# CPU Virtualization: Limited Direct Execution (LDE)

*CS 571: Operating Systems (Spring 2021)*
Lecture 2a
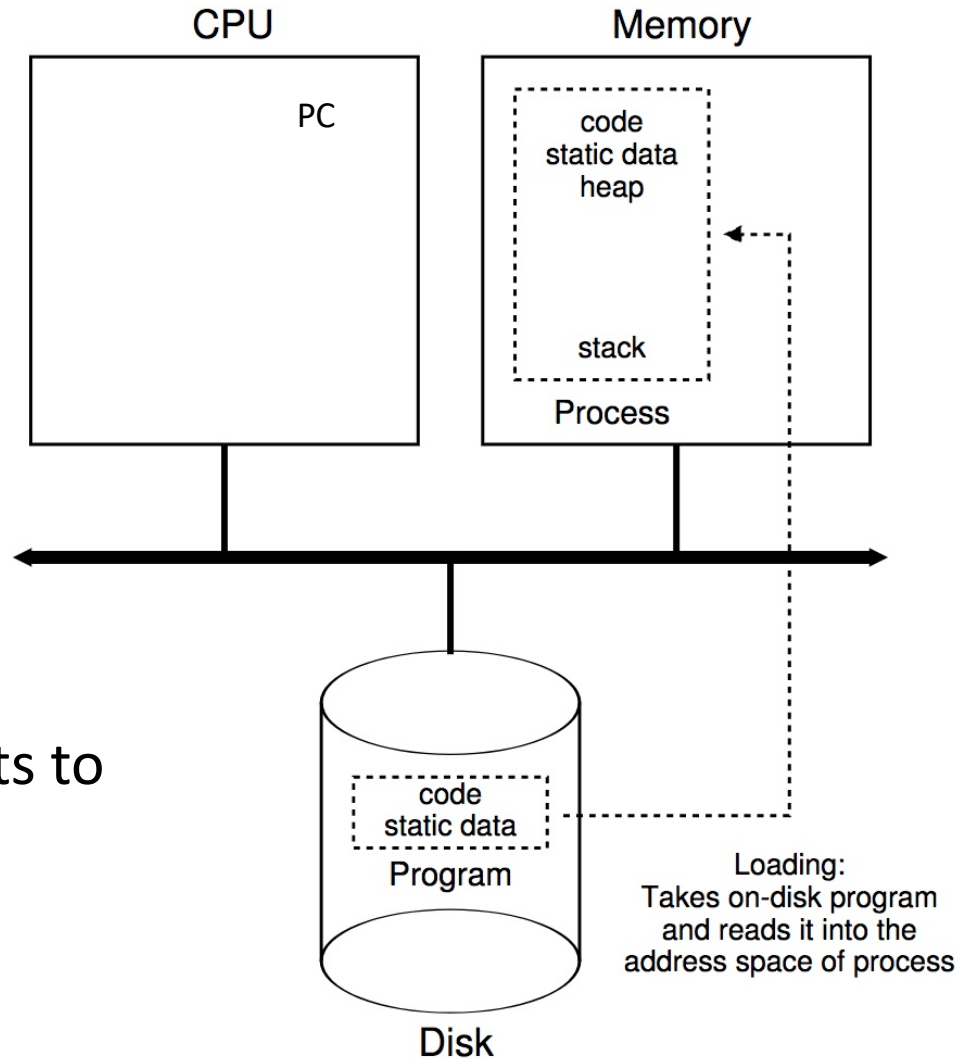
Yue Cheng

# Process Creation



Before, PC points to kernel code

# Process Creation



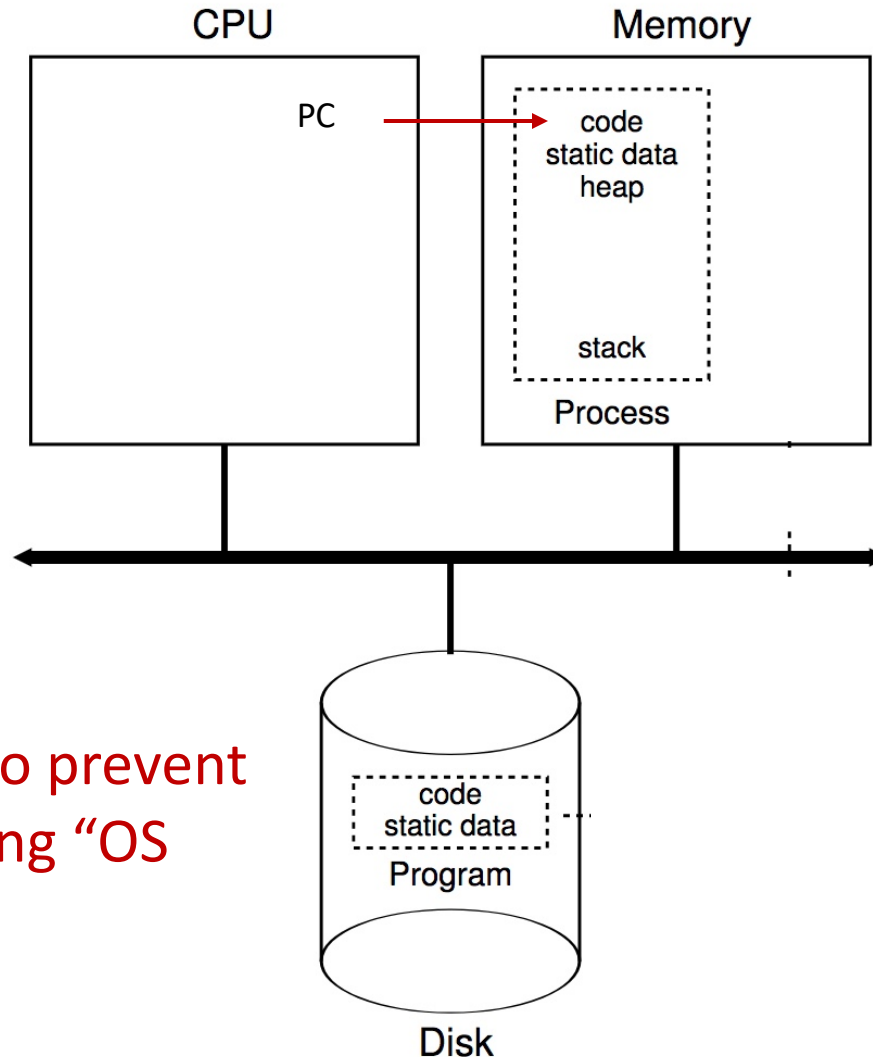Now, after process creation, CPU begins directly executing process code

# Process Creation



**Challenge**: how to prevent process from doing "OS kernel stuff"?

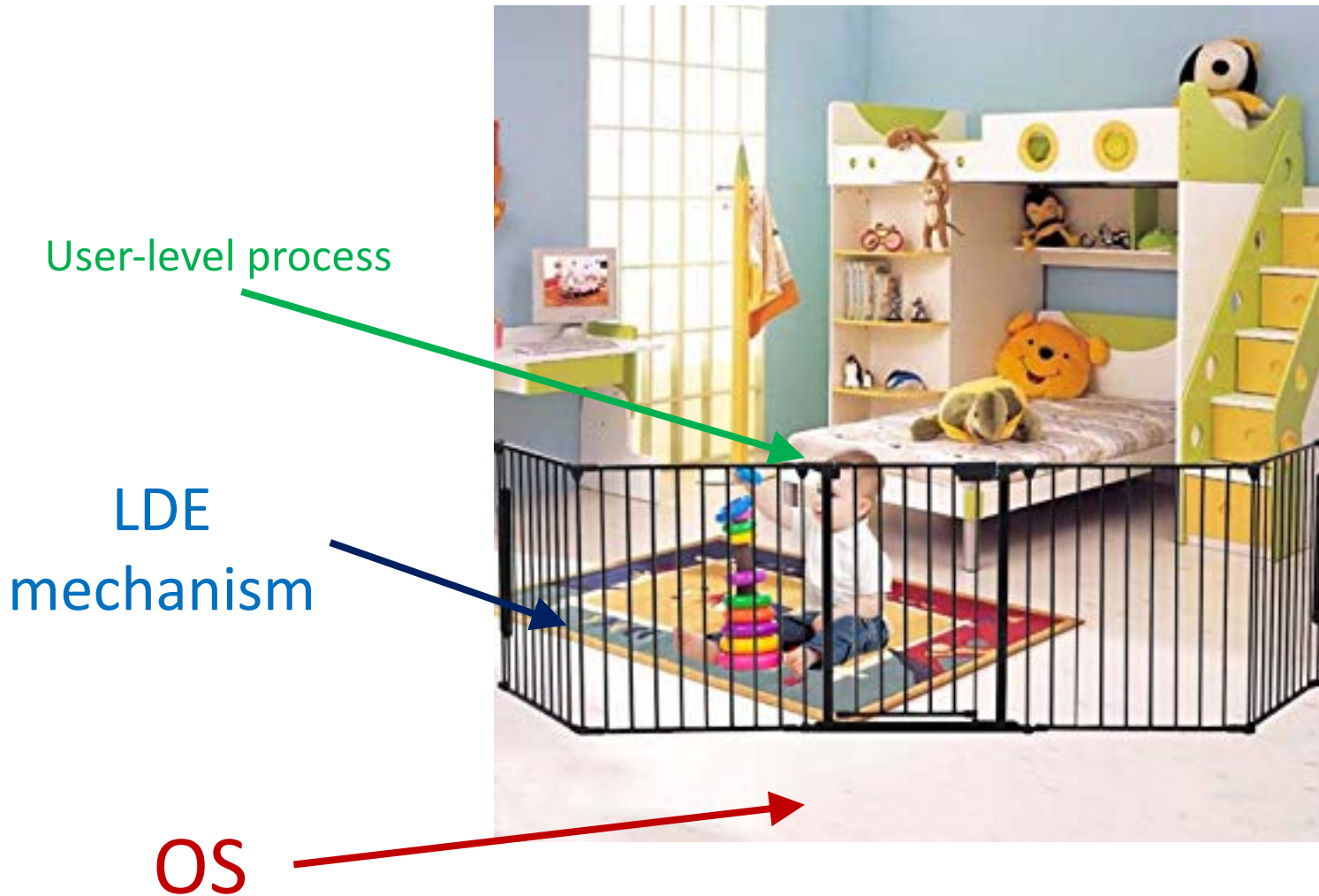# Limited Direct Execution (LDE)

# Limited Direct Execution (LDE)

- Low-level mechanism that implements the user-kernel space separation


- Usually let processes run with no OS involvement

- Limit what processes can do

- Offer privileged operations through well-defined channels with help of OS

# Limited Direct Execution (LDE)

# Limited Direct Execution (LDE)



User-level process

LDE
mechanism

OS

# What to limit?

- General memory access
- Disk I/O
- Certain x86 instructions

# How to limit?

- Need hardware support
- Add additional execution mode to CPU

- User mode: restricted, limited capabilities
- Kernel mode: privileged, not restricted

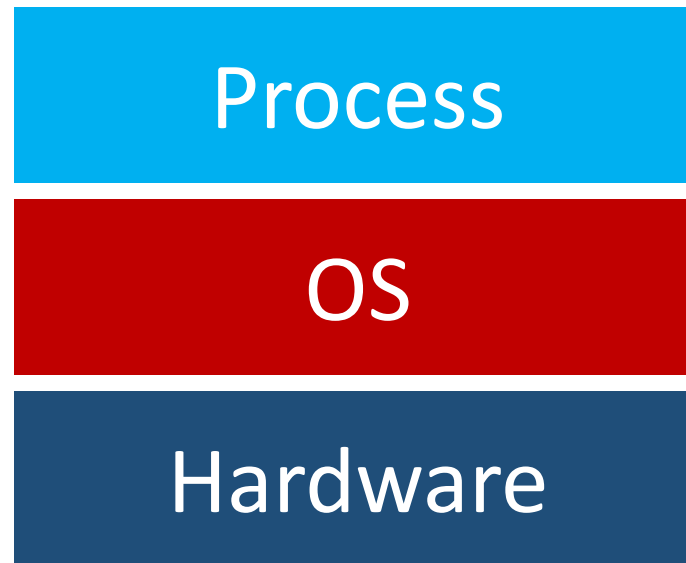- Processes start in user mode
- OS starts in kernel mode

# LDE: Remaining Challenges

1. What if process wants to do something privileged?

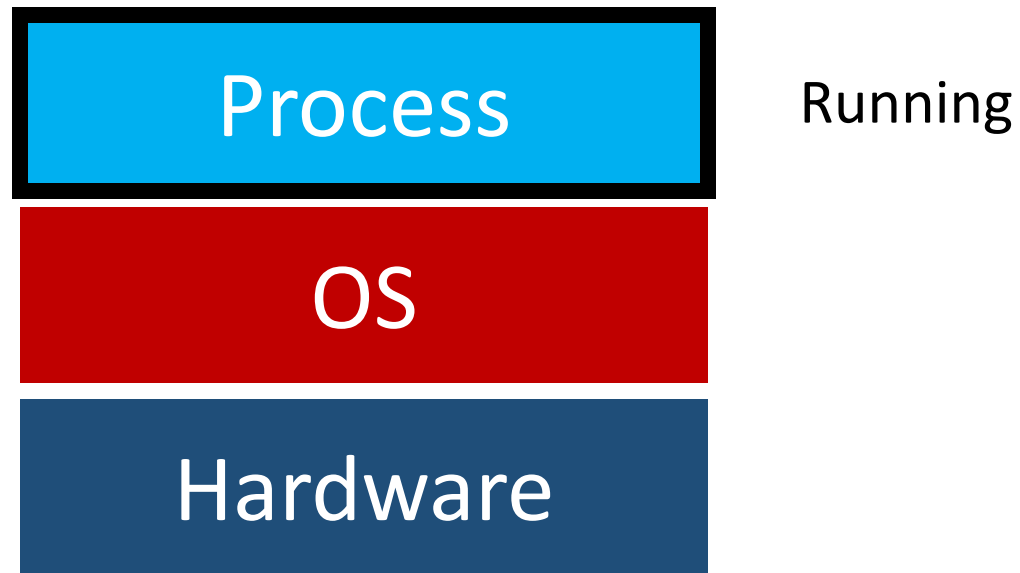2. How can OS switch processes (or do anything) if it's not running?

# LDE: Remaining Challenges

1. What if process wants to do something privileged?

2. How can OS switch processes (or do anything) if it's not running?
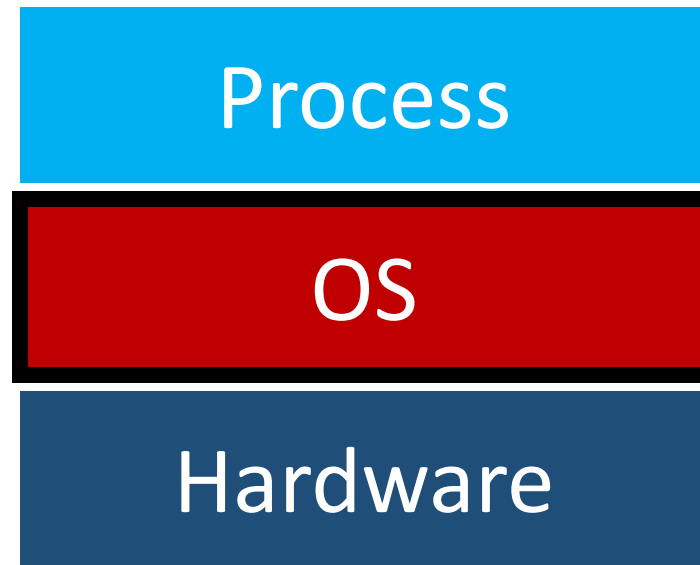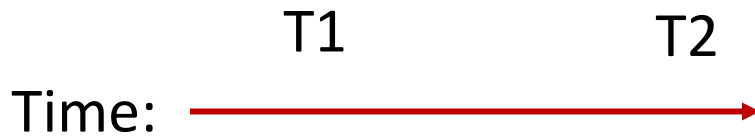
# Taking Turns

Process

OS

Hardware

# Taking Turns

Process

Running

OS

Hardware

T1

Time: →

# Taking Turns



Process

OS — Running

Hardware

T1          T2

Time: →

# Taking Turns

Process

Running

OS

Hardware

T1　　　　　T2　　　　　T3

Time:

# Taking Turns
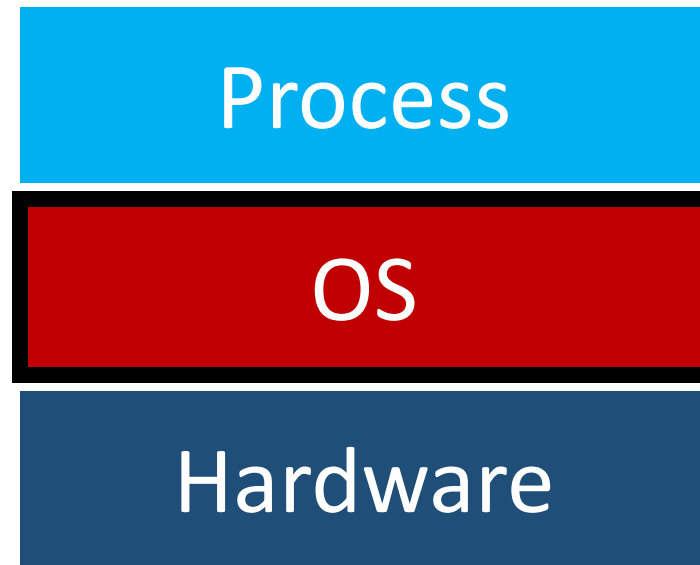
Process

OS                    Running
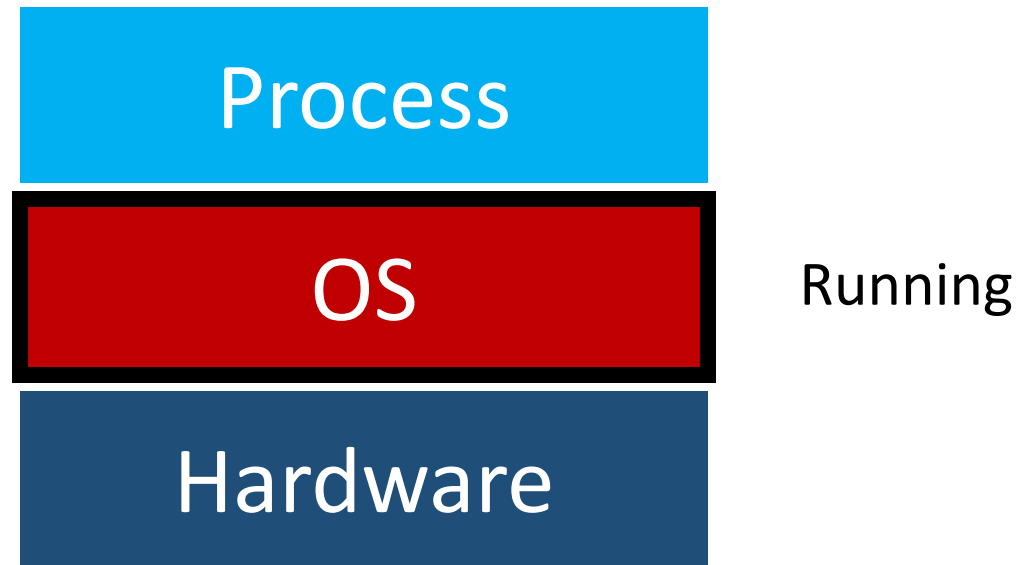
Hardware

T1            T2            T3            T4

Time:

# Taking Turns

**Question**: when/how do we switch to OS?

Process

OS

Hardware

Running

T1          T2          T3          T4

Time:

# Exceptions

# Interrupt

Process

OS

Hardware

# Interrupt

Process

OS

Hardware | key

# Interrupt

Process

OS    handler      Hardware interrupt

Hardware   key

# Interrupt

Process

OS

Hardware

# System Call

Process

OS

Hardware

# System Call

# System Call



Process | open

OS | handler — System call "trap"

Hardware

# System Call

Process

OS

Hardware

# Exception Handling

# Exception Handling: Implementation

- Goal: Processes and hardware should be able to call functions in the OS

- Corresponding OS functions should be:
  - At well-known locations
  - Safe from processes

Trap table

| | disk |
|---|---|
| | network |
| | timer |
| | keyboard |
| | system call |

Use array of function pointers to locate OS functions
(**Hardware knows where this is**)

Trap table

disk

network

**tick**

**timer**

keyboard

system call

Use array of function pointers to locate OS functions
(Hardware knows this through `lidt` instruction)

**Trap table**

disk

network

timer

keyboard

system call

How to handle variable number of system calls?

Trap table

disk

network

timer

keyboard

system call

syscall table

open

read

write

# Trap table



| |
|---|
| disk |
| network |
| timer |
| keyboard |
| **system call** |

**syscall**

# syscall table

| |
|---|
| open |
| read |
| write |

Trap table

disk

network

timer

keyboard

**system call**

**syscall**

syscall table

**open**

read

write

# Safe Transfers

- Only certain kernel functions should be callable
- Privileges should escalate at the moment of the call
  - Read/write disk
  - Kill processes
  - Access all memory
  - …

# LDE: Remaining Challenges

1.  What if process wants to do something privileged?

2.  How can OS switch processes (or do anything) if it's not running?

# Sharing (virtualizing) the CPU

# How does OS share...

- CPU?

- Memory?

- Disk?

# How does OS share...

- CPU? (a: time sharing)

- Memory? (a: space sharing)

- Disk? (a: space sharing)

# How does OS share...

• CPU? (a: time sharing)        **Today**

• Memory? (a: space sharing)

• Disk? (a: space sharing)

# How does OS share...

- CPU? (a: time sharing)    **Today**

- Memory? (a: space sharing)

- Disk? (a: space sharing)

**Goal:** processes should **not** know they are sharing **(each process will get its own virtual CPU)**

# What to do with processes that are not running?

- A: Store context in OS struct

# What to do with processes that are not running?

- A: Store context in OS struct

- Context:
    - CPU registers
    - Open file descriptors
    - State (sleeping, running, etc.)

# What to do with processes that are not running?

- A: Store context in OS struct

- Context:
  - CPU registers
  - Open file descriptors
  - State (sleeping, running, etc.)

# Process State Transitions

# Process State Transitions



Ready ──Scheduled──▶ Running
Running ──Descheduled──▶ Ready

I/O: done
I/O: initiate

Blocked

# Process State Transitions



View process state with "`ps xa`"

# How to transition? (mechanism)
# When to transition? (policy)

# Context Switch

- Problem: When to switch process contexts?

- Direct execution => OS can't run while process runs


- Can OS do anything while it's not running?

Y. Cheng

# Context Switch

- Problem: When to switch process contexts?
- Direct execution => OS can't run while process runs

- Can OS do anything while it's not running?
- A: it can't

# Context Switch

- Problem: When to switch process contexts?
- Direct execution => OS can't run while process runs


- Can OS do anything while it's not running?
- A: it can't


- Solution: Switch on interrupts
  - But what interrupt?

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call



P1

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call



P1

yield() call

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call

yield() call

OS

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call



OS

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call

yield() return

OS

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call



P2

yield() return

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call

P2

# Cooperative Approach

- Switch contexts for syscall interrupt
    - Special `yield()` system call

P2

yield() call

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call

yield() call
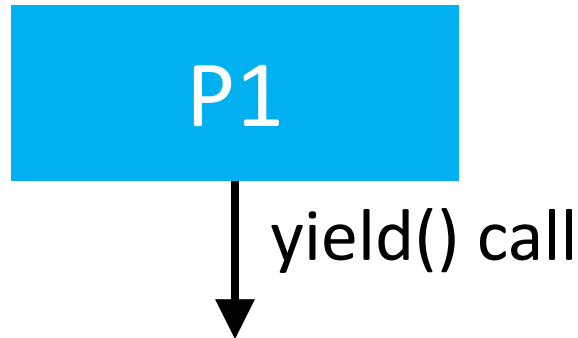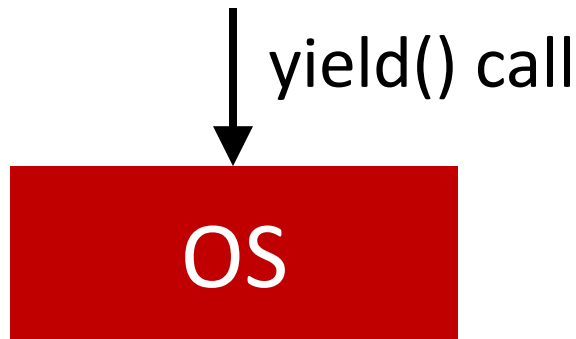
OS

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call

OS

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call

yield() return

OS

# Cooperative Approach

- Switch contexts for syscall interrupt
    - Special `yield()` system call

P1

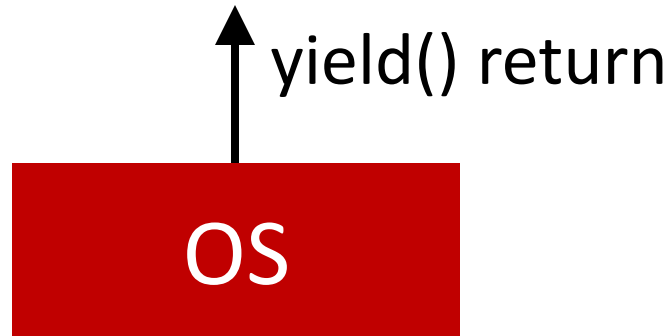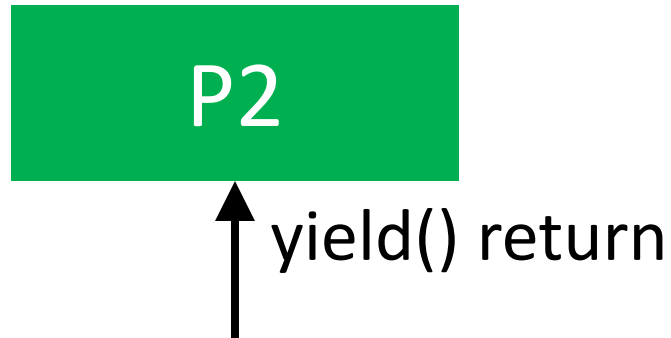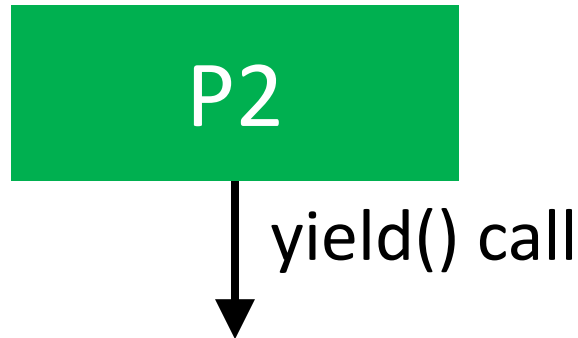yield() return

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call

P1

# Cooperative Approach

- Switch contexts for syscall interrupt
    - Special `yield()` system call
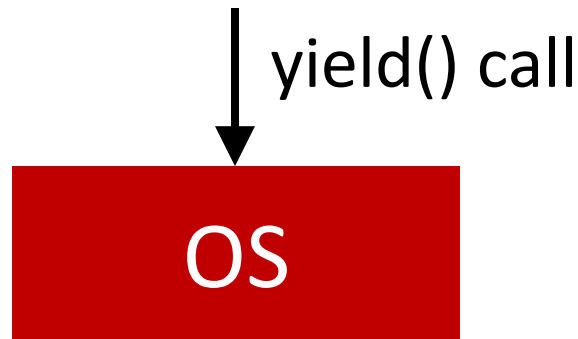
**P1**

**Critiques?**

# Cooperative Approach

- Switch contexts for syscall interrupt
  - Special `yield()` system call
- Cooperative approach is a passive approach

P1

**Critiques?**

**What if P1 never calls** `yield()`**?**

# Non-Cooperative Approach

- Switch contexts on timer (hardware) interrupt

- Set up before running any processes

- Hardware does not let processes prevent this
  - Hardware/OS enforces process preemption

# Non-Cooperative Approach

| OS @ run (kernel mode) | Hardware | Program (user mode) |
|---|---|---|
| | | Process A |
| | | ... |

# Non-Cooperative Approach

| OS @ run (kernel mode) | Hardware | Program (user mode) |
|---|---|---|
| | | Process A |
| | | ... |
| | timer interrupt | |
| | save regs(A) to k-stack(A) | |
| | move to kernel mode | |
| | jump to trap handler | |

# Non-Cooperative Approach

| OS @ run (kernel mode) | Hardware | Program (user mode) |
|---|---|---|
| | | Process A |
| | | ... |
| | **timer interrupt** | |
| | save regs(A) to k-stack(A) | |
| | move to kernel mode | |
| | jump to trap handler | |
| Handle the trap | | |
| Call `switch()` routine | | |
|   save regs(A) to proc-struct(A) | | |
|   restore regs(B) from proc-struct(B) | | |
|   switch to k-stack(B) | | |
| **return-from-trap (into B)** | | |

# Non-Cooperative Approach

| OS @ run (kernel mode) | Hardware | Program (user mode) |
|---|---|---|
| | | Process A |
| | | ... |
| | **timer interrupt** | |
| | save regs(A) to k-stack(A) | |
| | move to kernel mode | |
| | jump to trap handler | |
| Handle the trap | | |
| Call `switch()` routine | | |
|   save regs(A) to proc-struct(A) | | |
|   restore regs(B) from proc-struct(B) | | |
|   switch to k-stack(B) | | |
| **return-from-trap (into B)** | | |
| | restore regs(B) from k-stack(B) | |
| | move to user mode | |
| | jump to B's PC | |

# Non-Cooperative Approach

| OS @ run (kernel mode) | Hardware | Program (user mode) |
|---|---|---|
| | | Process A |
| | | ... |
| | **timer interrupt** | |
| | save regs(A) to k-stack(A) | |
| | move to kernel mode | |
| | jump to trap handler | |
| Handle the trap | | |
| Call `switch()` routine | | |
|   save regs(A) to proc-struct(A) | | |
|   restore regs(B) from proc-struct(B) | | |
|   switch to k-stack(B) | | |
| **return-from-trap (into B)** | | |
| | restore regs(B) from k-stack(B) | |
| | move to user mode | |
| | jump to B's PC | |
| | | Process B |
| | | ... |

# Preemptive Approach

P1

# Preemptive Approach

P1

tick

# Preemptive Approach
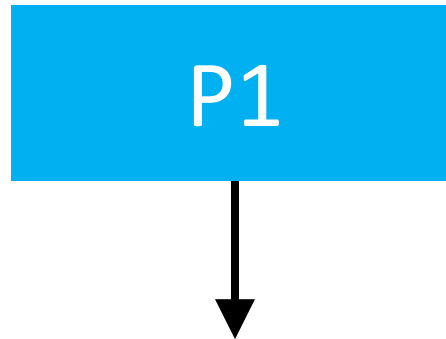
OS     tick

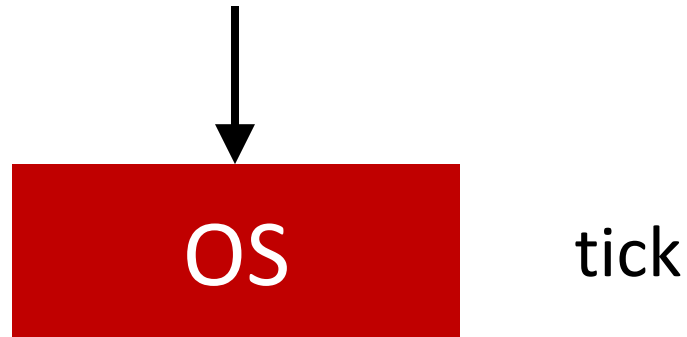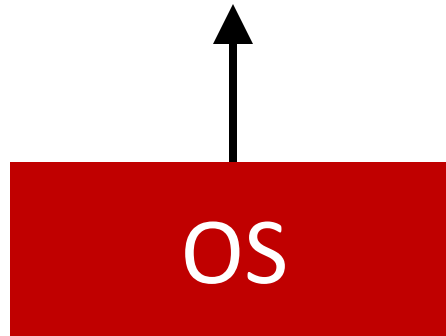# Preemptive Approach

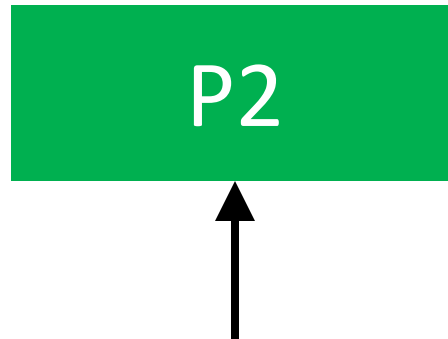<div style="text-align:center; background-color:#CC0000; color:white;">OS</div>

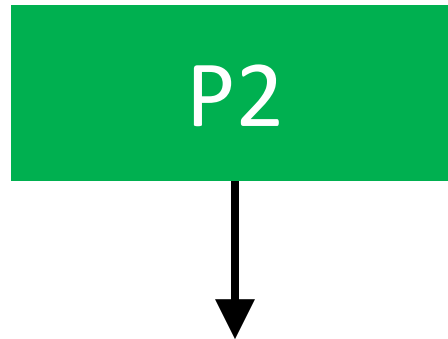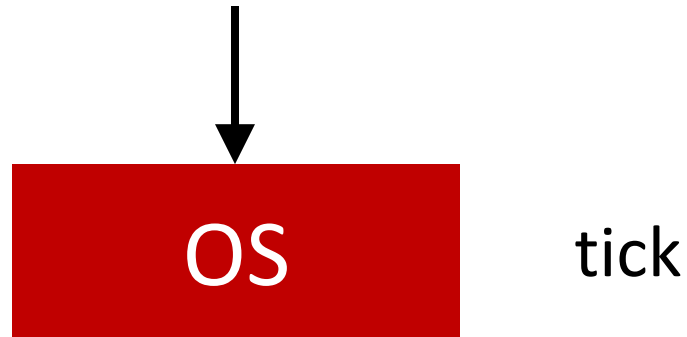# Preemptive Approach

# Preemptive Approach

# Preemptive Approach
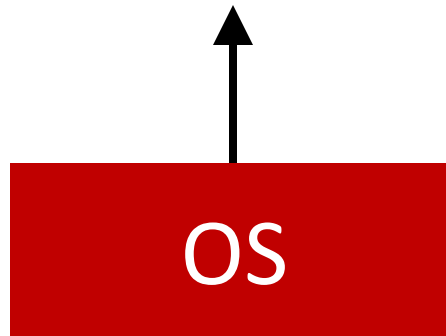
P2

# Preemptive Approach

P2

tick

# Preemptive Approach



OS       tick

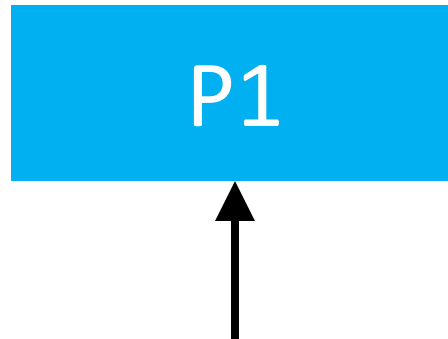# Preemptive Approach

OS

# Preemptive Approach

OS

# Preemptive Approach

# Preemptive Approach

P1

# LDE Summary

- Smooth <span style="color:red">context switching</span> makes each process think it has its own CPU (virtualization!)

- <span style="color:blue">Limited direct execution</span> makes processes fast

- Hardware provides a lot of OS support
  - Limited direct execution
  - Timer interrupt
  - Automatic register saving