Vivek Kumar Computer Science and Engineering IIIT Delhi vivekk@iiitd.ac.in



Last Lecture



void s wł	<pre>scheduler() { nile(true) {</pre>
	<pre>lock(process_table);</pre>
	<pre>foreach(Process p: scheduling_algorithm(process_table)) {</pre>
	if(p->state != READY) {
	continue;
	}
	p->state = RUNNING;
	<pre>unlock(process_table);</pre>
	<pre>swtch(scheduler_process, p);</pre>
	// p is done for now
	<pre>lock(process_table);</pre>
	}
	<pre>unlock(process_table);</pre>
}	
}	

Property	FCFS	SJF
Optimize Average Completion Time	×	
Prevent Starvation		*
Prevent Convoy Effect	*	*
Psychic Skills NOT Needed	\checkmark	*





Today's Class

• Process scheduling policies

- Shortest Time to Completion First
- o Round Robin
- o Multi-level feedback scheduling
- o Linux Completely Fair Scheduler (CFS)

Acknowledgements: Lecture content derived from the CS162 course offered at University of California, Berkeley

CSE231: Operating Systems

Shortest Time to Completion First (STCF) Introduce the notion of preemption

A running task can be de-scheduled before completion.

STCF

Schedule the task with the least amount of time left



CSE231: Operating Systems

	Process	Burst Time (left)	<u>Arrival Time</u>
	<i>P</i> ₁	3	10
Schedule the task with the	<i>P</i> ₂	6	1
least amount	P ₃	24	0
of time left	P_4	16	18



CSE231: Operating Systems

	<i>P</i> ₁	3	10
Schedule the task with the	P ₂	6	1
least amount	P ₃	23	0
of time left	P ₄	16	18





	FIULESS	buist fille (left)	Arrivari
	P_1	3	10
Schedule the task with the	P_2	0	1
least amount	<i>P</i> ₃	23	0
of time left	P_4	16	20





	<u>Process</u>	Burst Time (left)	<u>Arrival Time</u>
	P ₁	3	10
Schedule the task with the	P ₂	0	1
least amount	P ₃	20	0
of time left	P_4	16	18

	Р3	Р2	Р3	P1
0	1		7 1	0 13



<i>P</i> ₁	0	10
<i>P</i> ₂	0	1
P ₃	15	0
P ₄	16	18

Р3	P2	Р3	P1		Р3
0 1	1 7	7	10	13	18

CSE231: Operating Systems

<i>P</i> ₁	0	10
<i>P</i> ₂	0	1
<i>P</i> ₃	0	0
P_4	15	18

	Р3	P2	Р3	P1	РЗ
0	1		7 10) 13	33



Schedule the
task with the
least amount
of time left

P_1	0	10
<i>P</i> ₂	0	1
<i>P</i> ₃	0	0
P_4	15	18

F	P3	P2	Р3	P1	Р3	P4
0	1	7	7 10) 13	33	2



Are we done?

Can STCF lead to starvation?

Yes

Any scheduling policy that always favours a fixed property for scheduling leads starvation

No change!



Are we done?

Is STCF subject to the convoy effect?

No!

STCF is a preemptible policy



Today's Class

• Process scheduling policies

- o Shortest Time to Completion First
- o Round Robin
- o Multi-level feedback scheduling
- o Linux Completely Fair Scheduler (CFS)

Acknowledgements: Lecture content derived from the CS162 course offered at University of California, Berkeley

CSE231: Operating Systems

Round-Robin Scheduling

RR runs a job for a **time slice** (a **scheduling quantum**)

Once time slice over, Switch to next job in ready queue. => Called time-slicing



RR with Time Quantum = 20 $\frac{Process}{P_1} = \frac{Burst Time}{53}$ $\frac{P_2}{P_3} = \frac{68}{24}$



RR with Time Quantum = 20 $\frac{Process}{P_1} \begin{array}{r} Burst Time \\ 53 => 33 \\ P_2 \\ P_3 \\ P_4 \end{array} \begin{array}{r} 8 \\ 68 \\ 24 \end{array}$



RR with Time Quantum = 20 $\frac{Process}{P_1}$ $\frac{Burst Time}{33}$ P_2 P_3 P_4 24













RR with Time Quantum = 20 $\frac{Process}{P_1} \begin{array}{c} Burst Time \\ 33 => 13 \\ P_2 \\ P_3 \\ P_4 \end{array} \begin{array}{c} 0 \\ 48 \\ 4 \end{array}$











RR with Time Quantum = 20

Waiting time

- P₁= 0 + (68-20)+(112-88)=72
- P₂=(20-0)=20
- P₃=(28-0)+(88-48)+(125-108)+0=85
- P₄=(48-0)+(108-68)=88



$$\frac{72+20+85+88}{4} = 66.25$$

Average completion time

 $(\frac{125+28+153+112}{_4}=104.25)$



CSE231: Operating Systems

RR Quantum

- •Assume there is no context switching overhead
- What happens to average completion time when we *decrease Q*?



CS6456 University of Virginia

Switching is not free!

Small scheduling quantas lead to frequent context switches

- Context switch overhead
 - Trash cache-state

q must be large with respect to context switch, otherwise overhead is too high

Are we done?

Can RR lead to starvation?

No



Are we done?

Can RR suffer from convoy effect?

No

Only run a time-slice at a time



Summary

Property	FCFS	SJF	STCF	RR	High in case of
Optimize Average Completion Time	*			* -	long running processes
Optimize Average Waiting Time	*	*	*		Doesn't favor any fixed property
Prevent Starvation		*	*	\sim	
Prevent Convoy Effect	*	*		V -	Due to time slicing
Psychic Skills NOT Needed		*	*	\sim	



Today's Class

• Process scheduling policies

- o Shortest Time to Completion First
- o Round Robin
- Multi-level feedback scheduling
- o Linux Completely Fair Scheduler (CFS)

Acknowledgements: Lecture content derived from the CS162 course offered at University of California, Berkeley

CSE231: Operating Systems

What we Want?



CSE231: Operating Systems

Approximation of CPU Bursts using Priority Scheduling



- Priority scheduling always run the ready process (or job) with highest priority
- Systems may try to set priorities according to some policy goal
 Example: Give IO/interactive jobs higher priority than long calculation
- How to achieve fairness and avoid starvation?
 Elevate priority of threads that don't get CPU time

Multi-level Feedback Queue

[High Priority] Q8 \rightarrow A \rightarrow B Q7 Q6 Q5 Q4 \rightarrow C

Q2 [Low Priority] Q1 \longrightarrow D

Q3

Figure 8.1: MLFQ Example

- Create distinct queues for ready processes, each assigned a different priority level
- All processes belong to any one queue at any given time
- Processes can move between queues
- Use priorities to decide from which queue process should be picked next
- Individual queues run RR with increasing time slice

MLFQ (Naïve Version)

If Priority(A) > Priority(B) (different queues) A runs (B doesn't).

Rule 2

If Priority(A) = Priority(B), A & B run in RR.

Key question:

How do you set the priorities?

Vary the priority of a job based on its observed behavior



Learning behavior

Rule 3

When a job enters the system, it is placed at the highest priority (the topmost queue).

Rule 4a

If a job uses up an entire time slice while running, its priority is *reduced* (i.e., it moves down one queue).



CSE231: Operating Systems

Learning behavior



Long running CPU-bound



Learning behavior



Long running CPU-bound



Learning behavior



Long running CPU-bound



Learning behavior



Long running CPU-bound



Learning behavior

Rule 3

When a job enters the system, it is placed at the highest priority (the topmost queue).

Rule 4a

If a job uses up an entire time slice while running, its priority is *reduced* (i.e., it moves down one queue).

Rule 4b

If a job gives up the CPU before the time slice is up, it stays at the *same* priority level.



CSE231: Operating Systems

Learning behavior

Where do IO-bound/interactive jobs end up?

a) High Priority Queue b) Low Priority Queue

Ideally in high priority, but if it finishes up the time slice then moved to a lower priority queue



Learning behavior



Long running CPU-bound



Schedule



Learning behavior



Long running CPU-bound



Schedule





Learning behavior



Long running CPU-bound

Computes for 1 ms and then IO for 1 ms



Schedule

Learning behavior



Long running CPU-bound

Computes for 1 ms and then IO for 1 ms



Schedule



Learning behavior



Long running CPU-bound

Computes for 1 ms and then IO for 1 ms



Schedule



Learning behavior



Long running CPU-bound

Computes for 1 ms and then IO for 1 ms



Schedule



Learning behavior



Long running CPU-bound

Computes for 1 ms and then IO for 1 ms



Schedule



Learning behavior



Long running CPU-bound

Computes for 1 ms and then IO for 1 ms



Schedule



Learning behavior



Long running CPU-bound

Computes for 1 ms and then IO for 1 ms



Schedule



MLFQ (Naïve Version): Pseudocode

foreach(queue : PrioritizedReadyQueues) { // RR foreach(job : queue) { // RR int timeslice = queue.getSlice(); // Execute job for timeslice associated with this queue // If job used full timeslice then push job in // ready queue with one lower priority level



Are we done?

Can MLFQ can be gamed?

Intentionally insert IO request just before time quanta to stay on queue.

MLFQ is subject to starvation?

Yes!



CSE231: Operating Systems

Are we done?

MLQF can be gamed:

Intentionally insert IO request just before time quanta to stay on queue.

Rule 4

Once a job uses up its time allotment at a given levels (regardless of how many times gave up CPU), reduce priority

MLQF is subject to starvation:

Systematically prioritize higherpriority queues Rule 5 After some time period S, move all jobs in system to the topmost queue.



Today's Class

• Process scheduling policies

- o Shortest Time to Completion First
- o Round Robin
- o Multi-level feedback scheduling
- Linux Completely Fair Scheduler (CFS)

Acknowledgements: Lecture content derived from the CS162 course offered at University of California, Berkeley

CSE231: Operating Systems

Linux Completely Fair Scheduler

- Goal: Each process gets an equal share of CPU
- Attempts to divide the CPU time fairly (equally) among all the processes
 - If the latency to resume a scheduler is M units of time, N processes will get M/N time slice of CPU





• **Constraint 1**: Scheduling Latency

Target Latency: 20ms, 4 Processes
 Each process gets 5ms time slice

• Target Latency: 20 ms, 200 Processes

- Each process gets 0.1ms time slice
- Recall Round-Robin: Huge context switching overhead

• **Constraint 2**: *Minimum Granularity for scheduling latency*

- Target Latency 20ms, Minimum Granularity 1ms, 100 processes
 - \odot Each process gets 1ms time slice

- Which process to pick from the ready queue?
- •**Constraint 3**: Pick the process from the ready queue with minimum vruntime (helps the IO/interactive jobs)
 - Uses virtual runtime (vruntime) variable in PCB to keep track of time a process has executed on the CPU, and it is updated at every context switch
 - vruntime += t * weight based on process priority

<u>Process</u>	<u>Burst Time</u>
<i>P</i> ₁	5
<i>P</i> ₂	10
<i>P</i> ₃	16
P_4	16

Scheduling latency is 20ms, minimum granularity is 1ms, and weight=1

• Time slice in beginning = 20/4 = 5ms



0

<u>Process</u>	<u>Burst Time</u>	
<i>P</i> ₁	0	
<i>P</i> ₂	10	
<i>P</i> ₃	16	
P_4	16	



- Time slice in beginning = 20/4 = 5ms
- Time slice after P1 terminates = 20/3 ~ 6ms



<u>Process</u>	<u>Burst Time</u>
<i>P</i> ₁	0
<i>P</i> ₂	4
<i>P</i> ₃	16
P_4	16

- Scheduling latency is 20ms, minimum granularity is 1ms, and weight=1
- Time slice in beginning = 20/4 = 5ms
- Time slice after P1 terminates = 20/3 ~ 6ms



<u>Process</u>	<u>Burst Time</u>
<i>P</i> ₁	0
<i>P</i> ₂	4
<i>P</i> ₃	10
P_4	16

- Scheduling latency is 20ms, minimum granularity is 1ms, and weight=1
- Time slice in beginning = 20/4 = 5ms
- Time slice after P1 terminates = 20/3 ~ 6ms

	P ₁	P ₂	P ₃	P ₄
0	5	1	11	17

<u>Process</u>	<u>Burst Time</u>
<i>P</i> ₁	0
<i>P</i> ₂	4
<i>P</i> ₃	10
P_4	10

- Scheduling latency is 20ms, minimum granularity is 1ms, and weight=1
- Time slice in beginning = 20/4 = 5ms
- Time slice after P1 terminates = 20/3 ~ 6ms



<u>Process</u>	<u>Burst Time</u>
<i>P</i> ₁	0
<i>P</i> ₂	0
<i>P</i> ₃	10
P_4	10

- Scheduling latency is 20ms, minimum granularity is 1ms, and weight=1
- Time slice in beginning = 20/4 = 5ms
- Time slice after P1 terminates = 20/3 ~ 6ms
- Time slice after P2 terminates = 20/2 = 10ms



Process	<u>Burst Time</u>		
<i>P</i> ₁	0		
<i>P</i> ₂	0		
<i>P</i> ₃	0		
P_4	10		

- Scheduling latency is 20ms, minimum granularity is 1ms, and weight=1
- Time slice in beginning = 20/4 = 5ms
- Time slice after P1 terminates = 20/3 ~ 6ms
- Time slice after P2 terminates = 20/2 = 10ms

	<i>P</i> ₁	P ₂	P ₃		<i>P</i> ₄	P ₂		<i>P</i> ₃
0	5	1	1	17	2	3	27	

<u>Process</u>	<u>Burst Time</u>		
<i>P</i> ₁	0		
<i>P</i> ₂	0		
<i>P</i> ₃	0		
P_4	0		

- Scheduling latency is 20ms, minimum granularity is 1ms, and weight=1
- Time slice in beginning = 20/4 = 5ms
- Time slice after P1 terminates = 20/3 ~ 6ms
- Time slice after P2 terminates = 20/2 = 10ms

	<i>P</i> ₁	P ₂	P ₃	P_4	P ₂	P ₃	P_4	
0	5	1	11	17 2	.3 2	27 3	37 4	- 47

Next Lecture

• Dynamic memory allocation

