

CS330: Operating Systems

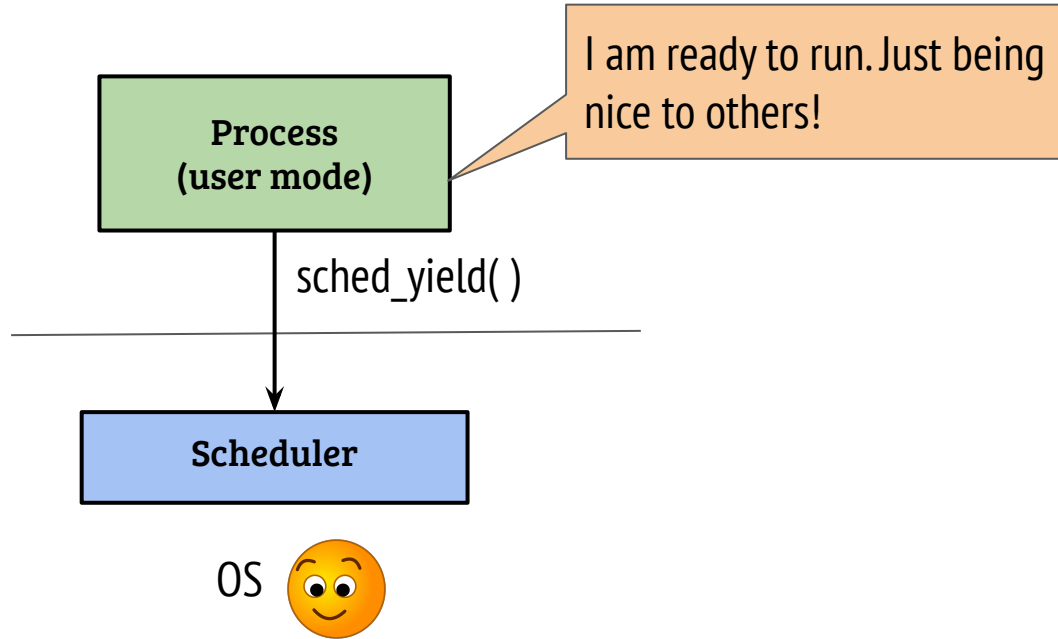
Process scheduling

Recap: OS execution, user-OS mode switch

- Which stack is used by the OS for kernel-mode execution?
- The hardware switches the SP to point it to a pre-configured per-process OS stack on mode switch
- How the user process state preserved and restored?
- The user execution state is saved/restored using the kernel stack by the hardware (and OS)
- Which address space the OS uses?
- A part of the process address space is reserved for OS and is protected

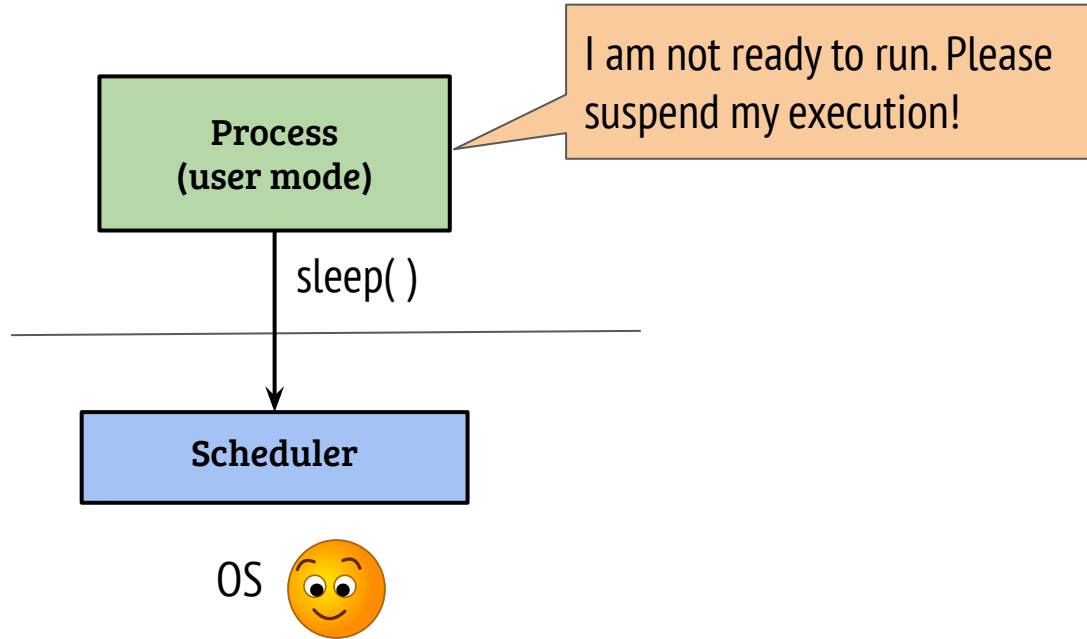
Today's agenda: Process context switch and scheduling

Triggers for process context switch



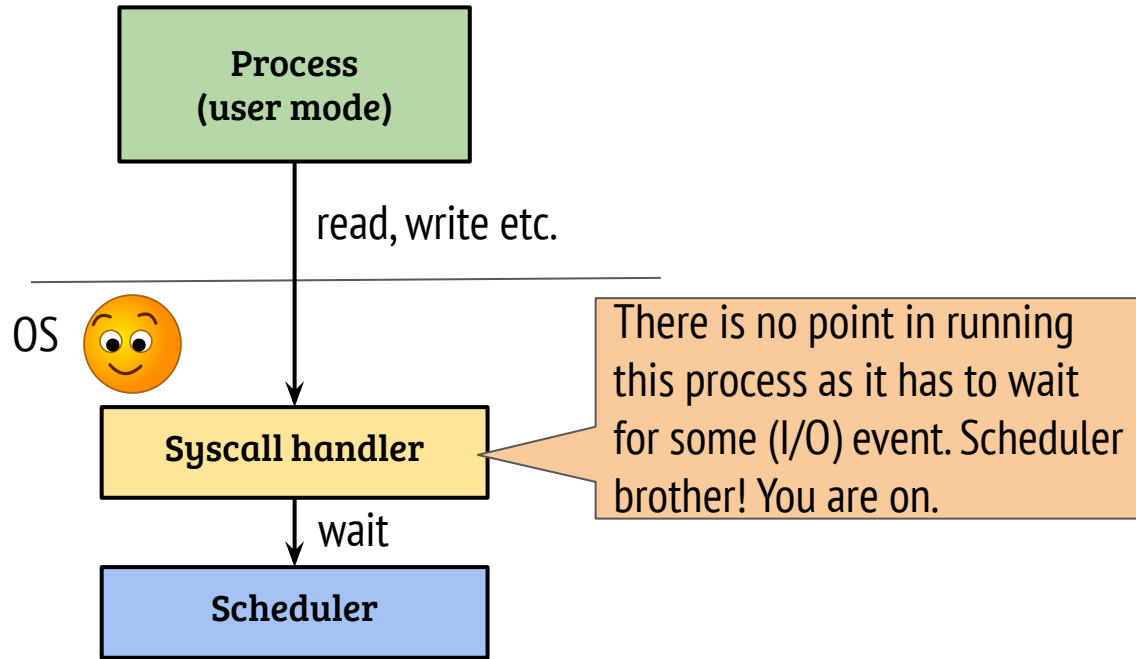
- The user process can invoke the scheduler through explicit system calls like `sched_yield` (see man page)

Triggers for process context switch



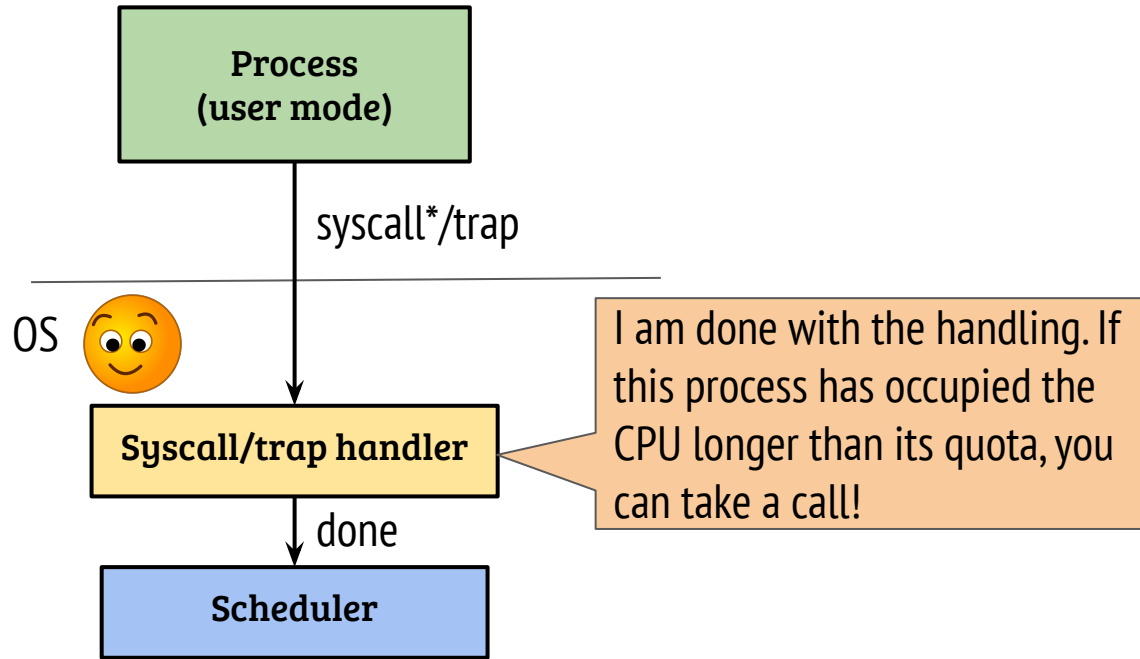
- The user process can invoke `sleep()` to suspend itself
 - `sleep()` is not a system call in Linux, it uses `nanosleep()` system call

Triggers for process context switch



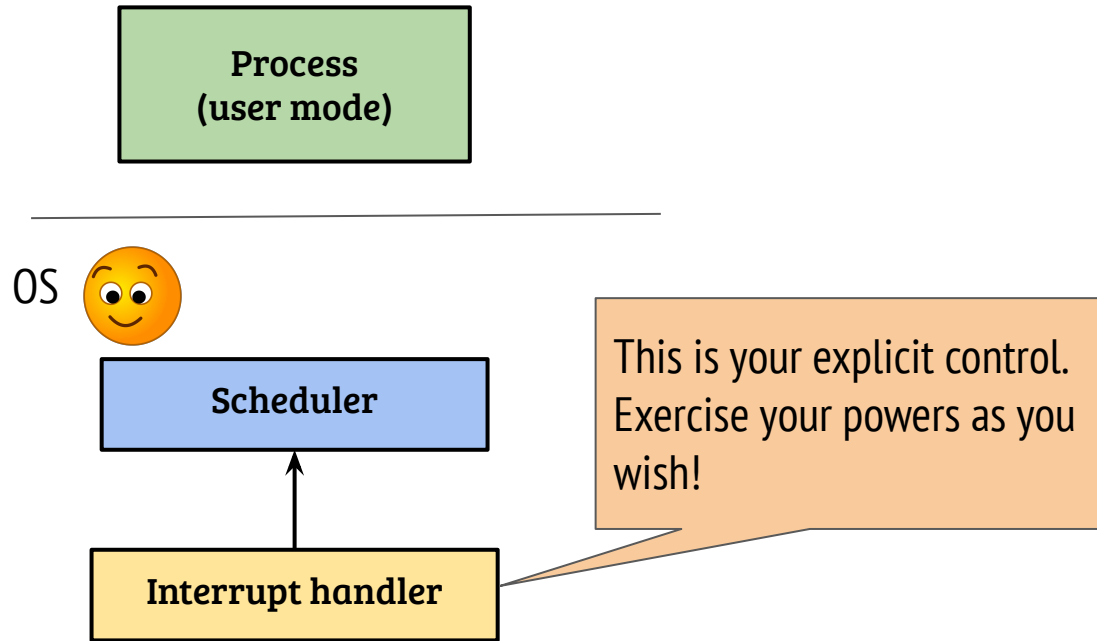
- This condition arises mostly during I/O related system calls
 - Example: `read()` from a file on disk

Triggers for process context switch



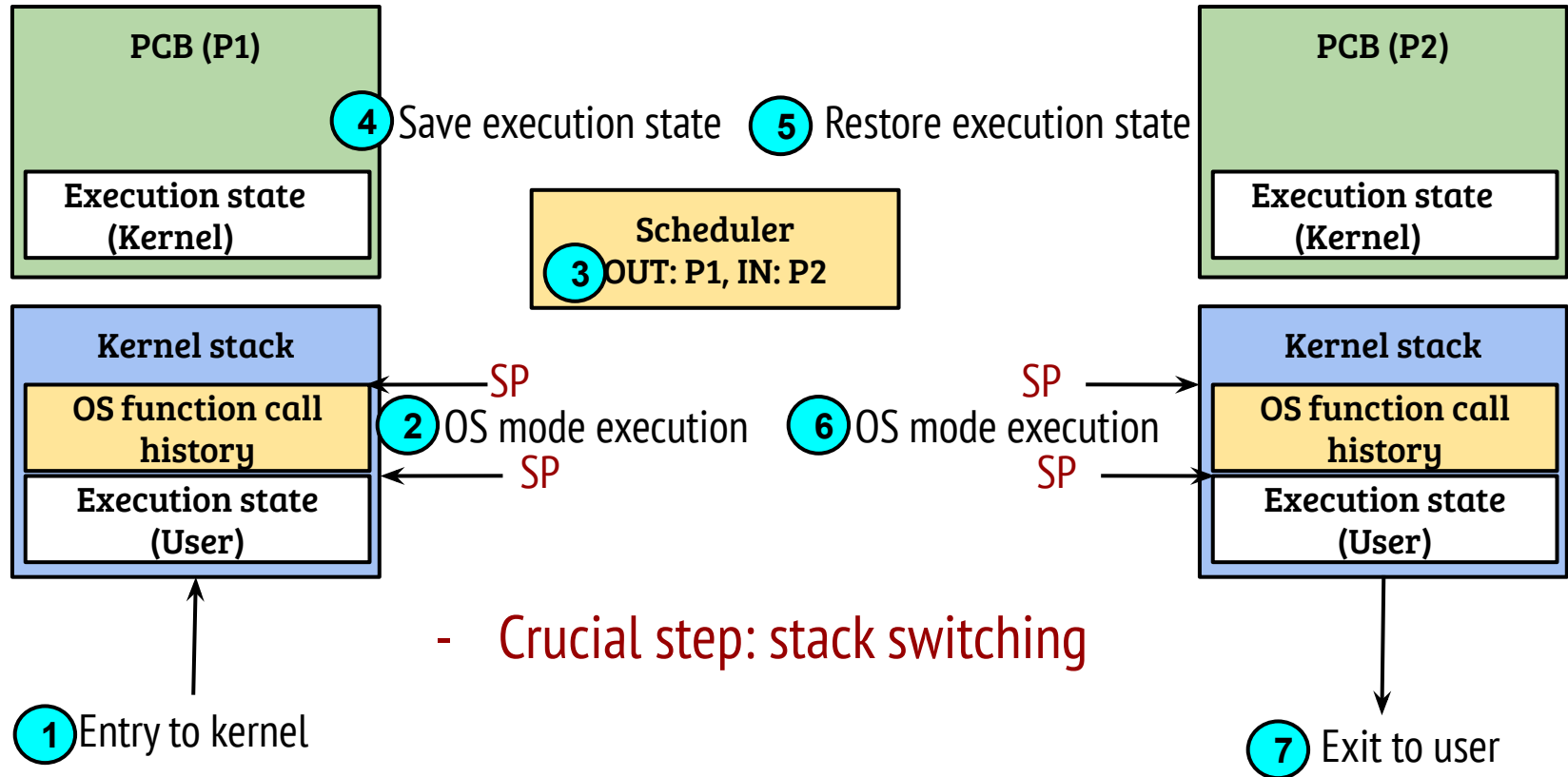
- The OS gets the control back on every system call and exception
- Before returning from syscall, the scheduler can deschedule

Triggers for process context switch

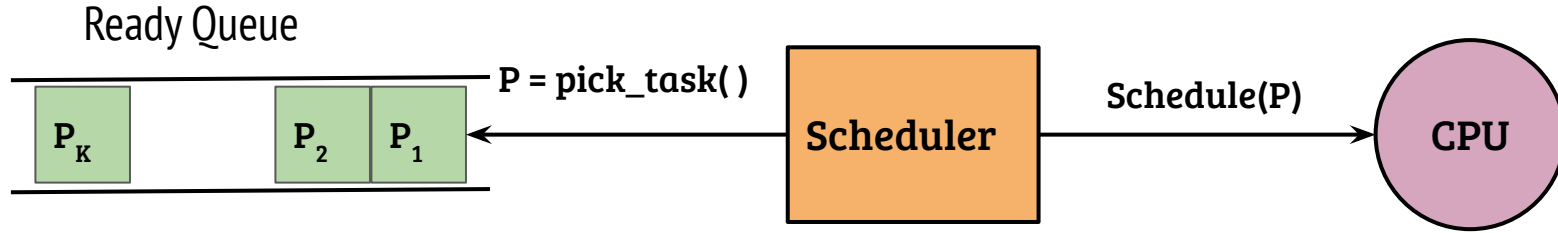


- Timer interrupts can be configured to generate interrupts periodically or after some configured time
- The OS can invoke the scheduler after handling any interrupt

Process context switch



Scheduling



- A queue of processes ready to execute is maintained
- The scheduler decides to pick the next process based on some scheduling policy and performs a context switch
- The outgoing process is put back to ready queue (if required)

Scheduling

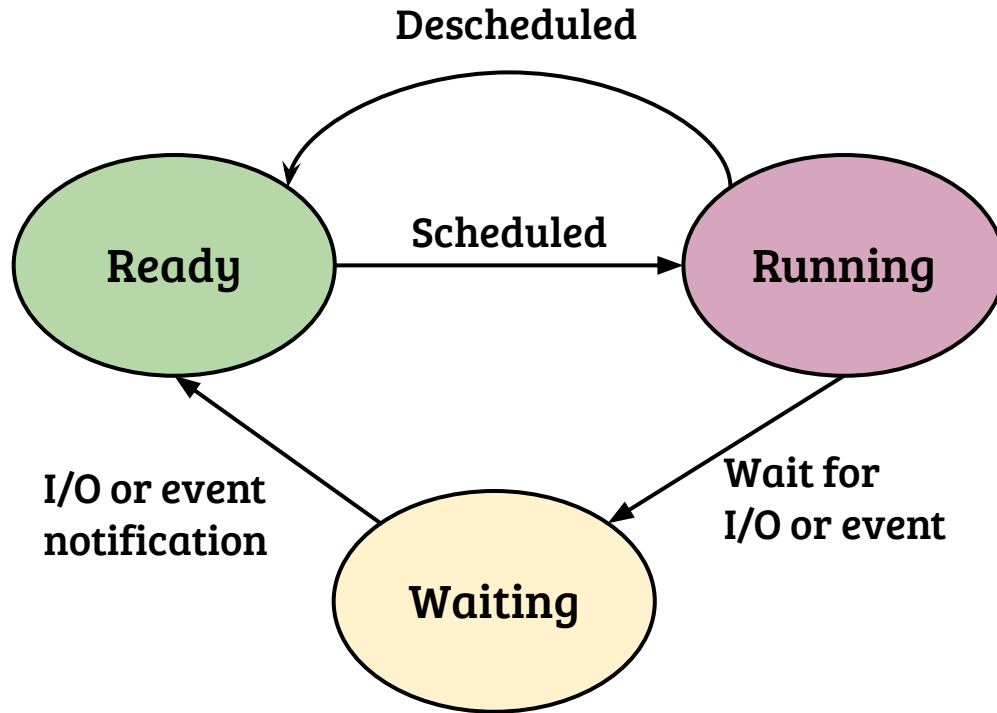


- How is the list of ready processes managed?
- What if there are no processes in ready queue? Can that happen?
- Can we classify the schedulers based on how they are invoked?
- What is a good scheduling strategy?

policy and performs a context switch

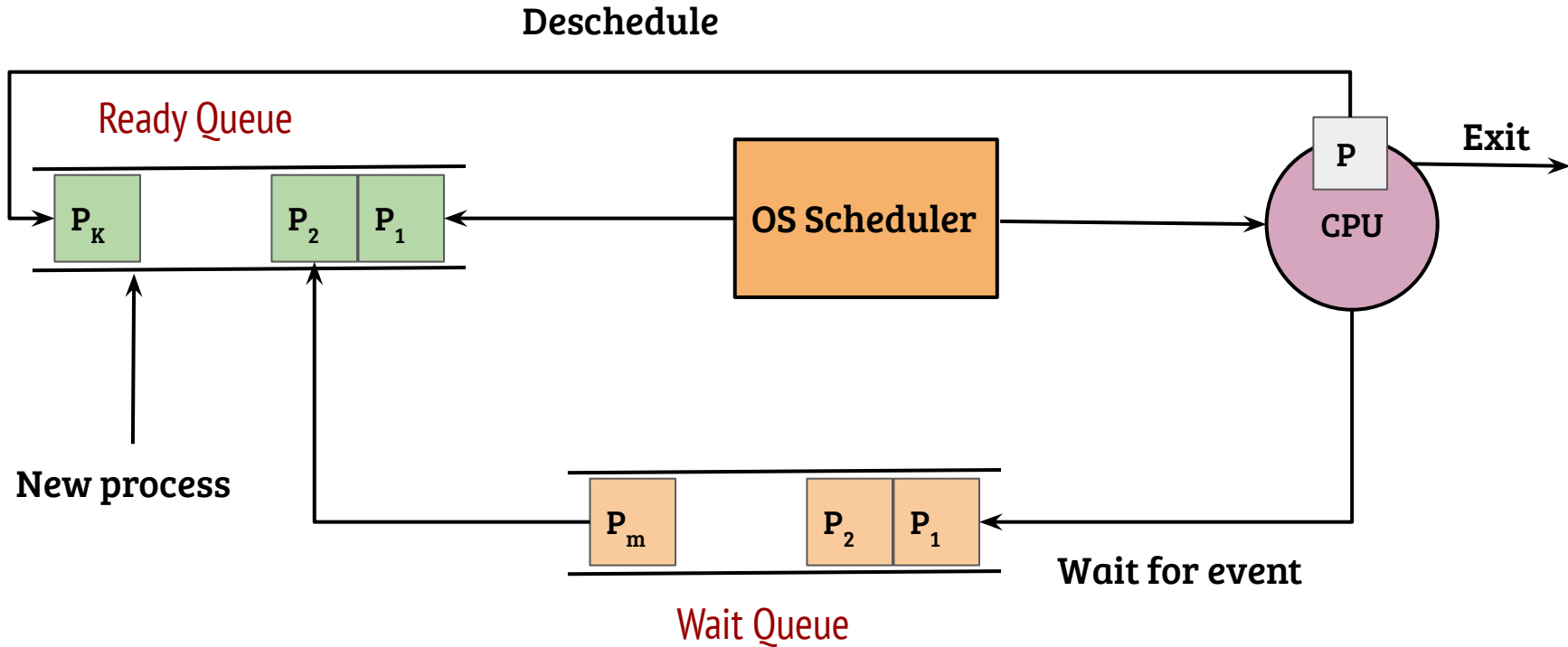
- The outgoing process is put back to ready queue (if required)

Process states and transitions



- Most processes perform a mixture of CPU and I/O activities
- When the process is waiting for an I/O, it is moved to waiting state
- A process becomes ready again when the event completion is notified (e.g., a device interrupt)

Scheduler overview

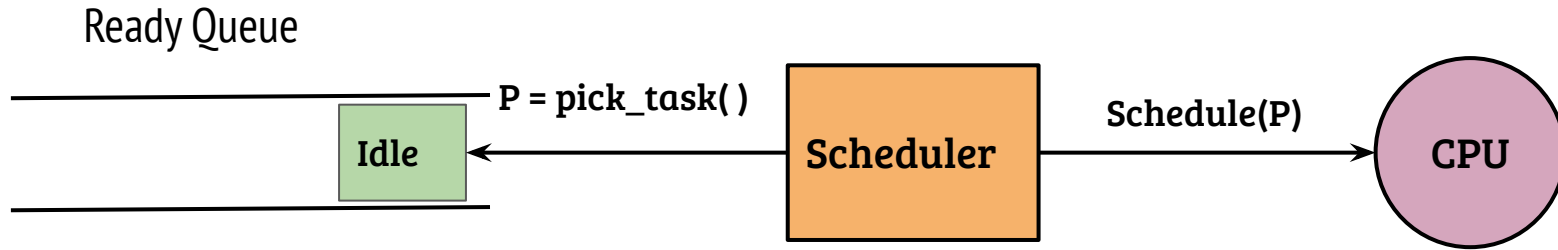


Scheduling

Ready Queue

- How is the list of ready processes managed?
- Each process is associated with three primary states: Running, Ready and Waiting. A process can be moved to waiting state from running state, if needed.
- What if there are no processes in ready queue? Can that happen?
- Can we classify the schedulers based on how they are invoked?
- What is a good scheduling strategy?
- The outgoing process is put back to ready queue (if required)

System idle process



- There can be an instance when there are zero processes in ready queue
- A special process (system idle process) is always there
- The system idle process halts the CPU
- HLT instruction on X86_64: Halts the CPU till next interrupt

Scheduling

Ready Queue

- How is the list of ready processes managed?
- Each process is associated with three primary states: Running, Ready and Waiting. A process can be moved to waiting state from running state, if needed.
- What if there are no processes in ready queue? Can that happen?
- There is always an idle process which executes HLT
- Can we classify the schedulers based on how they are invoked?
- What is a good scheduling strategy?

Scheduling: preemptive vs. non-preemptive

- There are scheduling points which are triggered because of the current process execution behavior (non-preemptive)
 - Process termination
 - Process explicitly yields the CPU
 - Process waits/blocks for an I/O or event

Scheduling: preemptive vs. non-preemptive

- There are scheduling points which are triggered because of the current process execution behavior (non-preemptive)
 - Process termination
 - Process explicitly yields the CPU
 - Process waits/blocks for an I/O or event
- The OS may invoke the scheduler in other conditions (preemptive)
 - Return from system call (specifically `fork()`)
 - After handling an interrupt (specifically timer interrupt)

Scheduling

Ready Queue

- How is the list of ready processes managed?
- Each process is associated with three primary states: Running, Ready and Waiting. A process can be moved to waiting state from running state, if needed.
- What if there are no processes in ready queue? Can that happen?
- There is always an idle process which executes HLT
- Can we classify the schedulers based on how they are invoked?
- Non-preemptive: triggered by the process, Preemptive: OS interjections
- What is a good scheduling strategy?

Scheduling metrics

- Turnaround time: Time of completion - Time of arrival
 - Objective: *Minimize turnaround time*

Scheduling metrics

- Turnaround time: Time of completion - Time of arrival
 - Objective: *Minimize turnaround time*
- Waiting time: Sum of time spent in ready queue
 - Objective: *Minimize waiting time*

Scheduling metrics

- Turnaround time: Time of completion - Time of arrival
 - Objective: *Minimize turnaround time*
- Waiting time: Sum of time spent in ready queue
 - Objective: *Minimize waiting time*
- Response time: Waiting time before first execution
 - Objective: *Minimize response time*

Scheduling metrics

- Turnaround time: Time of completion - Time of arrival
 - Objective: *Minimize turnaround time*
- Waiting time: Sum of time spent in ready queue
 - Objective: *Minimize waiting time*
- Response time: Waiting time before first execution
 - Objective: *Minimize response time*
- Average value of above metrics represent the average efficiency

Scheduling metrics

- Turnaround time: Time of completion - Time of arrival
 - Objective: *Minimize turnaround time*
- Waiting time: Sum of time spent in ready queue
 - Objective: *Minimize waiting time*
- Response time: Waiting time before first execution
 - Objective: *Minimize response time*
- Average value of above metrics represent the average efficiency
- Standard deviation represents fairness across different processes

First Come First Served (FCFS)

- FIFO queue based non-preemptive scheduling
- Example

First Come First Served (FCFS)

- FIFO queue based non-preemptive scheduling
- Example
- Advantages
 - Easy to implement
- Issues with FCFS
 - Convoy effect
 - Not suitable for interactive applications

Shortest Job First (SJF)

- Select the process with shortest CPU burst
- Pick the next process only when the current process is finished (non-preemptive)
- Example

Shortest Job First (SJF)

- Select the process with shortest CPU burst
- Pick the next process only when the current process is finished (non-preemptive)
- Example
- Optimal on waiting time and turnaround time
- Not realistic (how can we know the execution time?)

Shortest Time to Completion First (STCF)

- Pick the process with shortest remaining time when a new process arrives in the ready queue (SRTF)
- Example
- Improves the efficiency of SJF at the cost of more context switches

Round-robin scheduling

- Preemptive scheduling with time slicing
- Ready queue is maintained as a circular queue
- At end of the time quantum, If there are other processes in the queue
 - Current process goes to the TAIL of the queue
 - Next process is picked up from the HEAD of the queue
- New processes are added to the TAIL of the queue
- Design choice: size of time quantum

Priority scheduling

- Select the process with highest priority
- Can be preemptive and non-preemptive
- SJF: priority defined by job length
- Advantages: practical (no assumptions)
- Disadvantages: Starvation

Problem formulation with I/O bursts

Process	Arrival Time	CPU bursts	I/O bursts
P1	0	0-3, 7-9, 14-15	3-7,9-14
P2	2	2-10, 12-15	10-12
P3	3	3-4, 10-11	4-10

- Most processes goes through a series of CPU and I/O bursts
- Looks complicated for analysis, can it be simplified?

Problem formulation with I/O bursts

Process	Arrival Time	CPU bursts	I/O bursts
P1	0	0-3, 7-9, 14-15	3-7,9-14
P2	2	2-10, 12-15	10-12
P3	3	3-4, 10-11	4-10

- Most processes require a series of CPU and I/O bursts
- Looks complicated for analysis, can it be simplified?
- Every CPU burst is treated as a new process where the CPU burst start is the process arrival time and burst length is the execution time

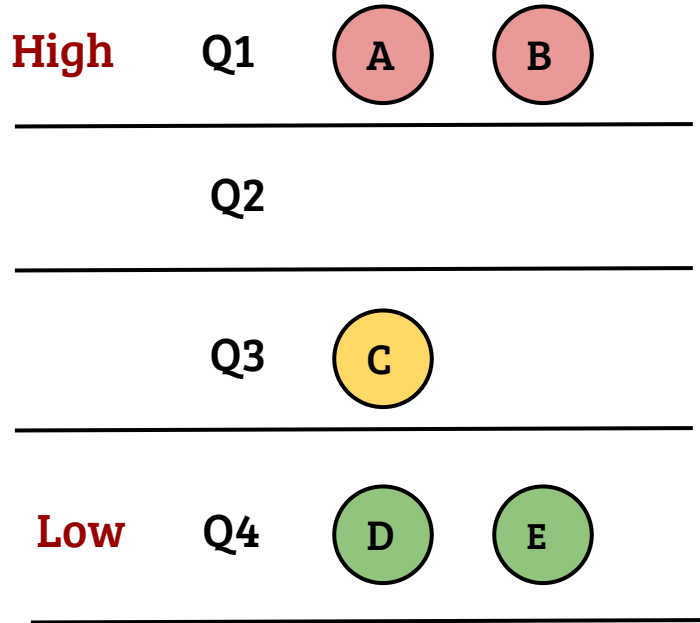
CS330: Operating Systems

Process scheduling policies

Recap: basic scheduling policies

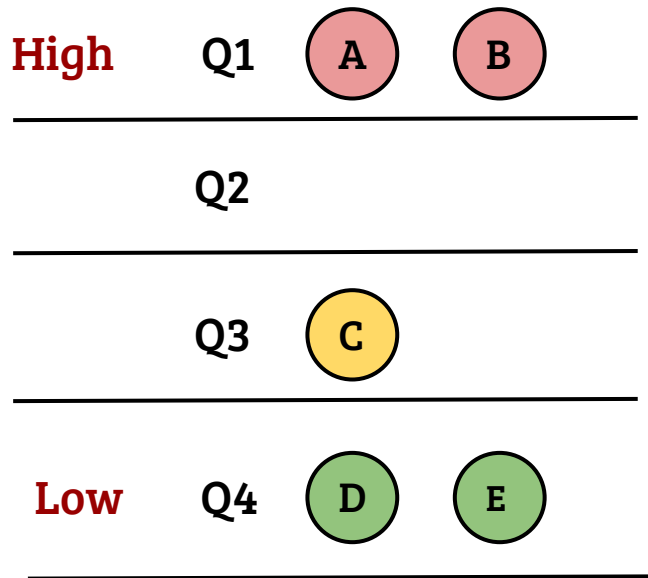
- Scheduling metrics: turnaround time, waiting time, response time
- Fast come first serve (FCFS)
 - Simple but inefficient (convoy effect)
- Shortest job first (SJF) and Shortest time to completion first (STCF)
 - Optimal and efficient. Issues: unrealistic, starvation
- Round robin (RR)
 - Good response time, Issues: scheduling overheads
- Priority scheduling
 - Starvation

Static priority based scheduling



- Processes are assigned to different queues based on their priority
- Process from the non-empty highest priority queue is always picked
- Different queues may implement different schemes within a queue
- Main concern: Starvation
 - Ex: Low priority processes hug the CPU

Multilevel feedback queue

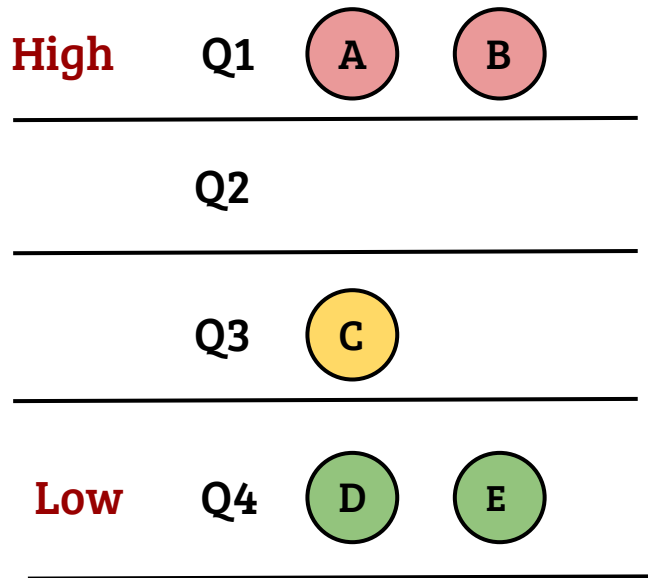


OS



- Dynamically adjust priorities such that
1. Interactive applications are responsive
 2. Short jobs do not suffer
 3. No starvation
 4. No user can trick the scheduler

Multilevel feedback queue



Dynamically adjust priorities such that

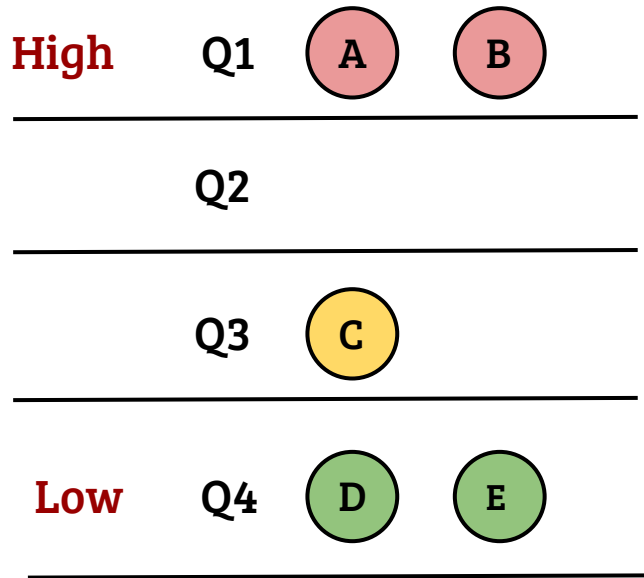
1. Interactive applications are responsive
2. Short jobs do not suffer
3. No starvation
4. No user can trick the scheduler

OS



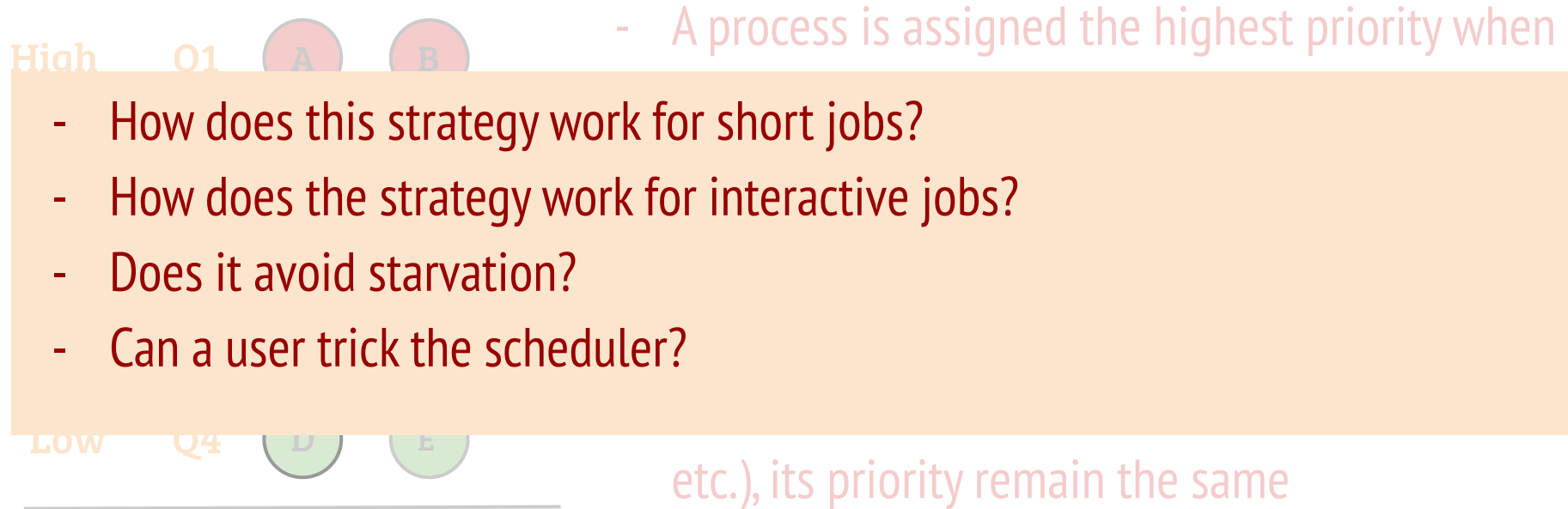
- Basic multi level strategy
 - Pick a process from highest priority queue
 - Within a queue, apply RR

Multilevel feedback queue: Dynamic priorities



- A process is assigned the highest priority when it is created
- If the process consumes the slice (scheduler invoked because of timer), its priority is reduced
- If the process relinquishes the CPU (I/O wait etc.), its priority remain the same

Multilevel feedback queue: Dynamic priorities



MLFQ: Approximation of SJF

- MLFQ can approximate SJF because
 - Long running jobs are moved to low priority queues
 - New jobs are added to highest priority queue
- A shorter job may not get a chance to execute for a small duration. What is the upper bound?

MLFQ: Approximation of SJF

- MLFQ can approximate SJF because
 - Long running jobs are moved to low priority queues
 - New jobs are added to highest priority queue
- A shorter job may not get a chance to execute for a small duration. What is the upper bound?
- $(\# \text{ of jobs in the highest priority queue} + 1) \times (\text{time quantum})$

Multilevel feedback queue: Dynamic priorities

High

O1



- A process is assigned the highest priority when

- How does this strategy work for short jobs?
- Works nicely, approximates SJF
- How does the strategy work for interactive jobs?
- Does it avoid starvation?
- Can a user trick the scheduler?

etc.), its priority remain the same

MLFQ: Interactive jobs

- MLFQ favors interactive jobs because
 - Interactive jobs maintain the highest priority as they relinquish the CPU before quantum expires
 - Long running jobs are moved to low priority queues

MLFQ: Interactive jobs

- MLFQ favors interactive jobs because
 - Interactive jobs maintain the highest priority as they relinquish the CPU before quantum expires
 - Long running jobs are moved to low priority queues
- Conclusion: In a steady state, interactive jobs compete with short and other interactive jobs

Multilevel feedback queue: Dynamic priorities

- How does this strategy work for short jobs?
- Works nicely, approximates SJF
- How does the strategy work for interactive jobs?
- Works pretty well as interactive jobs retain priority
- Does it avoid starvation?
- Can a user trick the scheduler?

MLFQ: Starvation and other issues

- Long running processes may starve with the proposed scheme
- Additionally, permanent demotion of priority hurts processes which change their behavior
 - Example: A process performing a lot of computation only at start gets pushed to a low priority queue permanently
- How to avoid the above issues?

MLFQ: Starvation and other issues

- Long running processes may starve with the proposed scheme
- Additionally, permanent demotion of priority hurts processes which change their behavior
 - Example: A process performing a lot of computation only at start gets pushed to a low priority queue permanently
- How to avoid the above issues?
 - Periodic priority boost: all processes moved to high priority queue
 - Priority boost with aging: recalculate the priority based on scheduling history of a process

Multilevel feedback queue: Dynamic priorities

- How does this strategy work for short jobs?
- Works nicely, approximates SJF
- How does the strategy work for interactive jobs?
- Works pretty well as interactive jobs retain priority
- Does it avoid starvation?
- No. Requires additional mechanism like priority boost.
- Can a user trick the scheduler?

MLFQ: The tricky user

- A smart user can maintain highest priority for long running processes by exploiting the scheduling strategy. How?

MLFQ: The tricky user

- A smart user can maintain highest priority for long running processes by exploiting the scheduling strategy. How?
- Assumption: user knows the time quantum

MLFQ: The tricky user

- A smart user can maintain highest priority for long running processes by exploiting the scheduling strategy. How?
- Assumption: user knows the time quantum
- Strategy: Voluntarily release the CPU before time quantum expires
- Result: Batch process competes with other interactive processes!

MLFQ: The tricky user

- A smart user can maintain highest priority for long running processes by exploiting the scheduling strategy. How?
- Assumption: user knows the time quantum
- Strategy: Voluntarily release the CPU before time quantum expires
- Result: Batch process competes with other interactive processes!
- Core of the issue: binary history regarding a process

MLFQ: The tricky user

- A smart user can maintain highest priority for long running processes by exploiting the scheduling strategy. How?
- Assumption: user knows the time quantum
- Strategy: Voluntarily release the CPU before time quantum expires
- Result: Batch process competes with other interactive processes!
- Core of the issue: binary history regarding a process
 - MLFQ: Process consumed or not consumed the quantum
 - Advanced MLFQ: Better accounting, variable quanta

Multilevel feedback queue: Dynamic priorities

- How does this strategy work for short jobs?
- Works nicely, approximates SJF
- How does the strategy work for interactive jobs?
- Works pretty well as interactive jobs retain priority
- Does it avoid starvation?
- No. Requires additional mechanism like priority boost.
- Can a user trick the scheduler?
- Yes. Additional history regarding execution is required to be maintained

CS330: Operating Systems

Process scheduling policies

Recap: MLFQ

- MLFQ: Dynamic priority based on process characteristics
- How does this strategy work for short and interactive jobs?
- Approximates SJF for short jobs, interactive jobs retain higher priority levels
- Does it avoid starvation?
- No. Requires additional mechanism like priority boost.
- Can a user trick the scheduler?
- Yes. Additional history regarding execution is required to be maintained

Today's agenda: Overview of Linux scheduling

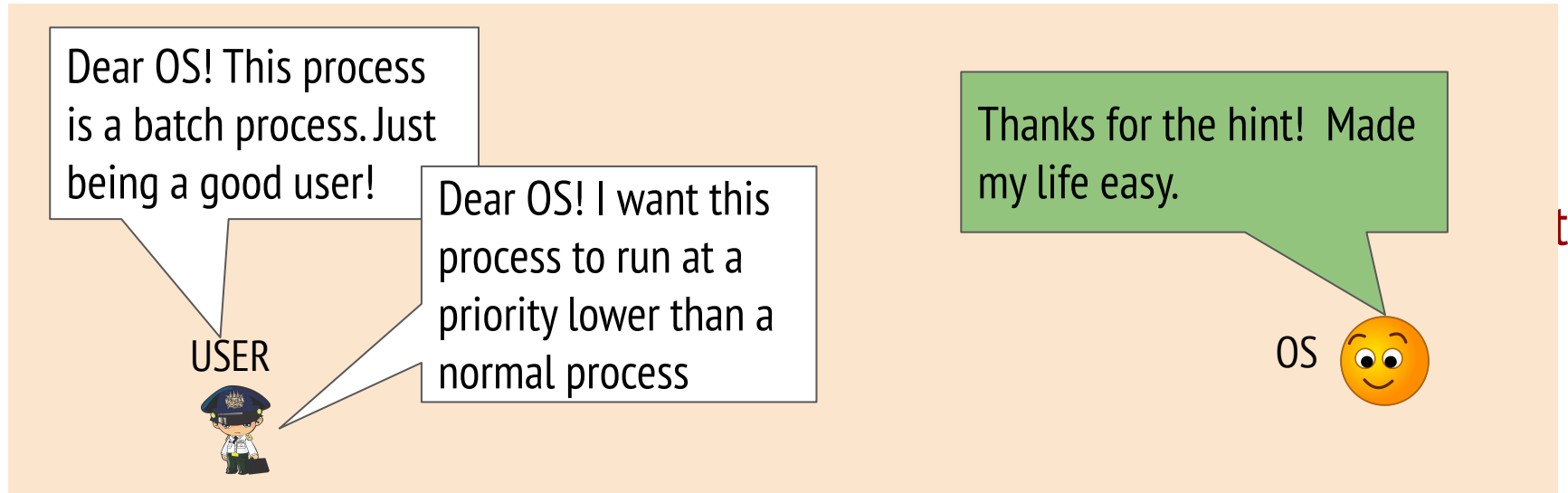
Scheduling is much more complex in a real OS!

- Scheduling requirement of different processes in the system are different
 - Real-time processes: Should meet strict deadlines
 - Interactive processes: Responsive scheduling
 - Batch processes: Starvation free scheduling

Scheduling is much more complex in a real OS!

- Scheduling requirement of different processes in the system are different
 - Real-time processes: Should meet strict deadlines
 - Interactive processes: Responsive scheduling
 - Batch processes: Starvation free scheduling
- Well intentioned users should be able to influence the scheduling policy in a positive manner

Scheduling is much more complex in a real OS!



- Well intentioned users should be able to influence the scheduling policy in a positive manner

Scheduling is much more complex in a real OS!

- Scheduling requirement of different processes in the system are different
 - Real-time processes: Should meet strict deadlines
 - Interactive processes: Responsive scheduling
 - Batch processes: Starvation free scheduling
- Well intentioned users should be able to influence the scheduling policy in a positive manner
- Greed of greedy users should be controlled by the OS

Scheduling is much more complex in a real OS!

Dear OS! This process requires higher priority than other normal processes. You know what, it is very interactive.

Not really! Just trying to fool you.

USER



Buddy! You can fool me for a little while. I will catch you eventually.

OS

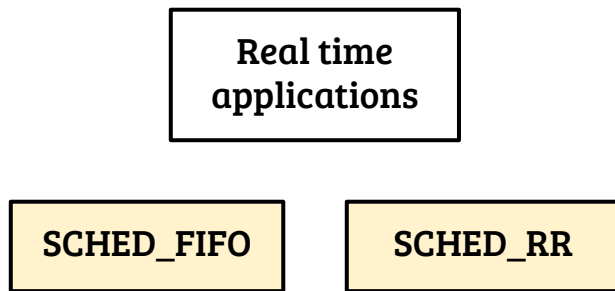


- Greed of greedy users should be controlled by the OS

Scheduling is much more complex in a real OS!

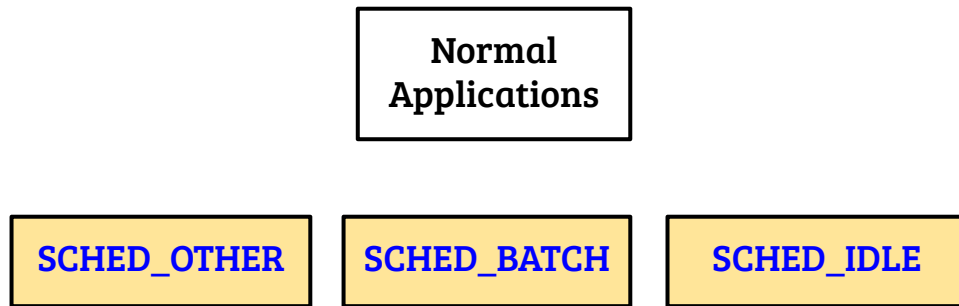
- Scheduling requirement of different processes in the system are different
 - Real-time processes: Should meet strict deadlines
 - Interactive processes: Responsive scheduling
 - Batch processes: Starvation free scheduling
- Well intentioned users should be able to influence the scheduling policy in a positive manner
- Greed of greedy users should be controlled by the OS
- Conclusion: OS scheduling should provide flexibility while being auto-tuning in nature

Linux scheduling classes: Real time applications



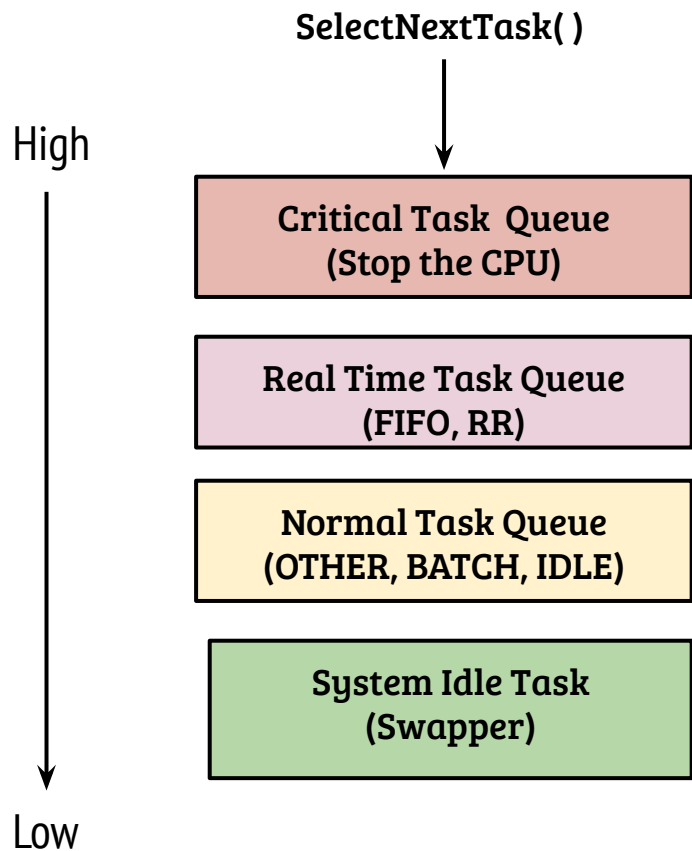
- Real time applications are always higher priority than normal processes
- Priority value: 1 to 99 (In Linux, lower value \Rightarrow higher priority)
- FIFO: Run to completion
- RR: Round robin within a given priority-level
- *sched_setscheduler* system call to define scheduling class and priorities

Linux scheduling classes: normal applications



- SCHED_OTHER: Default policy, OS dynamic priorities and variable time slicing comes into picture
- SCHED_BATCH: Assume CPU bound while calculating dynamic priorities
- SCHED_IDLE: Very low priority jobs

Selecting the next task



- A task is picked from the non-empty highest priority queue
- Critical task queue contains tasks which require immediate attention: hardware events, restart etc.
- Normal task queue (a.k.a fair scheduling class) implements the heuristics to self-adjust
- If all the queues are empty, swapper task is scheduled (HLT the CPU)

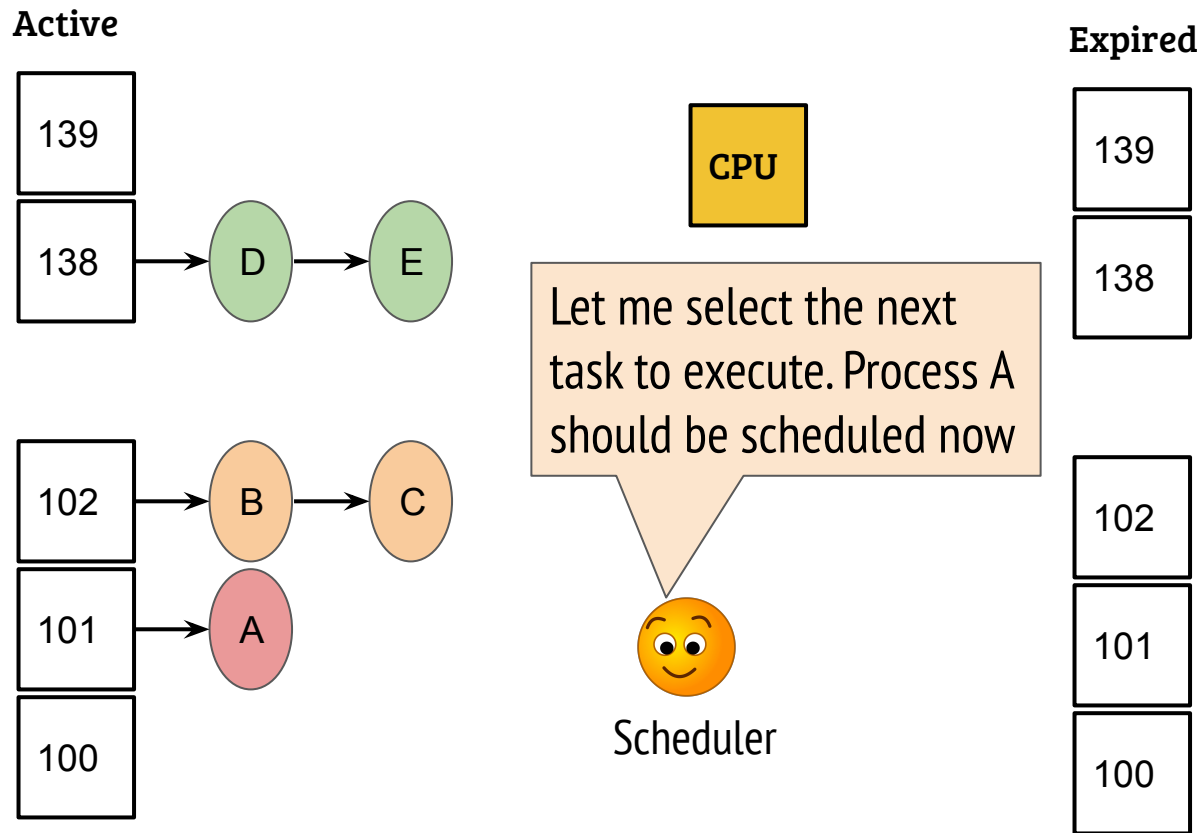
Normal (fair) scheduling class

- 40 priority levels (100 to 139)
- Every process starts with a default priority of 120
- Linux provides *nice* system call to adjust the static priority
 - *nice(int x)*, where x is between 19 to -20
 - *nice(19)* \Rightarrow Move the process to lowest priority queue i.e., 139
 - *nice(-20)* \Rightarrow Move the process to highest priority queue i.e., 100

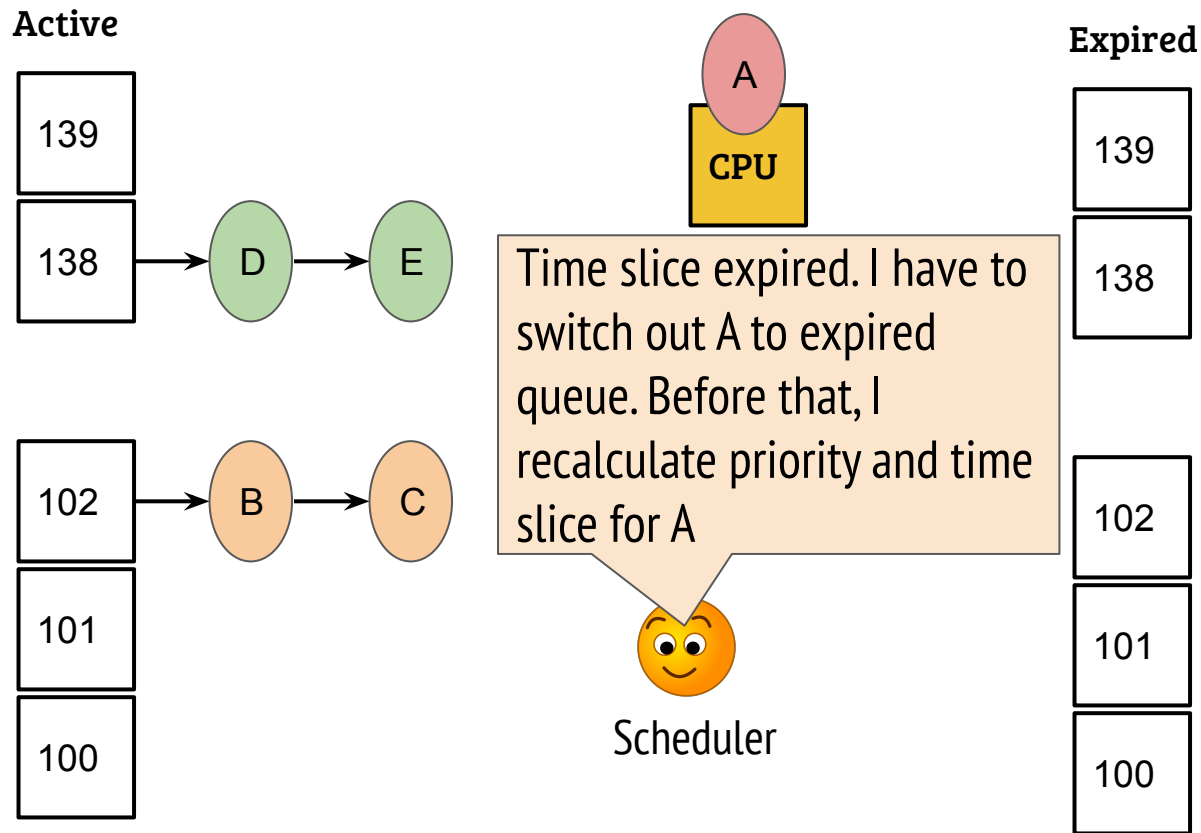
Normal (fair) scheduling class

- 40 priority levels (100 to 139)
- Every process starts with a default priority of 120
- Linux provides *nice* system call to adjust the static priority
 - *nice(int x)*, where x is between 19 to -20
 - *nice(19)* \Rightarrow Move the process to lowest priority queue i.e., 139
 - *nice(-20)* \Rightarrow Move the process to highest priority queue i.e., 100
- Dynamic priority is calculated by the Linux kernel considering the interactiveness of the process
 - More interactive processes move towards the priority level 100

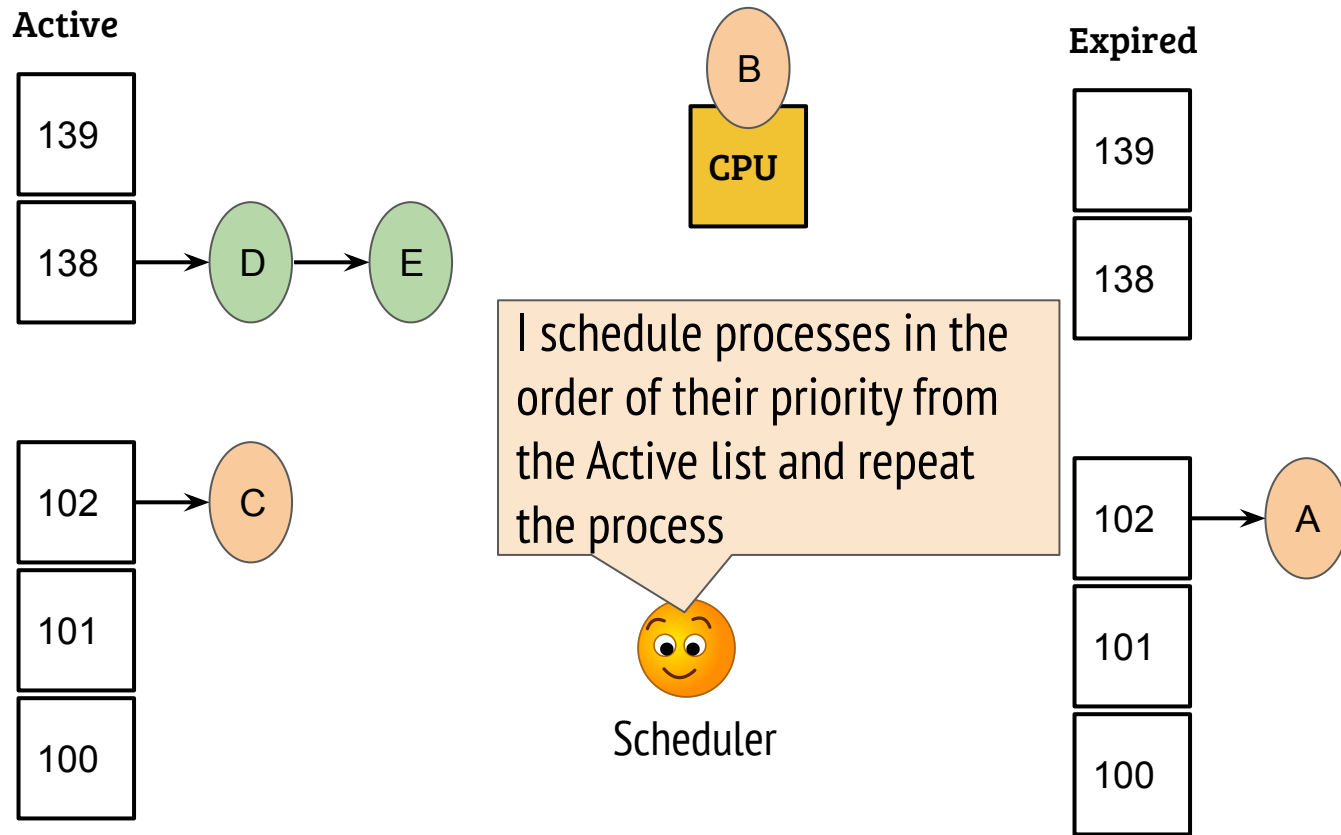
Linux O(1) scheduler



Linux O(1) scheduler

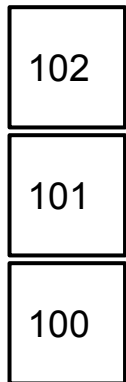
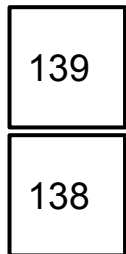


Linux O(1) scheduler



Linux O(1) scheduler

Active



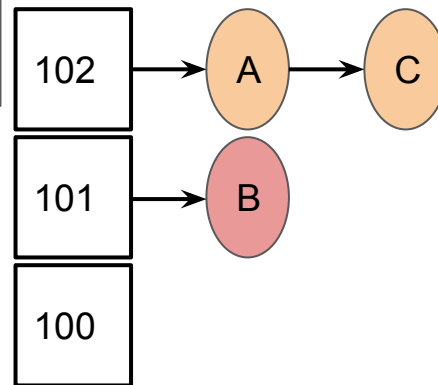
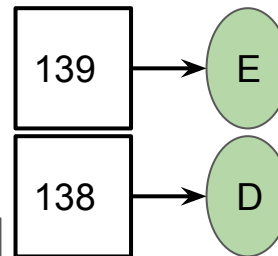
CPU

All the processes in *Active* list are finished. Let me swap the lists. *Expired* is now *Active*



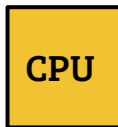
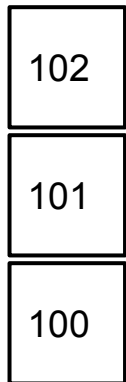
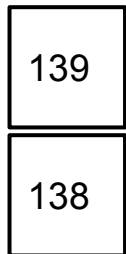
Scheduler

Expired



Linux O(1) scheduler

Expired

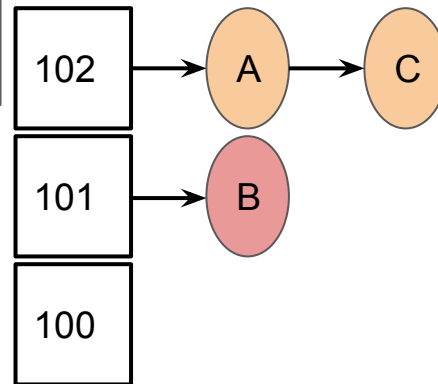
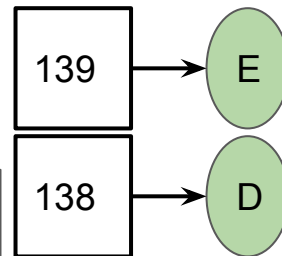


How is it O(1)? Because I do not search a global list of processes. Moreover, scanning the priority levels can be avoided if I maintain a bitmap of priority levels.



Scheduler

Active



$O(1)$ scheduler: value of time slice

- Objective: reduce timer interrupts (tickless system)
- High priority processes are given big time slices
 - Interactive processes relinquish CPU before the quantum expiry
- Low priority processes are given small time slices
 - Should not starve the interactive applications

$O(1)$ scheduler: value of time slice

- Objective: reduce timer interrupts (tickless system)
- High priority processes are given big time slices
 - Interactive processes relinquish CPU before the quantum expiry
- Low priority processes are given small time slices
 - Should not starve the interactive applications
- Result: In a busy system, low priority processes execute less frequently resulting in few timer interrupts