Scheduling

Johan Montelius

KTH

2021

メロト メタト メヨト メヨト 二日

1/46

<ロト < 部 > < 言 > < 言 > 言 の Q () 2/46

We have a set of processes: they all want to execute immediately and they do not want to be interrupted.

We have a set of processes: they all want to execute immediately and they do not want to be interrupted.

Solution:

We have a set of processes: they all want to execute immediately and they do not want to be interrupted.

Solution:

Let's keep some waiting and let's interrupt them.

We have a set of processes: they all want to execute immediately and they do not want to be interrupted.

Solution:

Let's keep some waiting and let's interrupt them.

Question:

We have a set of processes: they all want to execute immediately and they do not want to be interrupted.

Solution:

Let's keep some waiting and let's interrupt them.

Question:

• What metrics are important?

We have a set of processes: they all want to execute immediately and they do not want to be interrupted.

Solution:

Let's keep some waiting and let's interrupt them.

Question:

- What metrics are important?
- Does it matter in what order we schedule processes?

We have a set of processes: they all want to execute immediately and they do not want to be interrupted.

Solution:

Let's keep some waiting and let's interrupt them.

Question:

- What metrics are important?
- Does it matter in what order we schedule processes?
- Are there optimal solutions?

Assume we have a set of jobs.

• Each job takes an equal amount of time.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

Assume we have a set of jobs.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

This is unrealistic - we will relax these requirements.

...every now and then I get a little bit lonely

...every now and then I get a little bit lonely



$T_{\rm turnaround} = T_{\rm completion} - T_{\rm arrival}$

How long time does it take to complete the job?

Assume we have three tasks, all *arrive* at time 0 and take 10 ms to execute.

Assume we have three tasks, all *arrive* at time 0 and take 10 ms to execute.



4 ロ ト 4 伊 ト 4 王 ト 4 王 ト 王 の 9 0 0 6 / 46

Assume we have three tasks, all arrive at time 0 and take 10 ms to execute.



Assume we have three tasks, all *arrive* at time 0 and take 10 ms to execute.



4 ロ ト 4 日 ト 4 王 ト 4 王 ト 王 少 9 0 0
6 / 46

Assume we have three tasks, all *arrive* at time 0 and take 10 ms to execute.



What is the average $T_{turnaround}$?

Assume one task takes 30 ms to execute.

Assume one task takes 30 ms to execute.



<ロト < 団 > < 巨 > < 巨 > < 巨 > 三 の Q (~ 7/46

Assume one task takes 30 ms to execute.



Assume one task takes 30 ms to execute.



Assume one task takes 30 ms to execute.



What is the average $T_{turnaround}$?

Assume one task takes 30 ms to execute.



What is the average $T_{turnaround}$? Can we do better?

Always schedule the shortest job.

Always schedule the shortest job.



Always schedule the shortest job.



Always schedule the shortest job.



・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ うへで
Shortest Job First (SJF)

Always schedule the shortest job.



What is the average $T_{turnaround}$?

Shortest Job First (SJF)

Always schedule the shortest job.



What is the average $T_{turnaround}$? Problem solved!

<ロト < (日) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1) < (1)

Assume we have three tasks, one arrive at time 0 and takes 30 ms to execute. Two arrive at time 10 and take 10 ms each.

J1:

J2:

J3:



Assume we have three tasks, one arrive at time 0 and takes 30 ms to execute. Two arrive at time 10 and take 10 ms each.



Assume we have three tasks, one arrive at time 0 and takes 30 ms to execute. Two arrive at time 10 and take 10 ms each.



Assume we have three tasks, one arrive at time 0 and takes 30 ms to execute. Two arrive at time 10 and take 10 ms each.



Assume we have three tasks, one arrive at time 0 and takes 30 ms to execute. Two arrive at time 10 and take 10 ms each.



We need to preempt the execution of a job.

Shortest Time-to-Completion First (STCF)

Let's always schedule the task that has the shortest time left to completion.

J2:

J1:

J3:



Shortest Time-to-Completion First (STCF)

Let's always schedule the task that has the shortest time left to completion.



Shortest Time-to-Completion First (STCF)

Let's always schedule the task that has the shortest time left to completion.



The policy is also known as Preemptive Shortest Job First (PSJF)

Shortest Time-to-Completion First is an optimal policy.

Shortest Time-to-Completion First is an optimal policy.

The problem is that we do not know the total execution time aforehand.

Shortest Time-to-Completion First is an optimal policy.

The problem is that we do not know the total execution time aforehand.

There might be more important metrics than turnaround time.

Talk about ...



In an interactive environment we might want to minimize response time.

In an interactive environment we might want to minimize response time.

$T_{\rm response} = T_{\rm first \ scheduled} - T_{\rm arrival}$

In an interactive environment we might want to minimize response time.

$T_{\rm response} = T_{\rm first \ scheduled} - T_{\rm arrival}$

The response might not be completed unless the job completes but it's an ok metrics.

Assume we have three jobs that all arrive at time 0 and all take 40 ms to complete.

Assume we have three jobs that all arrive at time 0 and all take 40 ms to complete.



Assume we have three jobs that all arrive at time 0 and all take 40 ms to complete.



Assume we have three jobs that all arrive at time 0 and all take 40 ms to complete.



14 / 46

Assume we have three jobs that all arrive at time 0 and all take 40 ms to complete.



14 / 46

Preempt a job in order to improve response time, give each job a time-slice of 10 ms.



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Preempt a job in order to improve response time, give each job a time-slice of 10 ms.



<ロト < 部 > < 言 > < 言 > 言 の < で 15 / 46











Preempt a job in order to improve response time, give each job a time-slice of 10 ms.



What is the average response time?





^{15 / 46}

You can't



Assume we have two processes, each take 40 ms of CPU time but one will do $I/O\mbox{-}operations$ every 10 ms.

J1:

J2:


















deschedule when initiate I/O

An I/O-operation will take time to complete and we (the CPU) could do some useful work while a process is waiting.

deschedule when initiate I/O

An I/O-operation will take time to complete and we (the CPU) could do some useful work while a process is waiting.



deschedule when initiate I/O

An I/O-operation will take time to complete and we (the CPU) could do some useful work while a process is waiting.



A process is descheduled if it is preempted or if it initiates a I/O-operation.

J1:

J2:

0 10 20 30 40 50 60 70 80 90 100 110 120 ms





・ロト ・母 ・ ・ ヨ ・ ・ ヨ ・ りへぐ



・ロト ・母・ ・ヨ・ ・ヨ・ りへぐ



・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ うへで





・ロト ・四ト ・ヨト ・ヨト ・ヨー うへの





Ideal world:

• Each job takes an equal amount of time.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

Real world:

• Jobs take different amount of time.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

- Jobs take different amount of time.
- Jobs arrive at different time.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

- Jobs take different amount of time.
- Jobs arrive at different time.
- We can preempt job.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

- Jobs take different amount of time.
- Jobs arrive at different time.
- We can preempt job.
- Jobs do use I/O.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

- Jobs take different amount of time.
- Jobs arrive at different time.
- We can preempt job.
- Jobs do use I/O.
- Runt-time is not know.

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

- Jobs take different amount of time.
- Jobs arrive at different time.
- We can preempt job.
- Jobs do use I/O.
- Runt-time is not know.
- What do we do?

- Each job takes an equal amount of time.
- All jobs *arrive* at the same time.
- A job will run to completion.
- The jobs only use the CPU (no I/0 etc).
- The run-time of each job is known.

Real world:

- Jobs take different amount of time.
- Jobs arrive at different time.
- We can preempt job.
- Jobs do use I/O.
- Runt-time is not know.
- What do we do?

Can we design scheduling policies that give us good turn-around time and short response time?

Multi-level Feedback Queue (MLFQ)

Goals:

Multi-level Feedback Queue (MLFQ)

Goals:

• Good turnaround time - scheduled jobs so that jobs with short time to completion are not delayed too much.
Goals:

- Good turnaround time scheduled jobs so that jobs with short time to completion are not delayed too much.
- Improve responsiveness of interactive jobs schedule *interactive processes* more often.

Idea:

Goals:

- Good turnaround time scheduled jobs so that jobs with short time to completion are not delayed too much.
- Improve responsiveness of interactive jobs schedule *interactive processes* more often.

Idea:

• Multiple levels of priority - interactive jobs have higher priority.

Goals:

- Good turnaround time scheduled jobs so that jobs with short time to completion are not delayed too much.
- Improve responsiveness of interactive jobs schedule *interactive processes* more often.

Idea:

- Multiple levels of priority interactive jobs have higher priority.
- Each level uses round-robin to give processes an equal share.

Goals:

- Good turnaround time scheduled jobs so that jobs with short time to completion are not delayed too much.
- Improve responsiveness of interactive jobs schedule *interactive processes* more often.

Idea:

- Multiple levels of priority interactive jobs have higher priority.
- Each level uses round-robin to give processes an equal share.
- Processes can be moved to a higher or lower level depending on their behavior.

Goals:

- Good turnaround time scheduled jobs so that jobs with short time to completion are not delayed too much.
- Improve responsiveness of interactive jobs schedule *interactive processes* more often.

Idea:

- Multiple levels of priority interactive jobs have higher priority.
- Each level uses round-robin to give processes an equal share.
- Processes can be moved to a higher or lower level depending on their behavior.

How do we identify interactive processes and how do we make sure that they have high priority?

Rules of the game: MLFQ

<ロト <部ト < 言ト < 言ト < 言 > 3 () 22/46

• Rule 1: if Priority(A) > Priority(B) then A is scheduled for execution.

- Rule 1: if Priority(A) > Priority(B) then A is scheduled for execution.
- Rule 2: if Priority(A) = Priority(B) then A and B are scheduled in round-robin.

- Rule 1: if Priority(A) > Priority(B) then A is scheduled for execution.
- Rule 2: if Priority(A) = Priority(B) then A and B are scheduled in round-robin.
- Rule 3: when a new job is created it starts with the highest priority.

- Rule 1: if Priority(A) > Priority(B) then A is scheduled for execution.
- Rule 2: if Priority(A) = Priority(B) then A and B are scheduled in round-robin.
- Rule 3: when a new job is created it starts with the highest priority.

- Rule 1: if Priority(A) > Priority(B) then A is scheduled for execution.
- Rule 2: if Priority(A) = Priority(B) then A and B are scheduled in round-robin.
- Rule 3: when a new job is created it starts with the highest priority.

Change priority (let's try this)

• Rule 4a: a job that has to be preempted (time-slice consumed) is moved to a lower priority.

- Rule 1: if Priority(A) > Priority(B) then A is scheduled for execution.
- Rule 2: if Priority(A) = Priority(B) then A and B are scheduled in round-robin.
- Rule 3: when a new job is created it starts with the highest priority.

Change priority (let's try this)

- Rule 4a: a job that has to be preempted (time-slice consumed) is moved to a lower priority.
- Rule 4b: a job that initiates a I/O-operation (or yields) remains on the same level.

Q2:

Q1:

Q0:

0 10 20 30 40 50 60 70 80 90 100 110 120 ms

・ロト ・西ト ・ヨト ・ヨー うへの

Q2: Q1: Q0: ms









・ロト・日本・モート ヨー ちゃぐ



<ロト < 合 ト < 言 ト < 言 ト こ の へ () 23 / 46



・ロト ・日・ ・ヨ・ ・ヨ・ うへぐ



・ロト ・雪 ・ ・ ヨ ・ ・ ヨ ・ うへぐ





・ロト ・ 酉 ト ・ ヨ ト ・ ヨ ・ のへぐ



・ロト ・日・ ・ヨ・ ・ヨ・ うへぐ



・ロト・日本・ キョナ・ キョー うくぐ



・ロト・(聞)・(目)・(目)・ 目・ のへの



・ロト・(聞)・ (目)・ (目)・ (日)・ (の)の



・ロト・(聞)・ (目)・ (目)・ (日)・ (の)の



・ロト ・ 理 ・ ・ ヨ ・ ・ ヨ ・ うへぐ



・ロト・(型ト・(ヨト・ヨト) ヨー のへぐ



・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ うへで



<ロト < 部 ト < 言 ト < 言 ト 言 の Q (や 24 / 46



イロト イポト イヨト イヨト
















24 / 46





24 / 46













24 / 46



24 / 46

• Rule 1: if Priority(A) > Priority(B) then A is scheduled for execution.

- Rule 1: if Priority(A) > Priority(B) then A is scheduled for execution.
- Rule 2: if Priority(A) = Priority(B) then A and B are scheduled in round-robin.

- Rule 1: if Priority(A) > Priority(B) then A is scheduled for execution.
- Rule 2: if Priority(A) = Priority(B) then A and B are scheduled in round-robin.
- Rule 3: when a new job is created it starts with the highest priority.

- Rule 1: if Priority(A) > Priority(B) then A is scheduled for execution.
- Rule 2: if Priority(A) = Priority(B) then A and B are scheduled in round-robin.
- Rule 3: when a new job is created it starts with the highest priority.
- Rule 4a: a job that has to be preempted (time-slice consumed) is moved to a lower priority.

- Rule 1: if Priority(A) > Priority(B) then A is scheduled for execution.
- Rule 2: if Priority(A) = Priority(B) then A and B are scheduled in round-robin.
- Rule 3: when a new job is created it starts with the highest priority.
- Rule 4a: a job that has to be preempted (time-slice consumed) is moved to a lower priority.
- Rule 4b: a job that initiates a I/O-operation (or yields) remains on the same level.
- Rule 5: after some time, move a job to the highest priority.

A job is given a *allotted time*, to consume at each priority level.

A job is given a *allotted time*, to consume at each priority level.

• Rule 4: a job that has consumed its allotted time is moved to a lower priority.

A job is given a *allotted time*, to consume at each priority level.

- Rule 4: a job that has consumed its allotted time is moved to a lower priority.
- Rule 5: after some time, move all jobs to the highest priority.

Setting the parameters:

Setting the parameters:

- How long is a time slice?
- How many queues should there be?
- How long time should a allotted time be in a specified queue?
- How often should a job be boosted to the highest priority?

+ □ ト < @ ト < E ト < E ト E のQC</p>

• we stop focusing on *turnaround time* and *reaction* and

- we stop focusing on turnaround time and reaction and
- start treating every job in a fair manner.

- we stop focusing on turnaround time and reaction and
- start treating every job in a fair manner.

Give each job fair share.

Justice for all



Let's have a lottery:

Proportional share

Let's have a lottery:



We divide the tickets among the jobs: A - 35 tickets, B - 15 tickets and C - 50 tickets.

We divide the tickets among the jobs: A - 35 tickets, B - 15 tickets and C - 50 tickets.

The scheduler selects a winning ticket by random.

We divide the tickets among the jobs: A - 35 tickets, B - 15 tickets and C - 50 tickets.

The scheduler selects a winning ticket by random.

And the winner is: 23, 56, 13, 73, 8, 82, 17, 34,
We divide the tickets among the jobs: A - 35 tickets, B - 15 tickets and C - 50 tickets.

The scheduler selects a winning ticket by random.

And the winner is: 23, 56, 13, 73, 8, 82, 17, 34,





• A new job can be given a set of tickets as long as we keep track of how many tickets we have.

flexibility

- A new job can be given a set of tickets as long as we keep track of how many tickets we have.
- We can give a *user* a set of tickets and allow the user to distribute them among its jobs.

flexibility

- A new job can be given a set of tickets as long as we keep track of how many tickets we have.
- We can give a *user* a set of tickets and allow the user to distribute them among its jobs.
- Each user can have its local tickets and then have a local lottery.

flexibility

- A new job can be given a set of tickets as long as we keep track of how many tickets we have.
- We can give a *user* a set of tickets and allow the user to distribute them among its jobs.
- Each user can have its local tickets and then have a local lottery.
- We could allow each user to create new tickets, i.e. inflation, if we trust each other.

How to implement?

Stand in line

- Each job is given a number that represents the number of tickets it owns.
- All jobs are lined up i a row.
- Pick a random number from zero to the total number of tickets.
- Walk down the line and select the winner.



How does this work?

Why random?

• Each job is given a *stride value*, the higher the stride the lower the priority.

- Each job is given a *stride value*, the higher the stride the lower the priority.
- Each job keeps a *pass value* initially set to 0.

- Each job is given a *stride value*, the higher the stride the lower the priority.
- Each job keeps a *pass value* initially set to 0.
- In each round the job with the lowest pass value is selected and ...

- Each job is given a *stride value*, the higher the stride the lower the priority.
- Each job keeps a *pass value* initially set to 0.
- In each round the job with the lowest pass value is selected and ...
- ... the pass value is incremented by its stride value.

- Each job is given a *stride value*, the higher the stride the lower the priority.
- Each job keeps a *pass value* initially set to 0.
- In each round the job with the lowest pass value is selected and ...
- ... the pass value is incremented by its stride value.

A low stride value will make it more likely to be scheduled soon again.

• Hard : all deadlines should be met, missing a deadline is a failure.

- Hard : all deadlines should be met, missing a deadline is a failure.
- Soft : deadlines could be missed but the application should be notified and be able to take actions.

- Hard : all deadlines should be met, missing a deadline is a failure.
- Soft : deadlines could be missed but the application should be notified and be able to take actions.

We often have real-time requirements that are simply met since we happen to have the available resources.

<ロト <部ト < Eト < Eト E のQ(~ 37/46

• e: the worst case execution time for the task.

- e: the worst case execution time for the task.
- d: the deadline, when in the future do we need to finish.

- e: the worst case execution time for the task.
- d: the deadline, when in the future do we need to finish.
- p: the period, how often should the task be scheduled.

- e: the worst case execution time for the task.
- d: the deadline, when in the future do we need to finish.
- p: the period, how often should the task be scheduled.



- e: the worst case execution time for the task.
- d: the deadline, when in the future do we need to finish.
- p: the period, how often should the task be scheduled.



- e: the worst case execution time for the task.
- d: the deadline, when in the future do we need to finish.
- p: the period, how often should the task be scheduled.



- e: the worst case execution time for the task.
- d: the deadline, when in the future do we need to finish.
- p: the period, how often should the task be scheduled.



- e: the worst case execution time for the task.
- d: the deadline, when in the future do we need to finish.
- p: the period, how often should the task be scheduled.



- e: the worst case execution time for the task.
- d: the deadline, when in the future do we need to finish.
- p: the period, how often should the task be scheduled.



d < p: constrained, d = p default, d > p several out-standing

T1: T2: T3: 0 20 40 60 80 100 120 140 160 180 200






















Strategies

Given that p = d i.e. a task must be completed within its period.

- Rate Monotonic Scheduling (RMS):
 - Schedule the avilable task with the shortest period
 - Always works if utilization is < 69% (actually less than $n * (2^{1/n} 1)$), where n is the number of processes) could work for higher loads.
 - Simpler to reason about, easy to implement.

Strategies

Given that p = d i.e. a task must be completed within its period.

- Rate Monotonic Scheduling (RMS):
 - Schedule the avilable task with the shortest period
 - Always works if utilization is < 69% (actually less than $n * (2^{1/n} 1)$), where n is the number of processes) could work for higher loads.
 - Simpler to reason about, easy to implement.
- Earliest Deadline First (EDF):
 - Schedule based on the deadline, more freedom to choose tasks.
 - Always works if utilization is < 100%.
 - Used by Linux in the real-time extension (not in the regular system)

Strategies

Given that p = d i.e. a task must be completed within its period.

- Rate Monotonic Scheduling (RMS):
 - Schedule the avilable task with the shortest period
 - Always works if utilization is < 69% (actually less than $n * (2^{1/n} 1)$), where n is the number of processes) could work for higher loads.
 - Simpler to reason about, easy to implement.
- Earliest Deadline First (EDF):
 - Schedule based on the deadline, more freedom to choose tasks.
 - Always works if utilization is < 100%.
 - Used by Linux in the real-time extension (not in the regular system)

Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24

Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



◆□▶ ◆□▶ ◆目▶ ◆目▶ ●目 - のへで

Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - の�?

Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24

Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



・ロト ・ 母 ト ・ ヨ ト ・ ヨ ・ つへぐ

Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



^{41 / 46}

Assume we have tasks: T1: (10, 40, 40), T2: (20, 60, 60), T3: (30, 80, 80).

10/40 + 20/60 + 30/80 = 6/24 + 8/24 + 9/24 = 23/24



problems

Should we be conservative or take a chance?

Should we be conservative or take a chance?

Can we handle a dynamic set of tasks?

Should we be conservative or take a chance?

Can we handle a dynamic set of tasks?

What happens when we have critical resources that are protected by locks?

Multi-core architectures

<ロト < 合 > < 言 > < 言 > 言 の < で 43 / 46 Scheduling for a multi-core architecture more problematic (or rather more problematic to achieve high utilization).

Why?

• O(n) scheduler: the original scheduler, did not scale well.

- $\bullet~O(n)$ scheduler: the original scheduler, did not scale well.
- O(1) scheduler: multi-level feedback queues, dynamic priority, used up to version 2.6

- $\bullet~O(n)$ scheduler: the original scheduler, did not scale well.
- O(1) scheduler: multi-level feedback queues, dynamic priority, used up to version 2.6
- CFS: the completely fair scheduler, O(lg(n)), default today.

- $\bullet~O(n)$ scheduler: the original scheduler, did not scale well.
- O(1) scheduler: multi-level feedback queues, dynamic priority, used up to version 2.6
- CFS: the completely fair scheduler, O(lg(n)), default today.
- BF scheduler: no I will not tell you what it stands for.
• Similar to stride scheduler but uses a red-black tree to order processes.

- Similar to stride scheduler but uses a red-black tree to order processes.
- Will keep processes on the same core if it thinks it's a good choice.

- Similar to stride scheduler but uses a red-black tree to order processes.
- Will keep processes on the same core if it thinks it's a good choice.
- Scheduling classes:

- Similar to stride scheduler but uses a red-black tree to order processes.
- Will keep processes on the same core if it thinks it's a good choice.
- Scheduling classes:
 - SCHED_FIFO, SCHED_RR: high priority classes (often called real-time processes)

- Similar to stride scheduler but uses a red-black tree to order processes.
- Will keep processes on the same core if it thinks it's a good choice.
- Scheduling classes:
 - SCHED_FIFO, SCHED_RR: high priority classes (often called real-time processes)
 - SCHED_NORMAL: all the regular interactive processes

- Similar to stride scheduler but uses a red-black tree to order processes.
- Will keep processes on the same core if it thinks it's a good choice.
- Scheduling classes:
 - SCHED_FIFO, SCHED_RR: high priority classes (often called real-time processes)
 - SCHED_NORMAL: all the regular interactive processes
 - SCHED_BATCH: processes that only run if there are no interactive processes available.

- Similar to stride scheduler but uses a red-black tree to order processes.
- Will keep processes on the same core if it thinks it's a good choice.
- Scheduling classes:
 - SCHED_FIFO, SCHED_RR: high priority classes (often called real-time processes)
 - SCHED_NORMAL: all the regular interactive processes
 - SCHED_BATCH: processes that only run if there are no interactive processes available.
 - SCHED_IDLE: if we've got nothing else to do.

- Similar to stride scheduler but uses a red-black tree to order processes.
- Will keep processes on the same core if it thinks it's a good choice.
- Scheduling classes:
 - SCHED_FIFO, SCHED_RR: high priority classes (often called real-time processes)
 - SCHED_NORMAL: all the regular interactive processes
 - SCHED_BATCH: processes that only run if there are no interactive processes available.
 - SCHED_IDLE: if we've got nothing else to do.

• Bonnie Tyler: Turnaround, every now and then ...

- Bonnie Tyler: Turnaround, every now and then ...
- Bob Marley: Talking 'bout reaction

- Bonnie Tyler: Turnaround, every now and then ...
- Bob Marley: Talking 'bout reaction
- Rolling Stones: You can't always get what you want.

- Bonnie Tyler: Turnaround, every now and then ...
- Bob Marley: Talking 'bout reaction
- Rolling Stones: You can't always get what you want.
- Metallica: Justice for all.

- Bonnie Tyler: Turnaround, every now and then ...
- Bob Marley: Talking 'bout reaction
- Rolling Stones: You can't always get what you want.
- Metallica: Justice for all.
- Leif "Loket" Olsson: a lottery might work ok

- Bonnie Tyler: Turnaround, every now and then ...
- Bob Marley: Talking 'bout reaction
- Rolling Stones: You can't always get what you want.
- Metallica: Justice for all.
- Leif "Loket" Olsson: a lottery might work ok
- Real-time scheduling: if we actually know the maximum execution time, the deadline and the period.

- Bonnie Tyler: Turnaround, every now and then ...
- Bob Marley: Talking 'bout reaction
- Rolling Stones: You can't always get what you want.
- Metallica: Justice for all.
- Leif "Loket" Olsson: a lottery might work ok
- Real-time scheduling: if we actually know the maximum execution time, the deadline and the period.
- Multi-core schedulers: you have to think twice before selecting a process.

- Bonnie Tyler: Turnaround, every now and then ...
- Bob Marley: Talking 'bout reaction
- Rolling Stones: You can't always get what you want.
- Metallica: Justice for all.
- Leif "Loket" Olsson: a lottery might work ok
- Real-time scheduling: if we actually know the maximum execution time, the deadline and the period.
- Multi-core schedulers: you have to think twice before selecting a process.
- Linux: Completely Fair Scheduler, schedules in O(lg(n)) time, similar to stride scheduling.