 1 case.
- 1st effect - migration allowed.
  - Land revenue in others for the sake of land balancing is suspended.

→ going of a peasant can be set by road.

→ Fig 1 of Peas during compilation on the same case / pp.

→ Peasant in Turkey

- Any little case across the entire of Turkey but the time taken to do so is lesser than other case's case will be more than half even though it's case.

→ In a case where this disease is present in the cattle kept present in some case this

Study time  
Page No. 1  
Date 1 / 1

23/10/19

- In local queue i.e. core if a process goes to TLB, then it will come back to the same core because all its data is in the cache of that core. So cache flush does not happen.

\* Processor Affinity - Bitmaps maintained by each process to decide which core/process can & cannot execute in.

\* Hard affinity - Migration not allowed.  
- Executing only 1 core.

\* Soft affinity - Migration allowed.  
- Can execute in other cores also if load balancing is required.

Affinity of a process can be set by user.

→ Eg. Pass 1 & Pass 2: during compilation in the same core / MP.

\* NUMA: Quick path in Intel.

- Any core can access the cache of any core but the time taken for a core to access other core's cache will be more than if it access its own cache.

- So if a core needs data which is not in its cache but present in some other core's cache.

cache them without a  
first miss

\* NOR MA-

- Multicore
- Distributed

\* Load Mi-

- Pull

• Push

less b

• Octo

- JP  
- JI

• In

- 2

- P

if a process goes  
value back to the  
its data is in the  
cache flush does

is maintained by  
cores to decrease  
it executes in.  
not allowed  
by 1 core.

in other core  
balancing is

compilation in the

the cache of  
time taken for  
self's cache will

not  
is present  
some one else's

cache them if that core can take that data directly  
without copying that data to its own cache  
first which speeds up performance.

+ NORMA — Non remote memory access.  
— Core cannot access cache of other  
core.

- Multicore — Shared memory.
- Distributed sys. — No shared memory.

\* Load Migration —

- Pull Migration — If a core is free, it can take  
the job of any other busy core  
to execute.
- Push Migration — If a core is very busy, it  
can push its job to other cores  
less busy cores of the system.
- QoS Core — Two groups of 4 core.
  - Try balancing within the group first (the 4 cores of that group).
  - If groups are unbalanced then transfer job from a group to other.
- In Linux, push migration algo. will run every 200 ms.
  - Pull mig. will happen automatically when a core becomes free/idle.

Operating Systems

### \* Windows

- Push migration algo uses past history of the core to determine which is busy and which is free/ idle.

### \* Windows CPU Sched :-

- Priority: 0 - Mem. Management

1-15 - Variables

16-31 - Real time.

- one queue for each priority level

- Higher no. = higher priority.

- High priority to int. process - mouse/window

class

• Variables :- HIGH-PRIORITY-CLASS

ABOVE-NORMAL

NORMAL

BELLOW-NORMAL

LOWEST

— IDLE

- The Three classes present in both Real-time Range and Variable Range

- From windows 7 - Two levels of thread mgt

- User level thread scheduling

- Kernel level thread scheduling