

LDE and Threads

CS 571: Operating Systems (Spring 2020)

Lecture 2

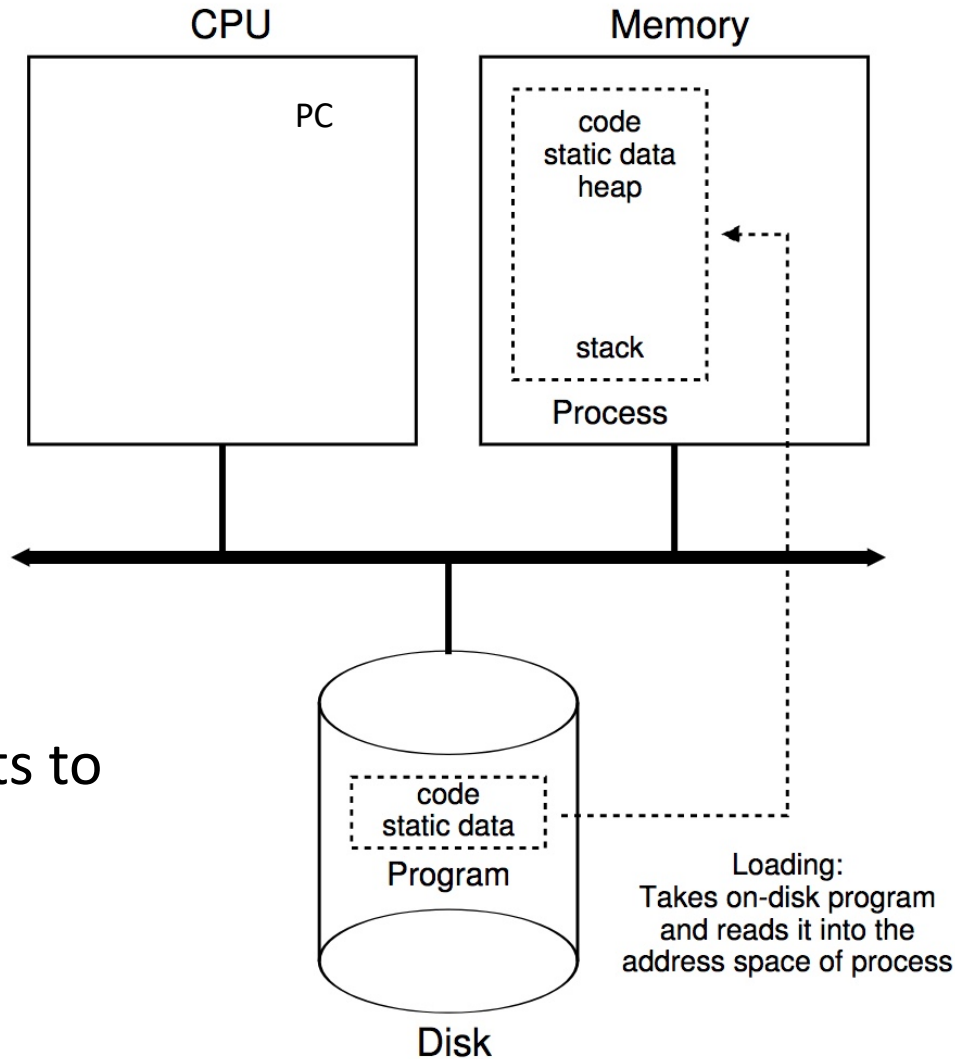
Yue Cheng

Some material taken/derived from:

- Wisconsin CS-537 materials created by Remzi Arpaci-Dusseau.

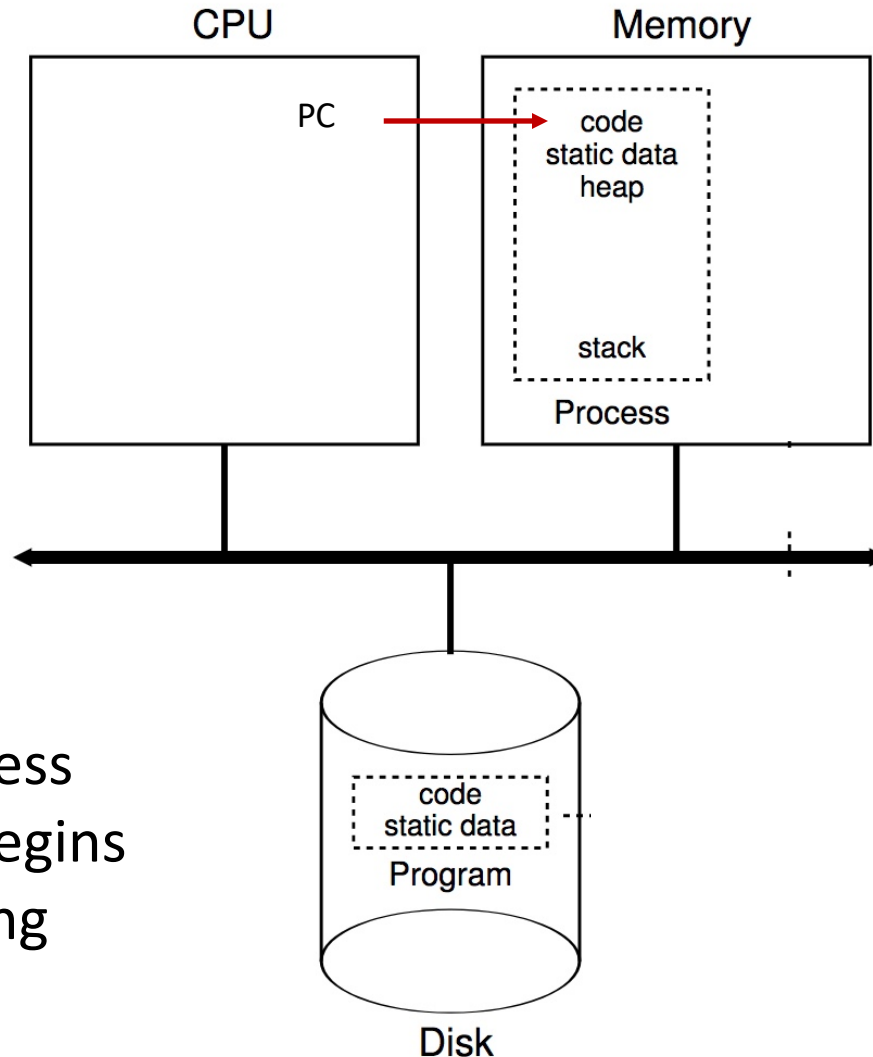
Licensed for use under a Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License.

Process Creation



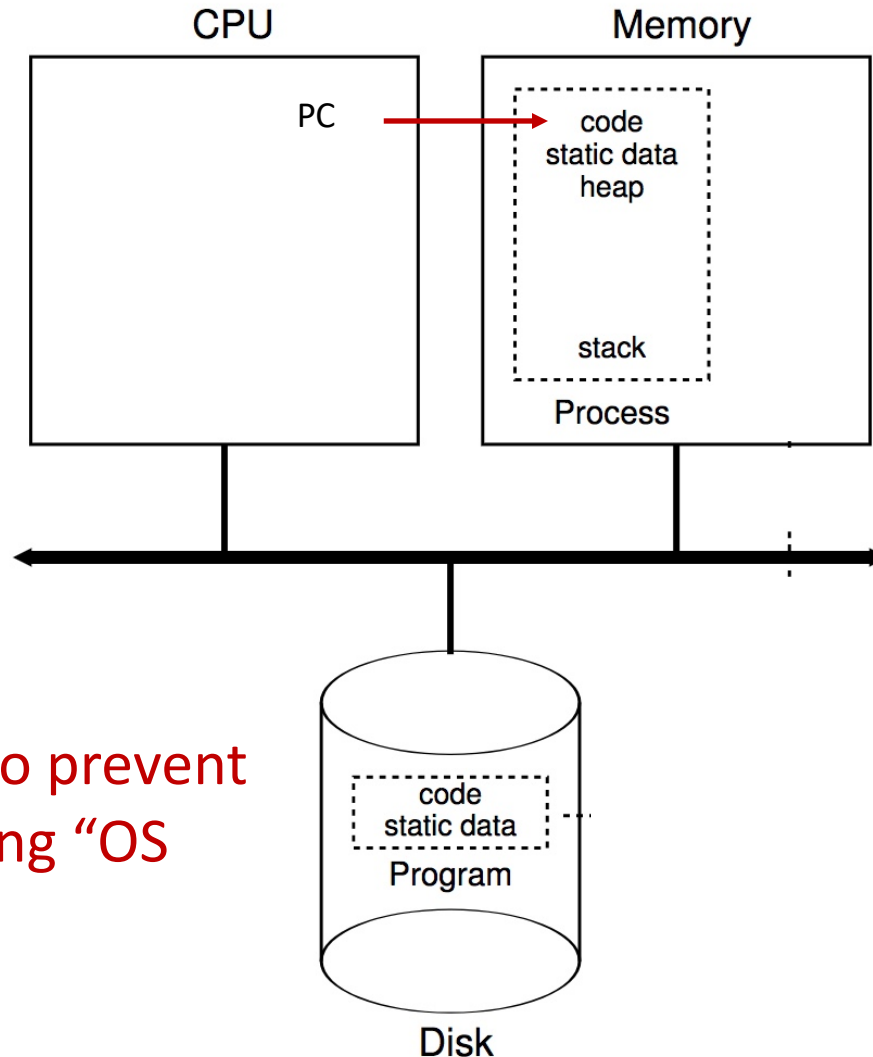
Before, PC points to kernel code

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)

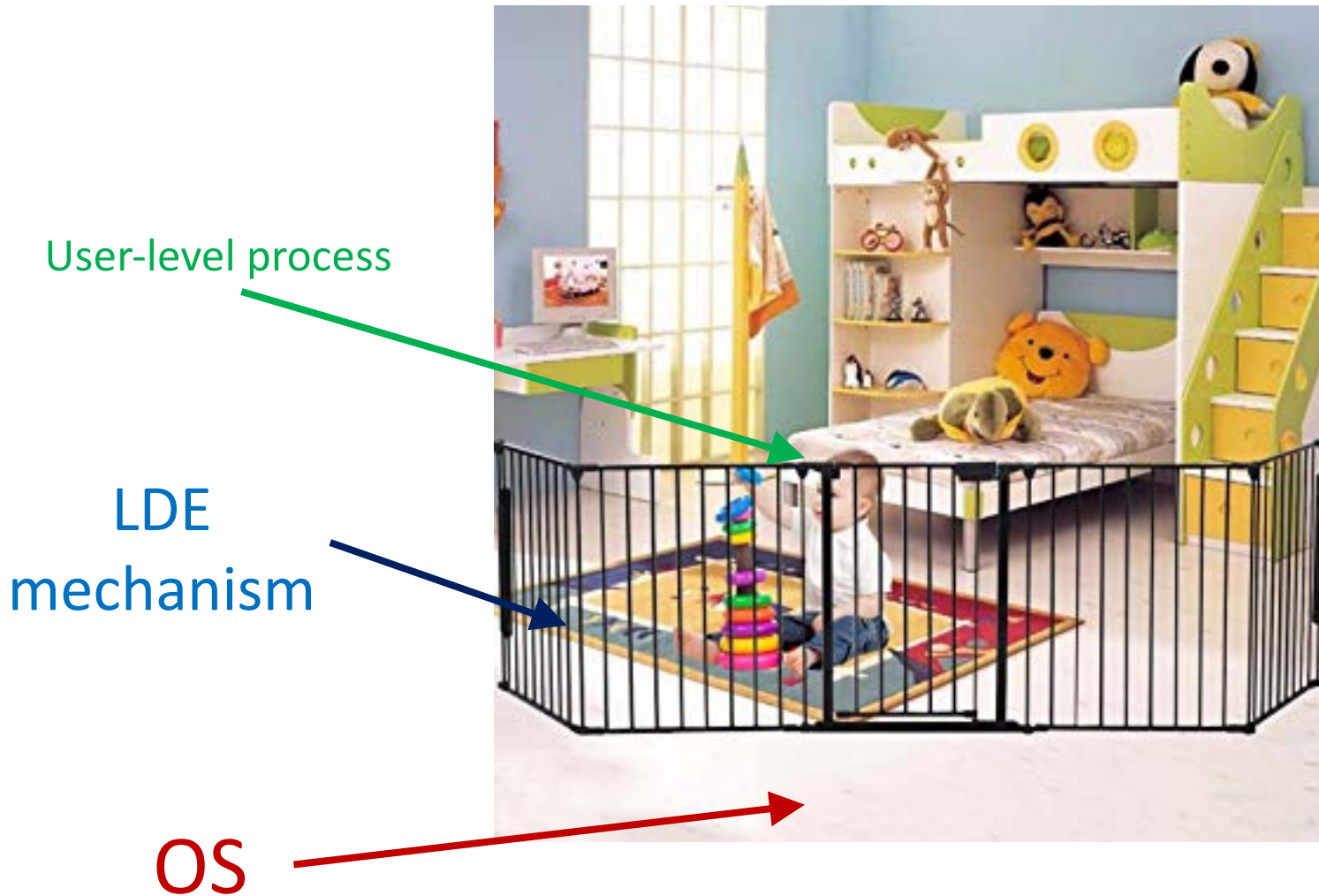
Limited Direct Execution (LDE)

- 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)



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**

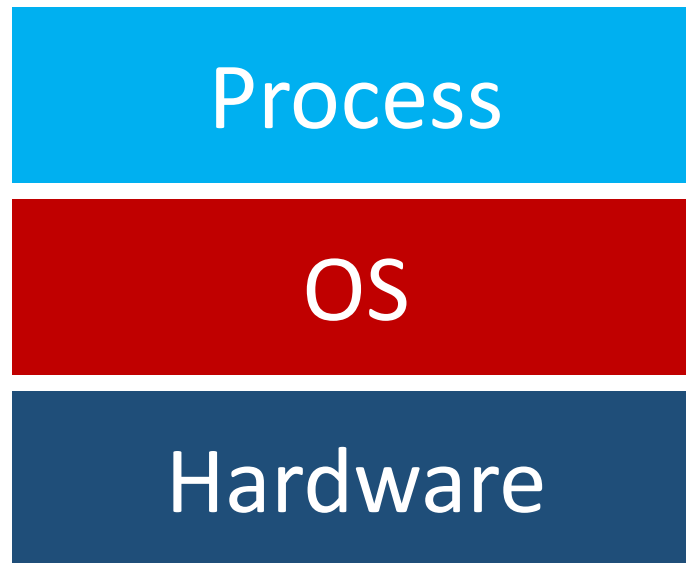
LDE: Remaining Challenges

1. What if process wants to do something privileged?
2. How can OS switch processes (or do anything) if it's not running?

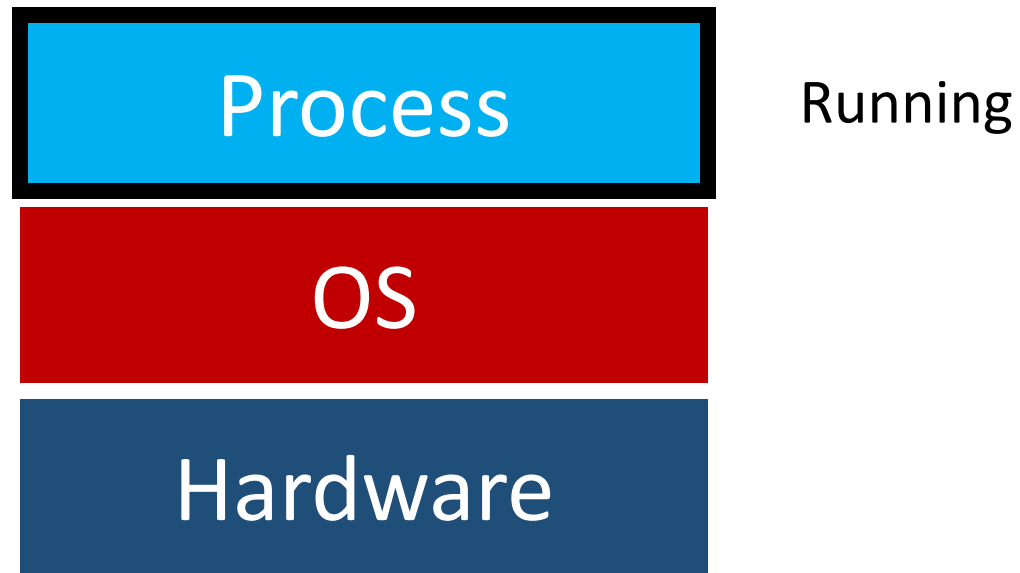
LDE: Remaining Challenges

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

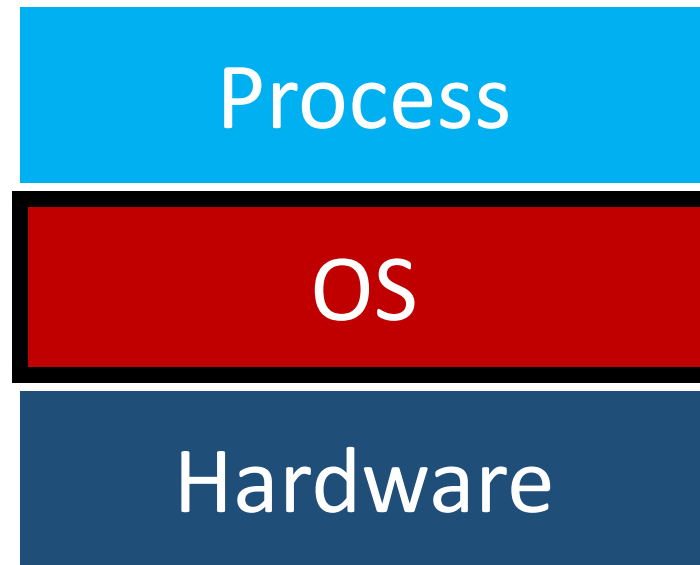


Taking Turns

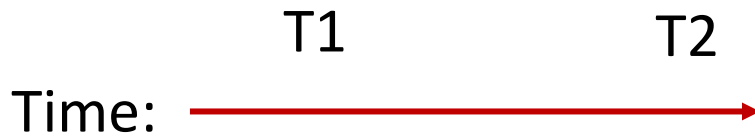


Time: $\xrightarrow{T1}$

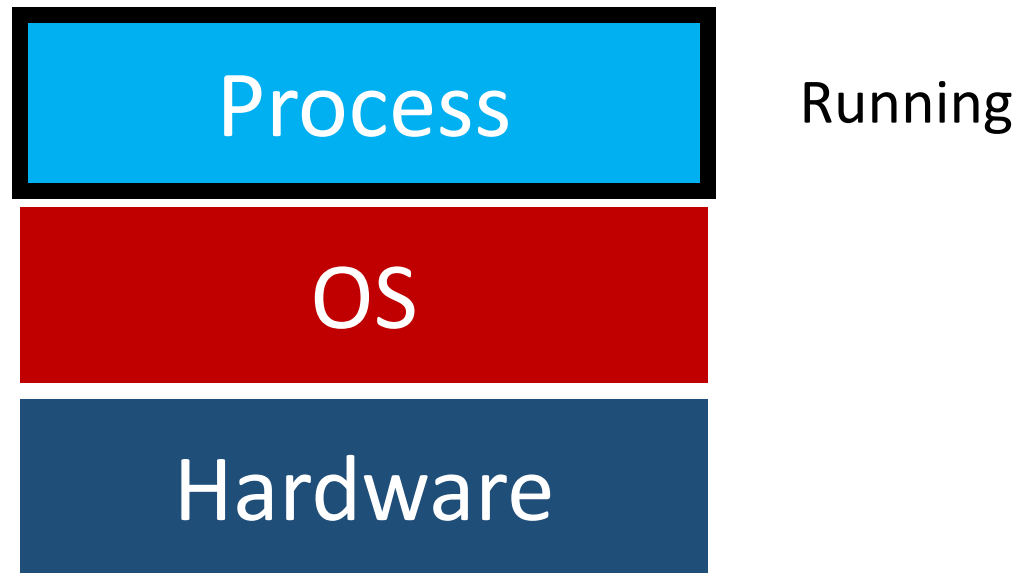
Taking Turns



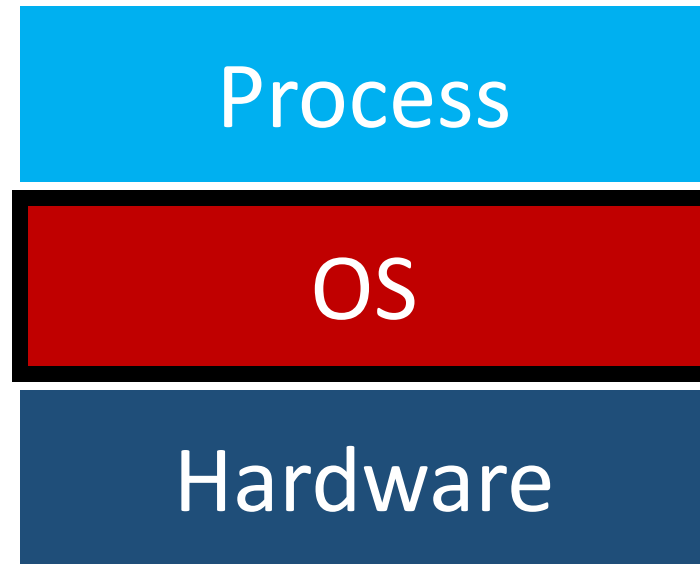
Running



Taking Turns



Taking Turns

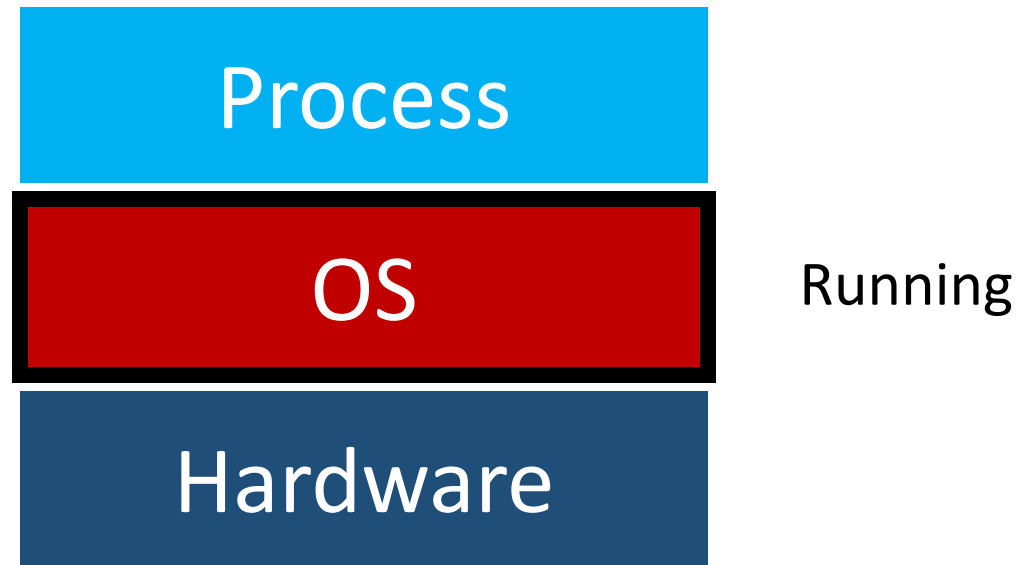


Running



Taking Turns

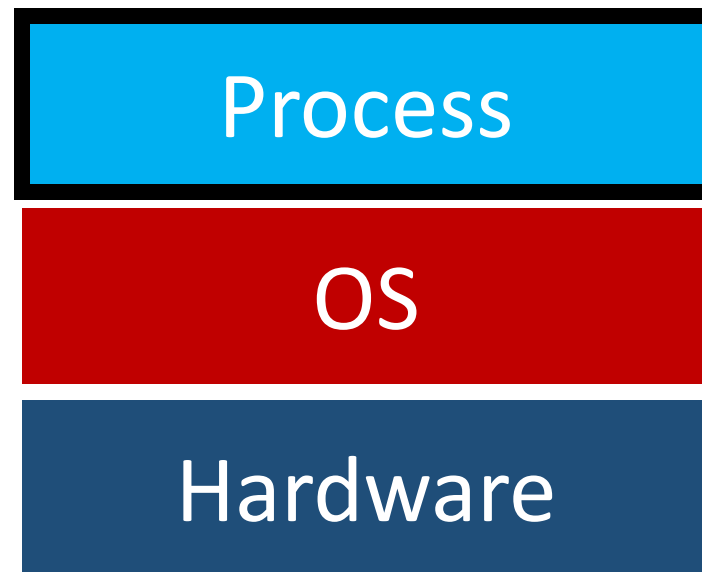
Question: when/how do we switch to OS?



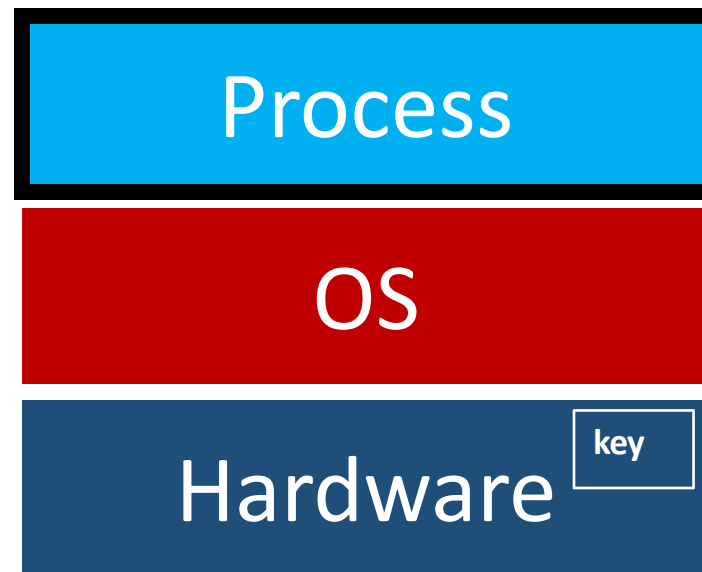
Time: 

Exceptions

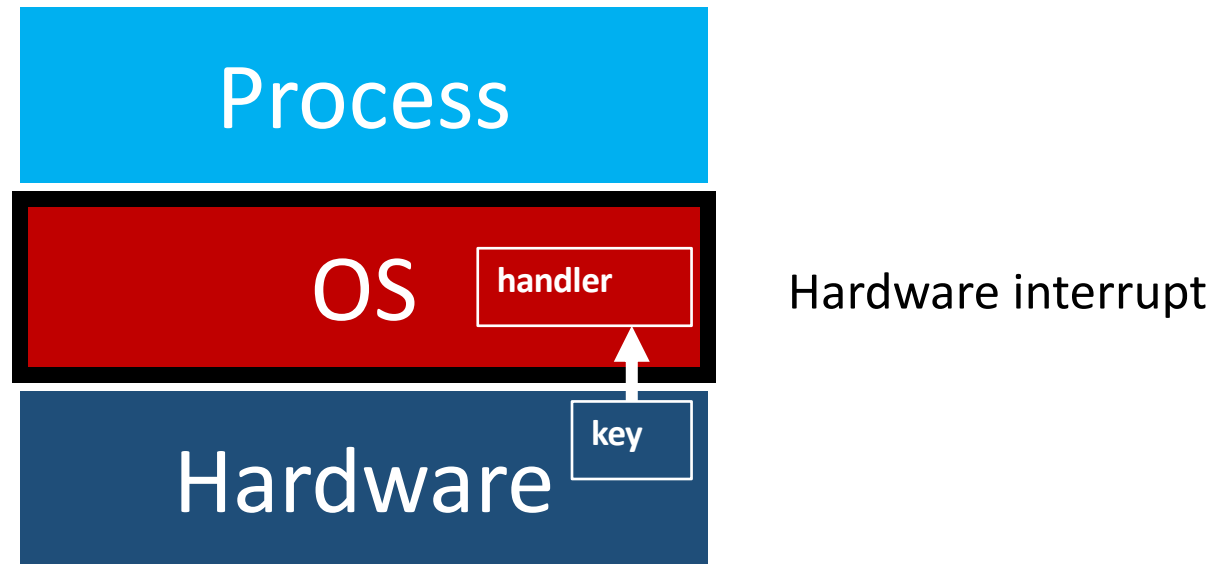
Interrupt



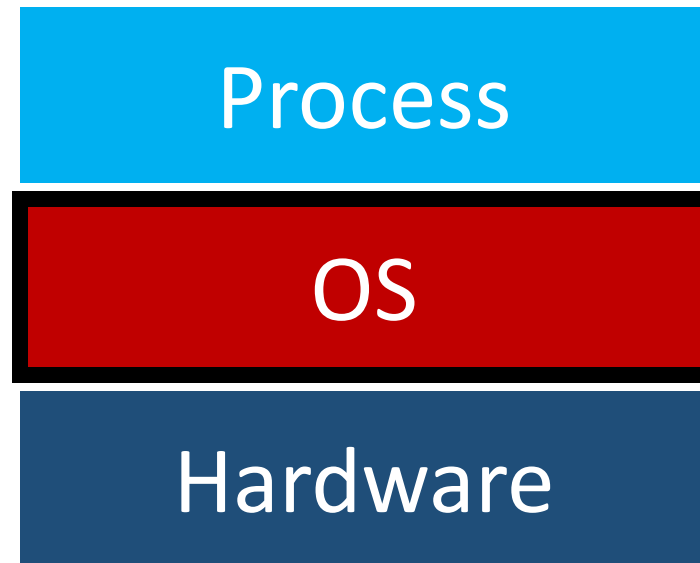
Interrupt



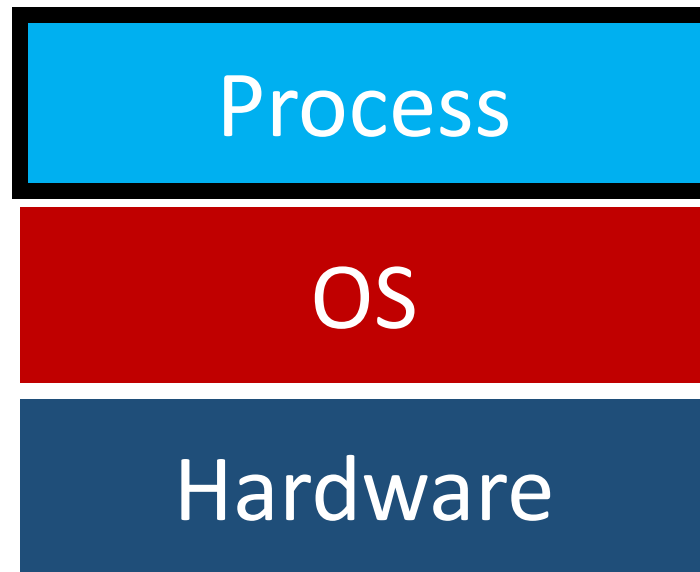
Interrupt



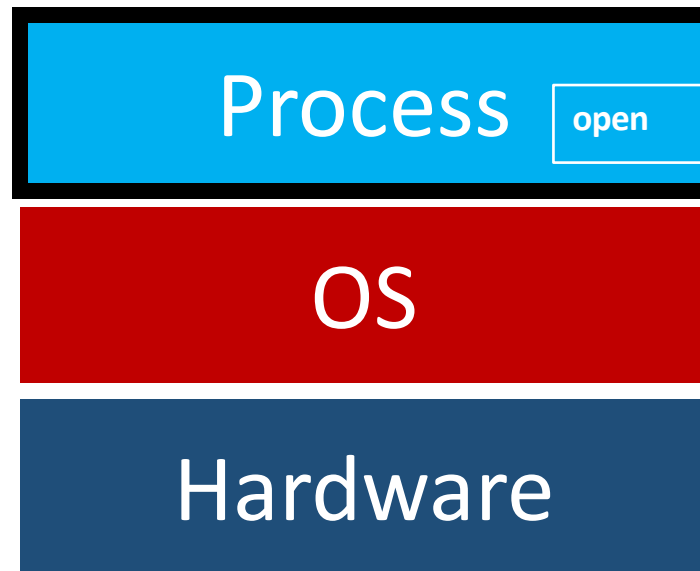
Interrupt



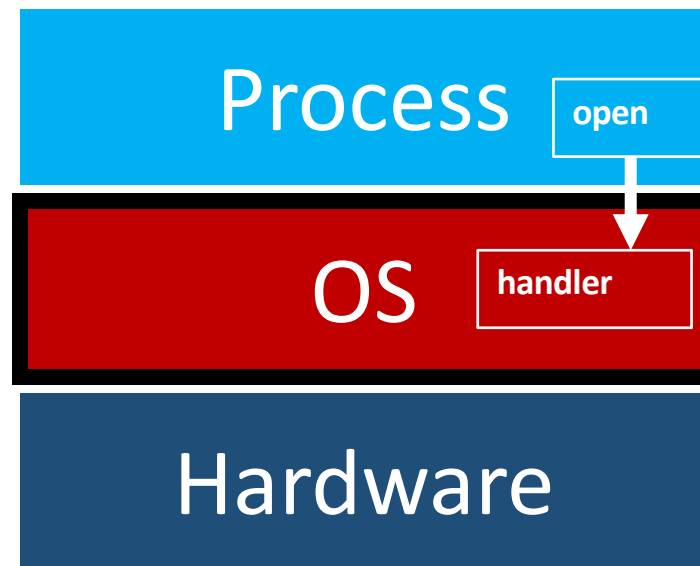
System Call



System Call

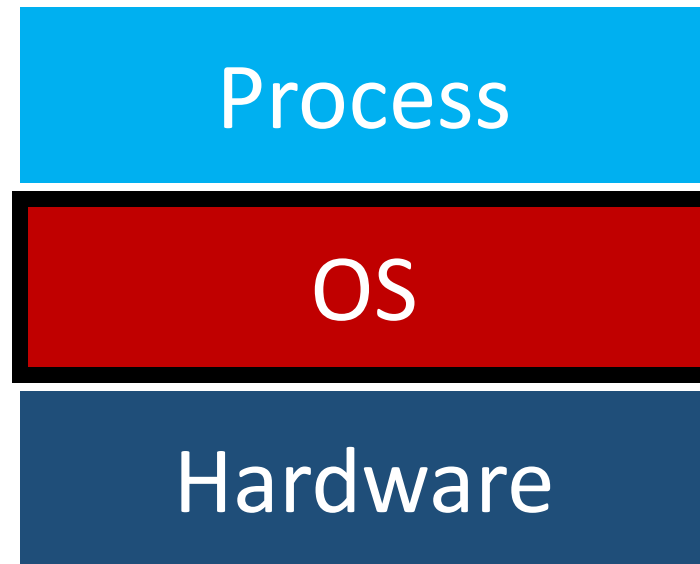


System Call



System call “trap”

System Call

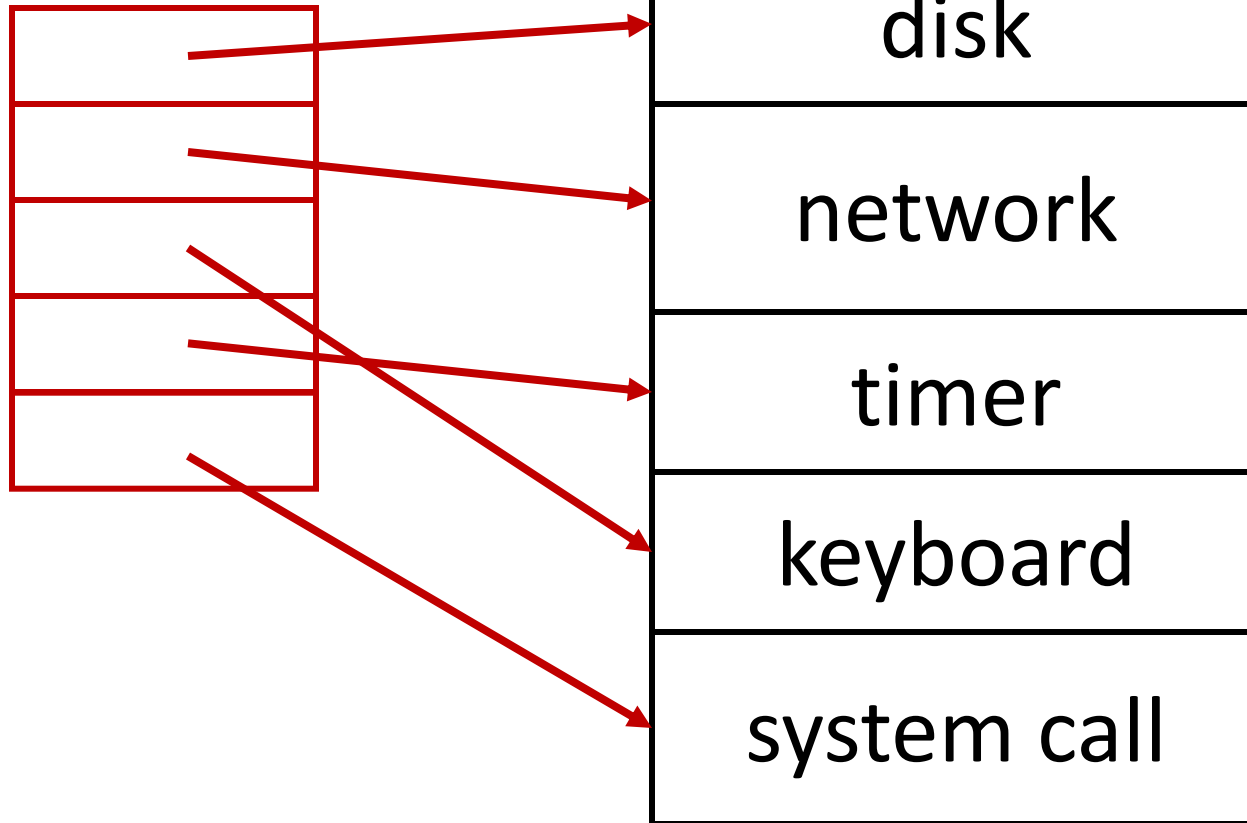


Exception Handling

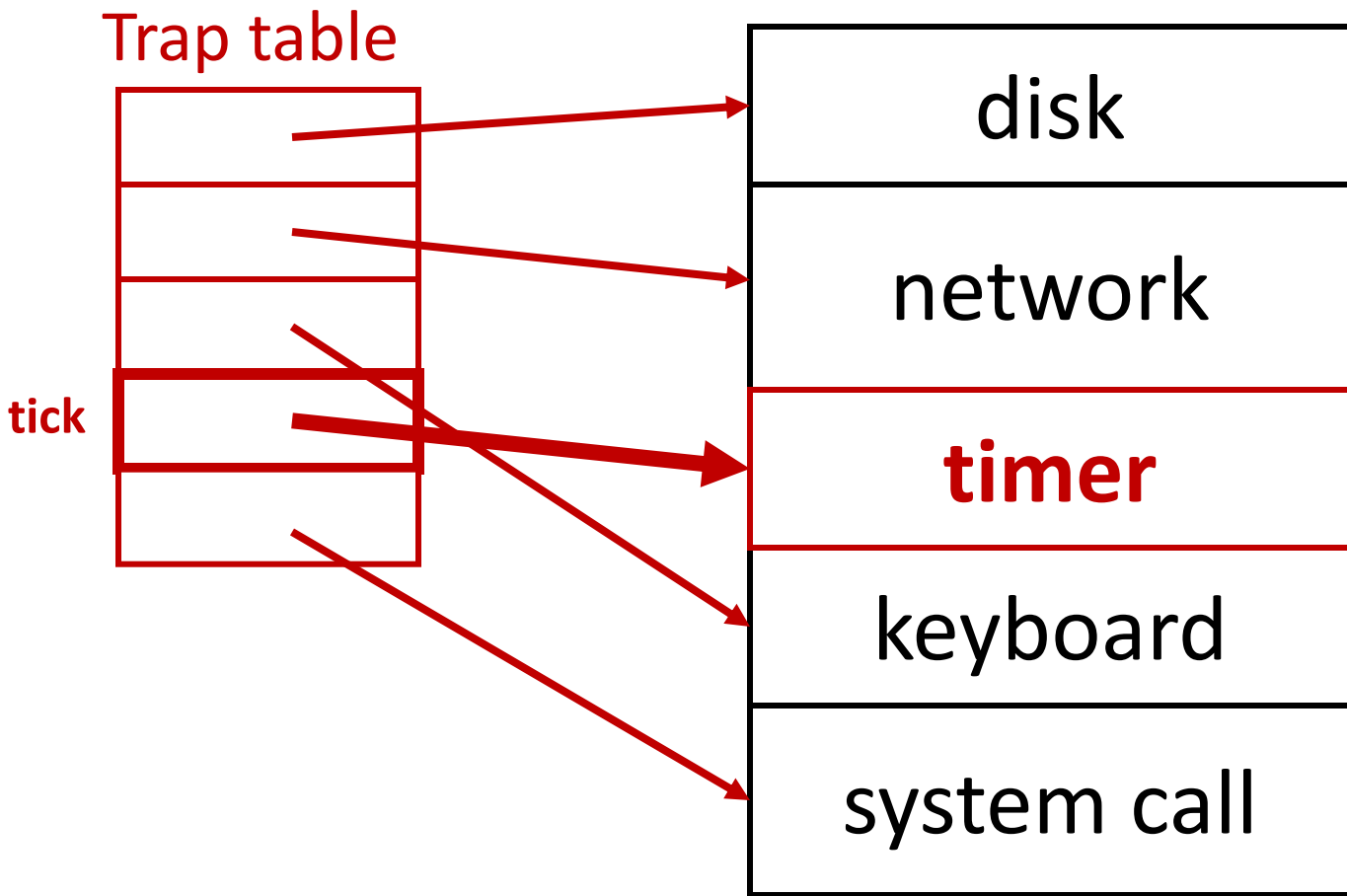
Exception Handling: Implementation

- 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

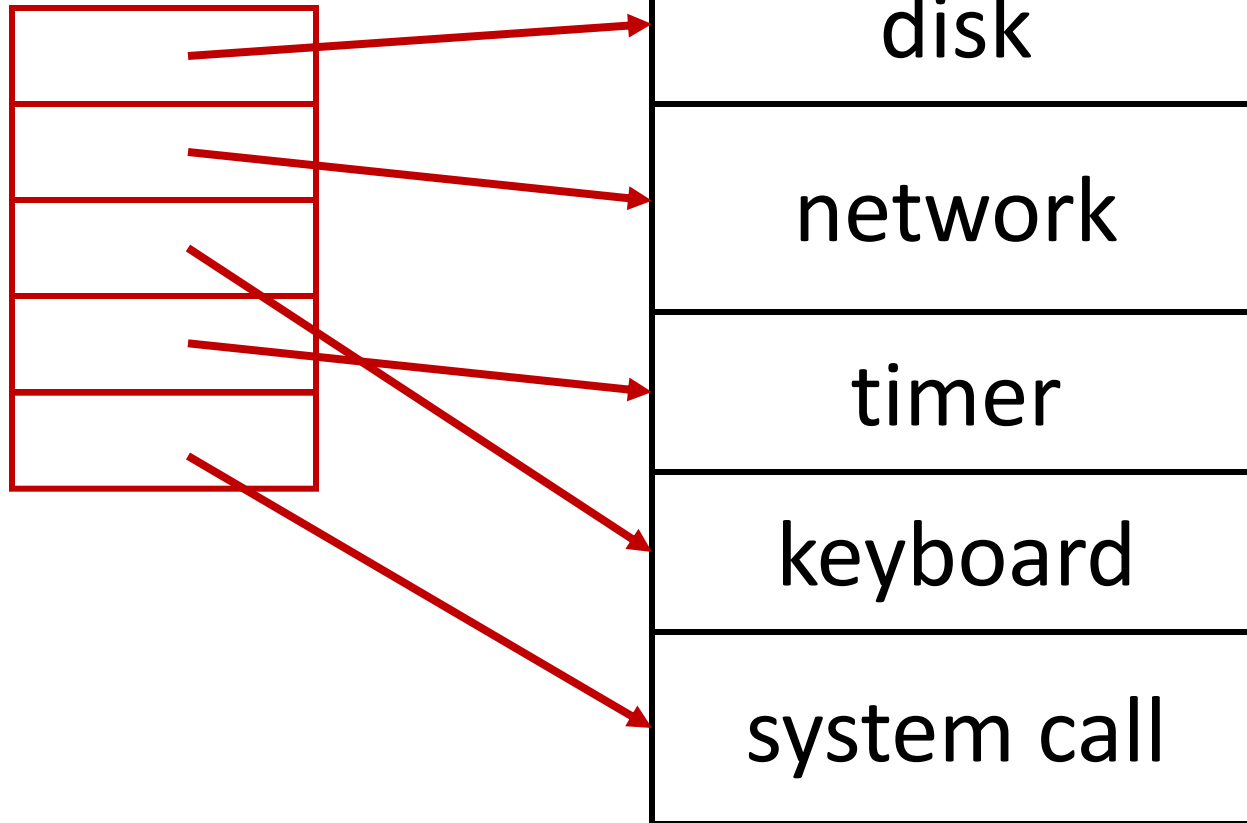


Use array of function pointers to locate OS functions
(Hardware knows where this is)



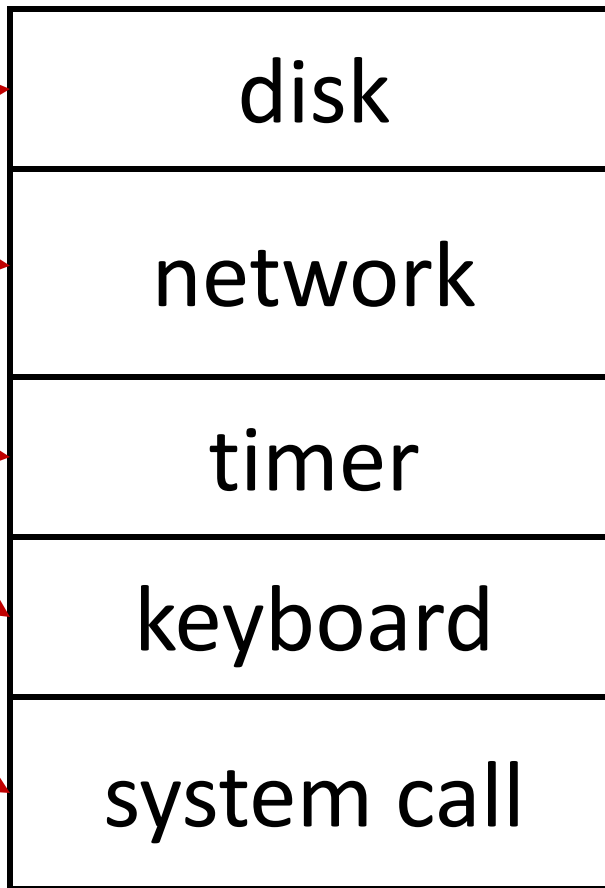
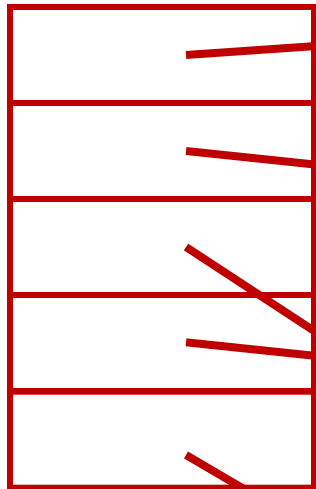
Use array of function pointers to locate OS functions
(Hardware knows this through **lidt** instruction)

Trap table

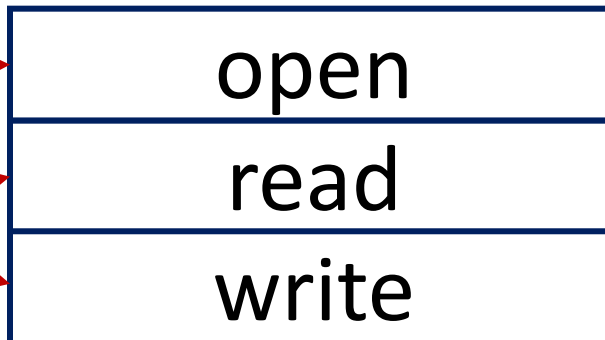
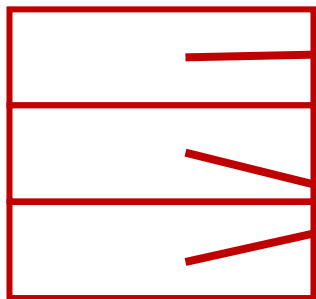


How to handle variable number of system calls?

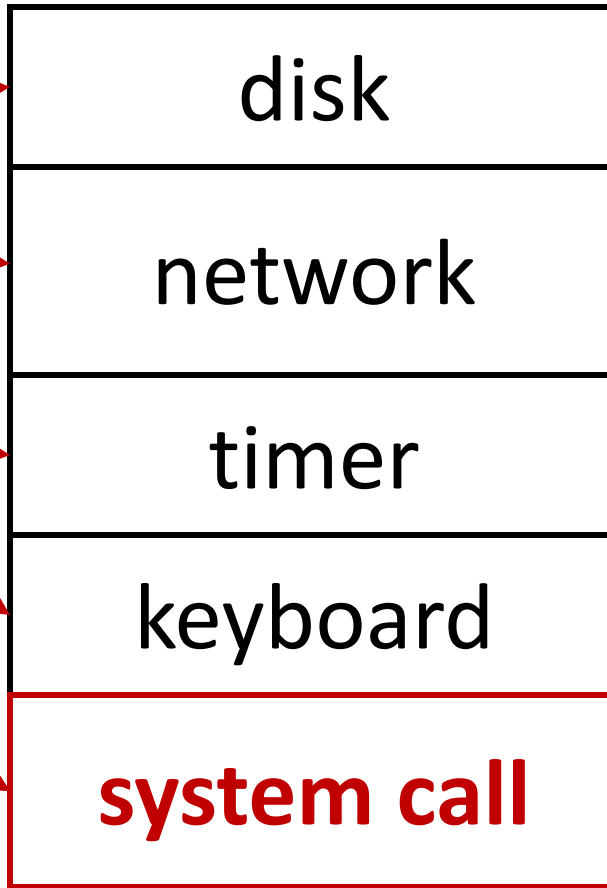
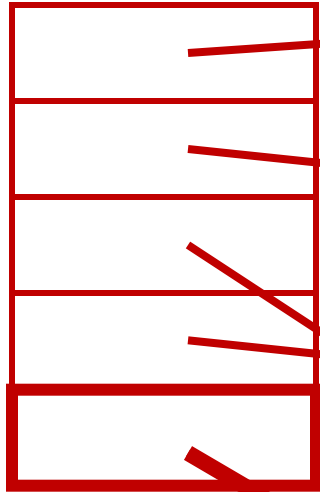
Trap table



syscall table

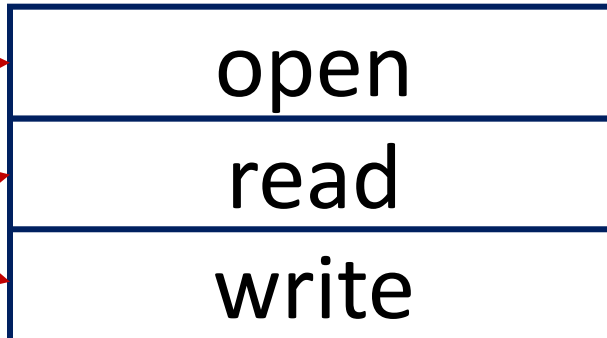
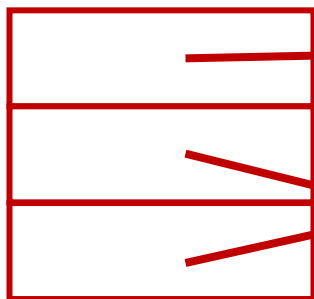


Trap table

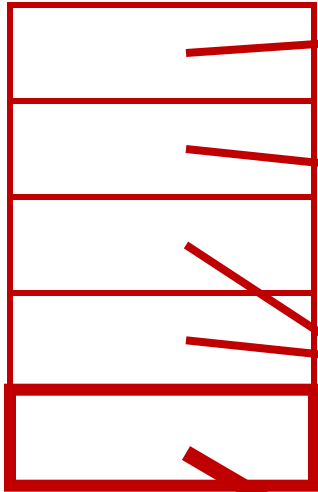


syscall

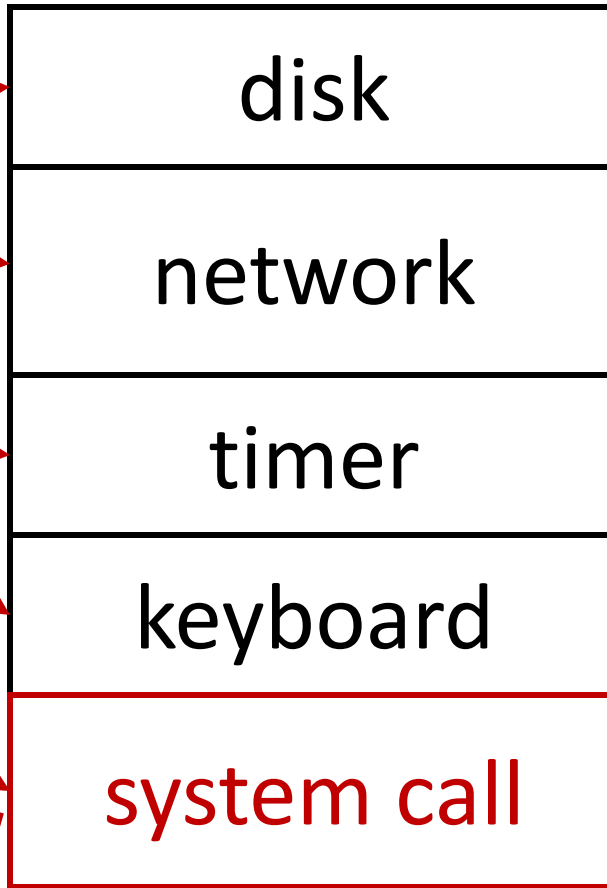
syscall table



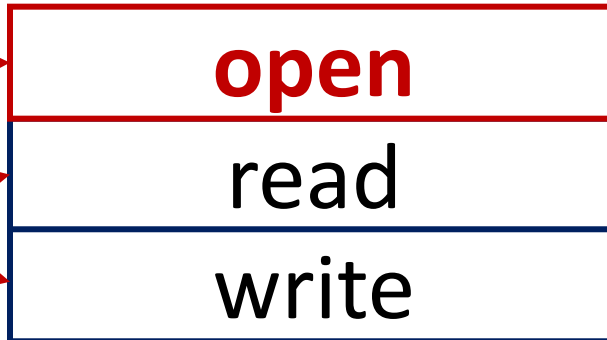
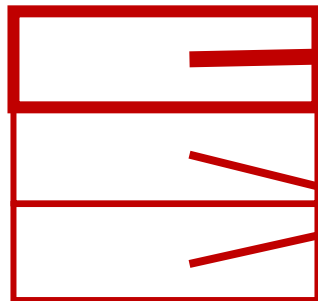
Trap table



syscall



syscall table



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**)
- Memory? (a: space sharing)
- Disk? (a: space sharing)

Today

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

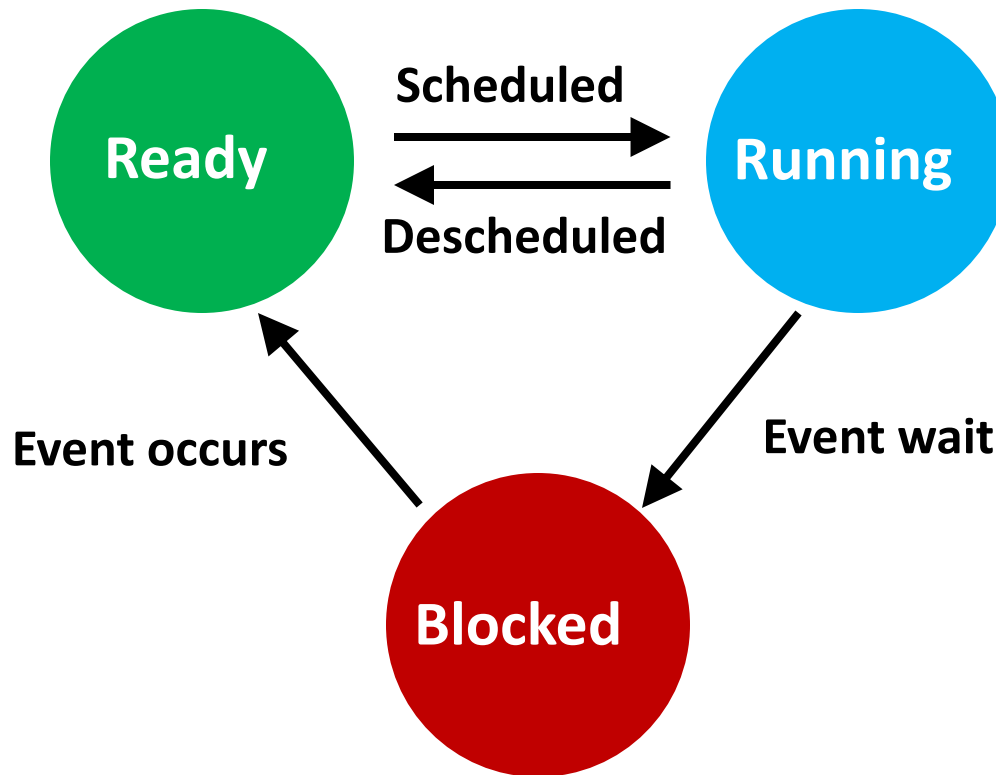
What to do with processes that are not running?

- A: Store context in OS struct
- Context:
 - CPU registers
 - Open file descriptors
 - State (sleeping, running, etc.)

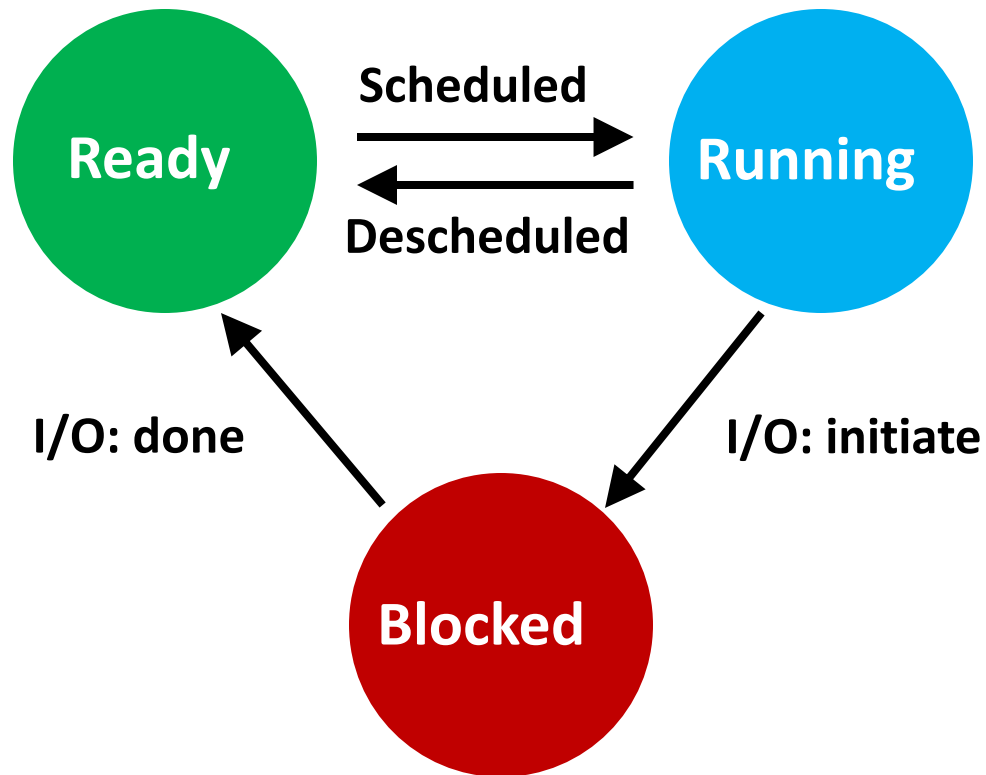
What to do with processes that are not running?

- A: Store context in OS struct
- Context:
 - CPU registers
 - Open file descriptors
 - **State** (sleeping, running, etc.)

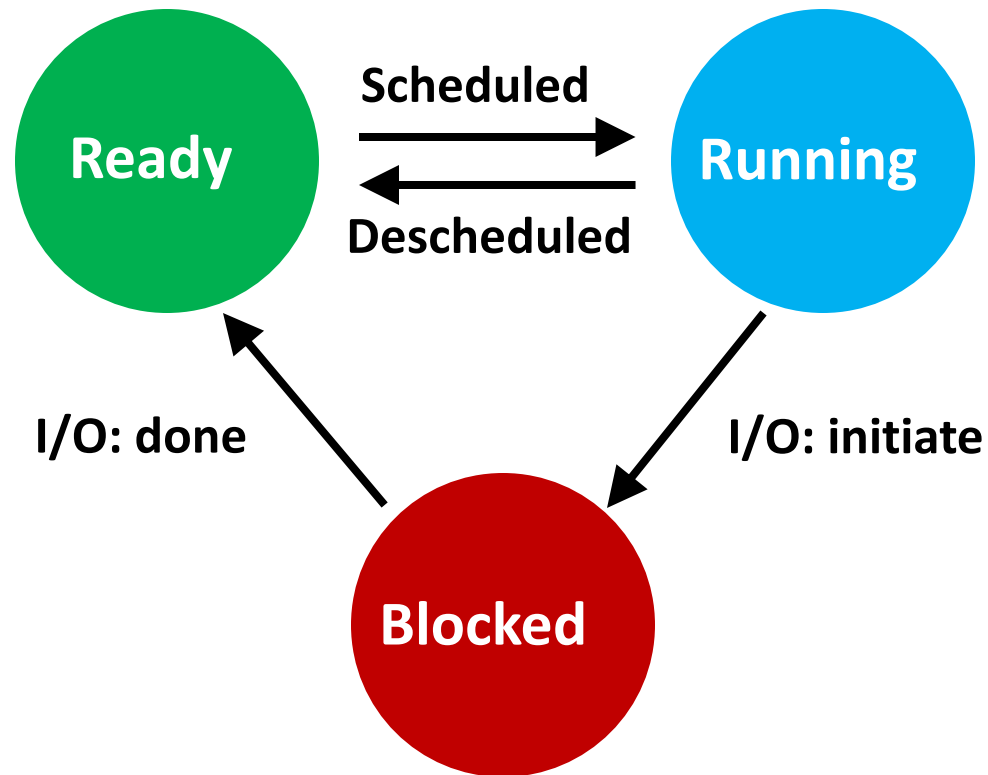
Process State Transitions



Process State Transitions



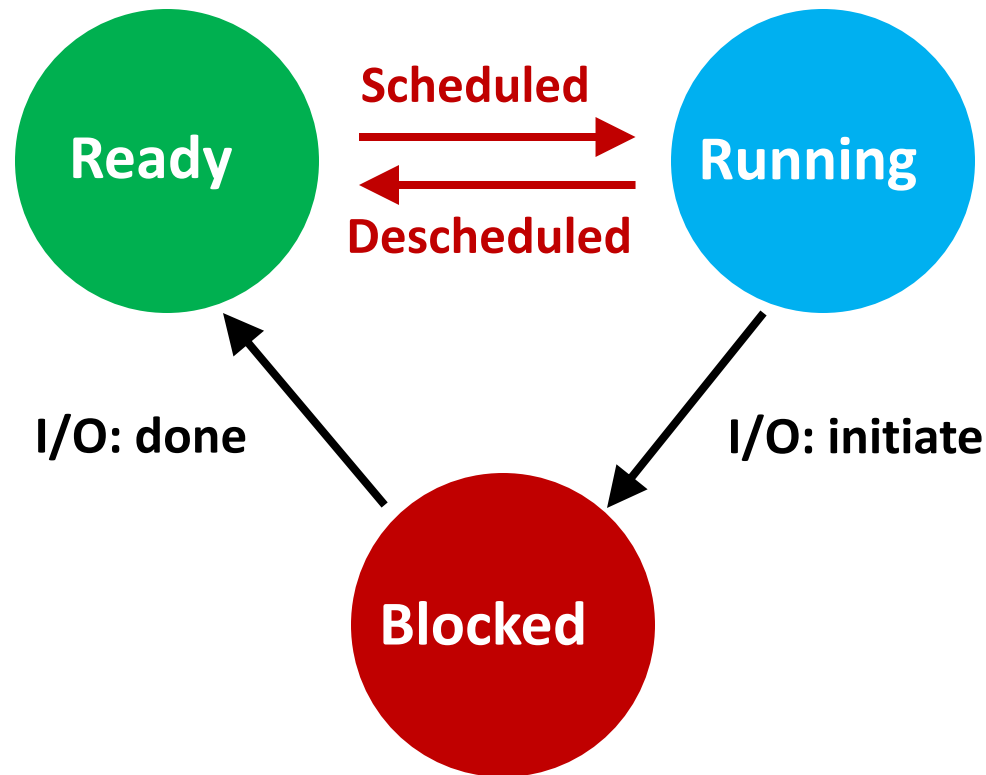
Process State Transitions



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?

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

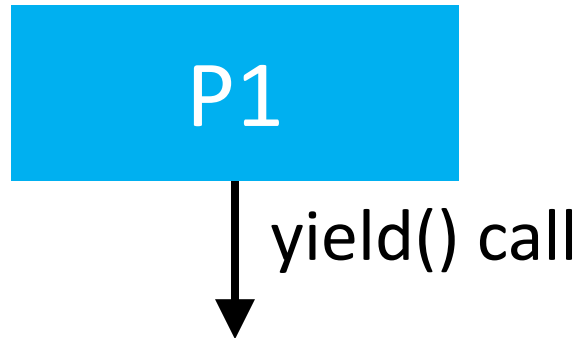
- Switch contexts for syscall interrupt
 - Special `yield()` system call



P1

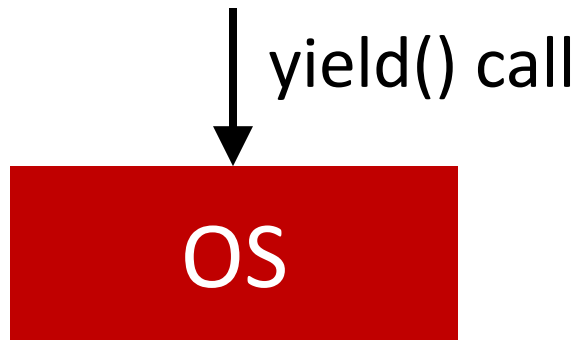
Cooperative Approach

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



Cooperative Approach

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



Cooperative Approach

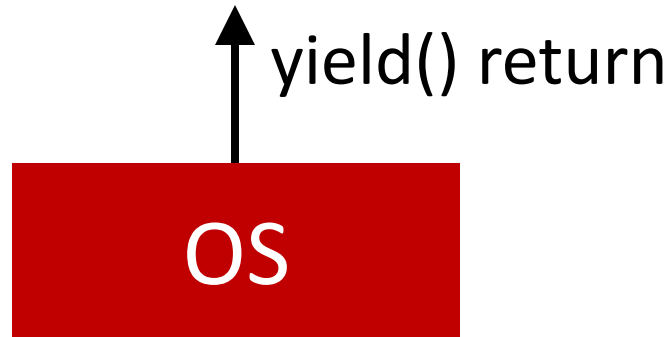
- Switch contexts for syscall interrupt
 - Special `yield()` system call



OS

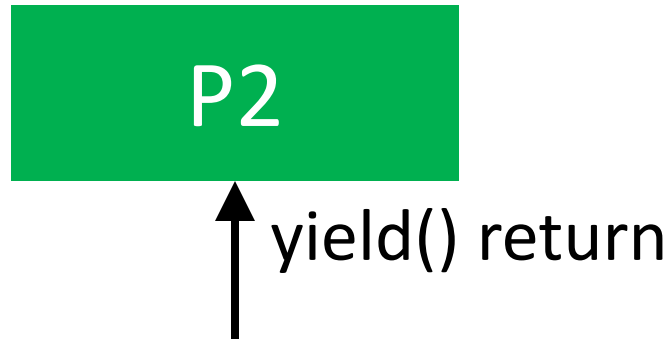
Cooperative Approach

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



Cooperative Approach

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



Cooperative Approach

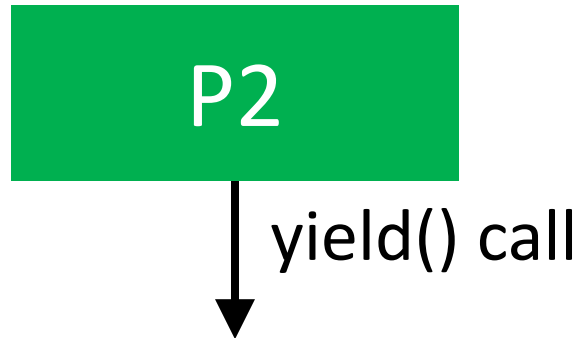
- Switch contexts for syscall interrupt
 - Special `yield()` system call



P2

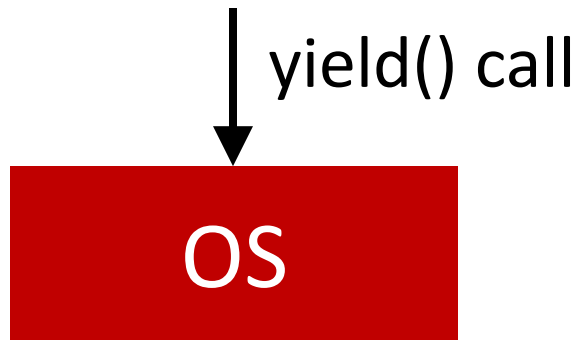
Cooperative Approach

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



Cooperative Approach

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



Cooperative Approach

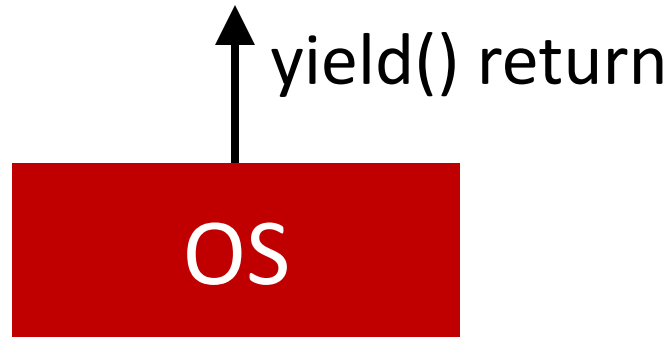
- Switch contexts for syscall interrupt
 - Special `yield()` system call



OS

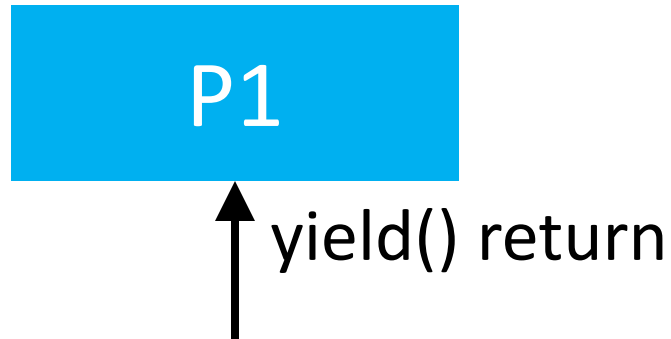
Cooperative Approach

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



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

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

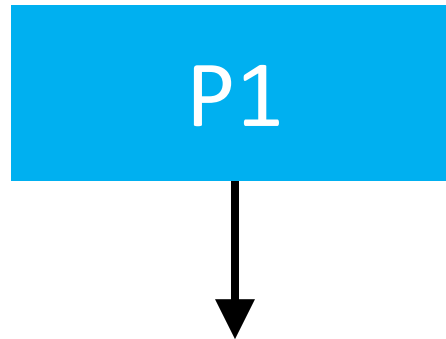
...

Preemptive Approach



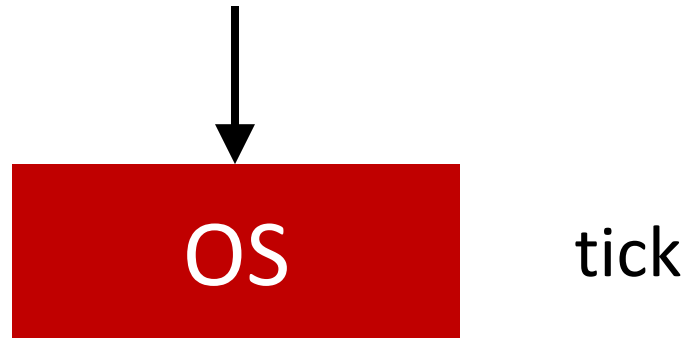
P1

Preemptive Approach



tick

Preemptive Approach

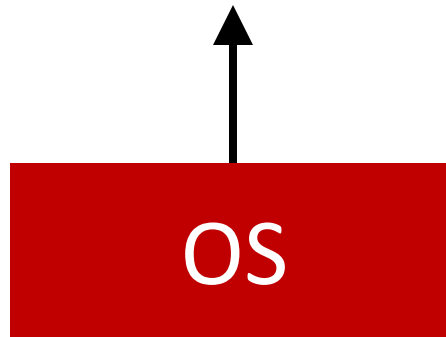


Preemptive Approach

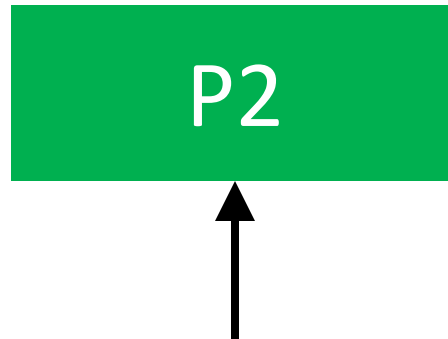


OS

Preemptive Approach



Preemptive Approach

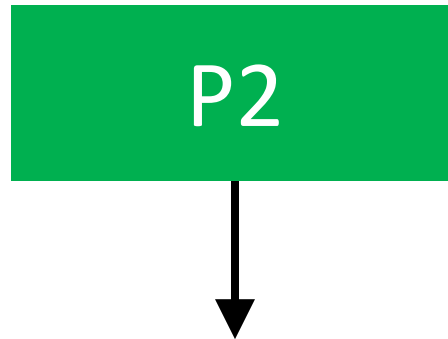


Preemptive Approach



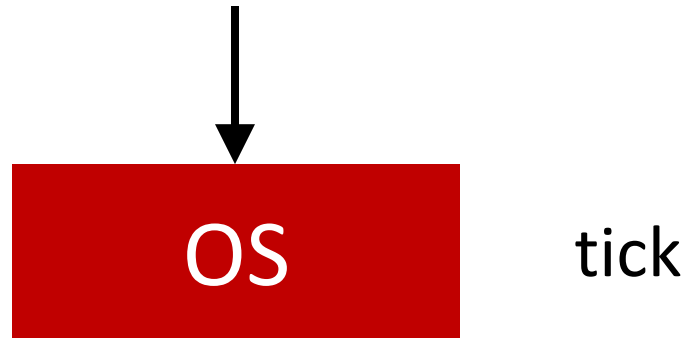
P2

Preemptive Approach



tick

Preemptive Approach

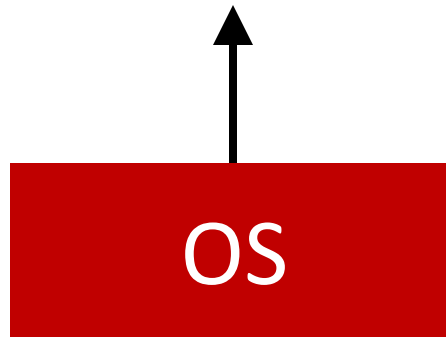


Preemptive Approach

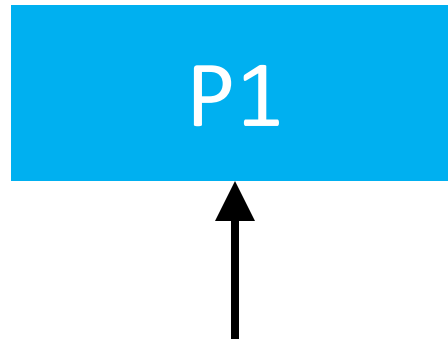


OS

Preemptive Approach



Preemptive Approach



Preemptive Approach



P1

LDE Summary

- Smooth **context switching** makes each process think it has its own CPU (virtualization!)
- **Limited direct execution** makes processes fast
- Hardware provides a lot of OS support
 - Limited direct execution
 - Timer interrupt
 - Automatic register saving

Threads

Why Thread Abstraction?

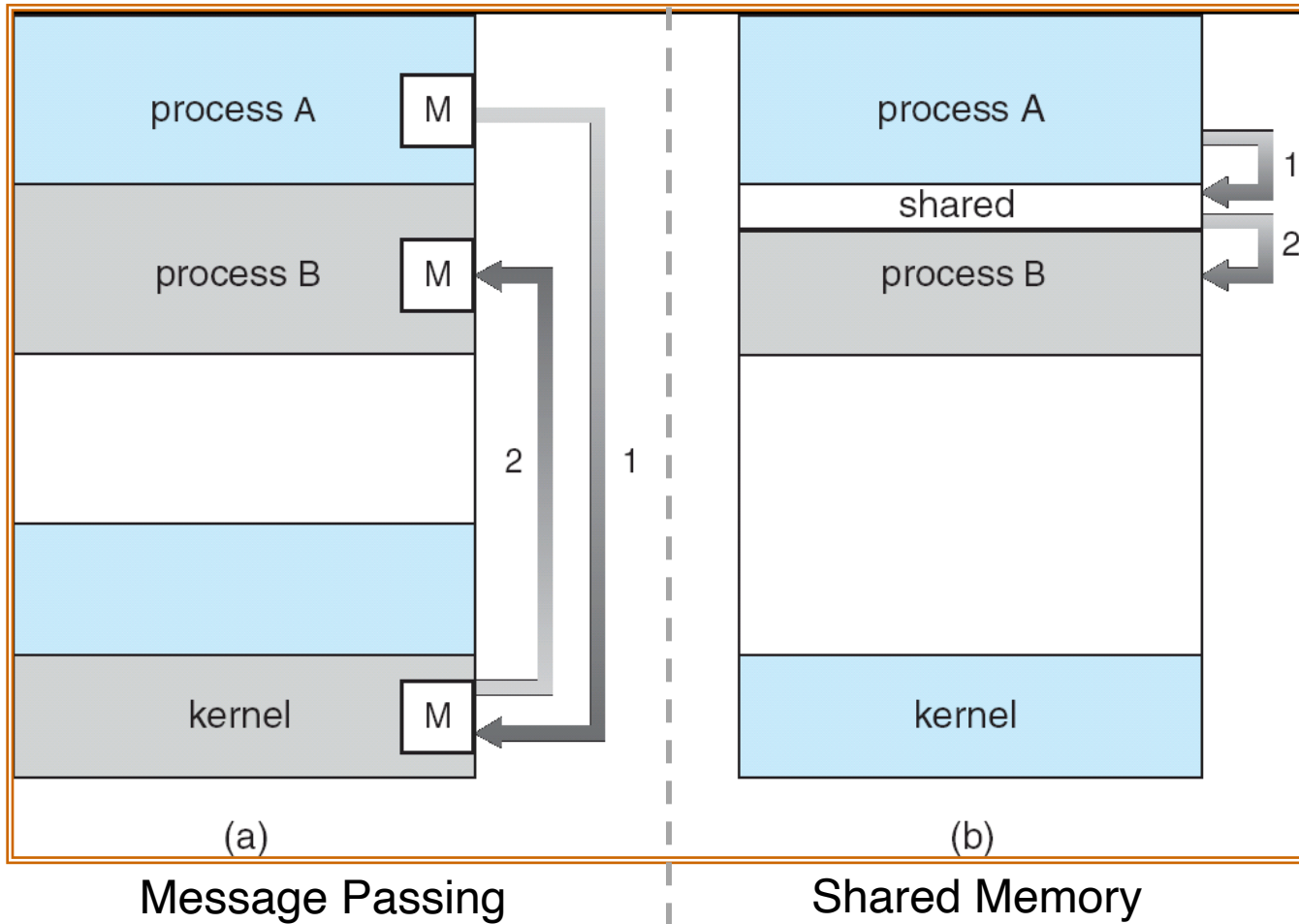
Process Abstraction: Challenge 1

- Inter-process communication (IPC)

Inter-Process Communication

- Mechanism for processes to communicate and to synchronize their actions.
- Two models
 - Communication through a shared memory region
 - Communication through message passing

Communication Models



Previously, in a distributed system, message-passing was the only possible communication model. However, remote direct memory access (RDMA) technique bridges this gap by providing remote memory access through network.

Communication through Message Passing

- Message system – processes communicate with each other **without** resorting to shared variables
- A message-passing facility must provide at least two operations:
 - `send(message, recipient)`
 - `receive(message, recipient)`
- With **indirect** communication, the messages are sent to and received from **mailboxes** (or, **ports**)
 - `send(A, message) /* A is a mailbox */`
 - `receive(A, message)`

Communication through Message Passing

- Message passing can be either **blocking** (**synchronous**) or **non-blocking** (**asynchronous**)
 - Blocking Send: The sending process is blocked until the message is received by the receiving process or by the mailbox
 - Non-blocking Send: The sending process resumes the operation as soon as the message is received by the kernel
 - Blocking Receive: The receiver blocks until the message is available
 - Non-blocking Receive: “Receive” operation does not block; it either returns a valid message or a default value (null) to indicate a non-existing message

Communication through Shared Memory

- The memory region to be shared must be explicitly defined
- System calls (Linux):
 - `shmget` creates a shared memory block
 - `shmat` maps/attaches an existing shared memory block into a process's address space
 - `shmdt` removes (“unmaps”) a shared memory block from the process's address space
 - `shmctl` is a general-purpose function allowing various operations on the shared block (receive information about the block, set the permissions, lock in memory, ...)
- Problems with `simultaneous access` to the shared variables
- Compilers for `concurrent programming languages` can provide direct support when declaring variables (e.g., “`shared int buffer`”)

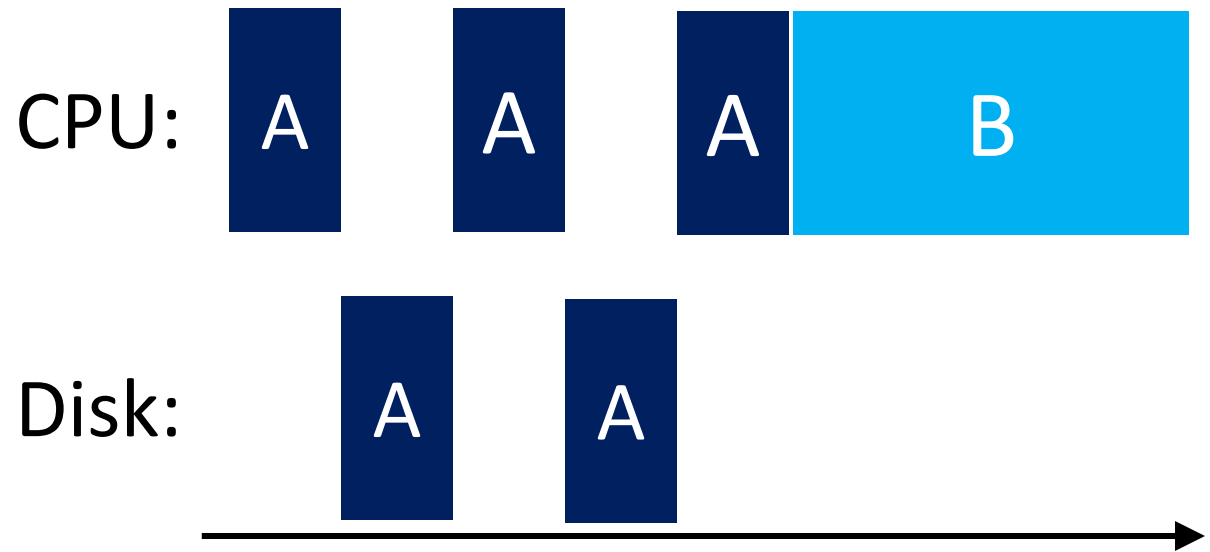
Process Abstraction: Challenge 1

- Inter-process communication (IPC)
 - Cumbersome programming!
 - Copying overheads (inefficient communication)
 - Expensive context switching (why expensive?)

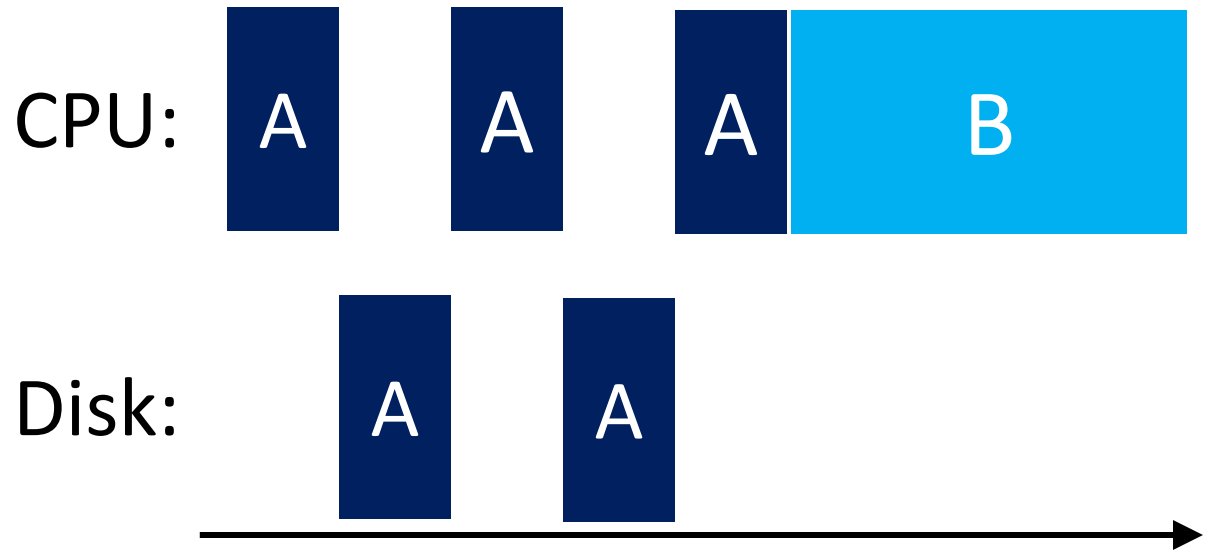
Process Abstraction: Challenge 2

- Inter-process communication (IPC)
 - Cumbersome programming!
 - Copying overheads (inefficient communication)
 - Expensive context switching (why expensive?)
- CPU utilization

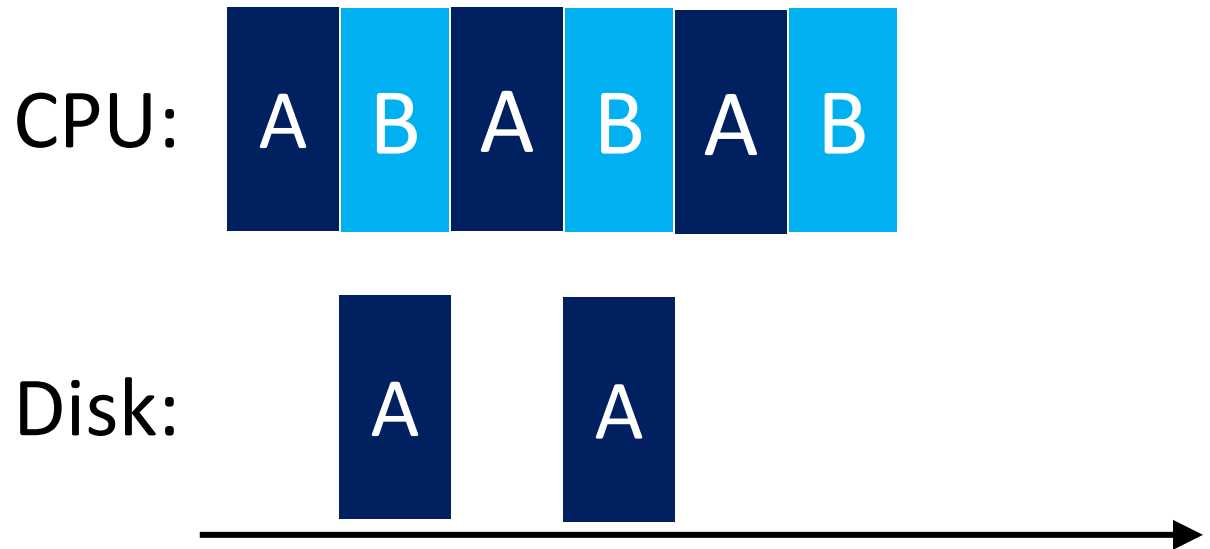
(a) Not interleaved



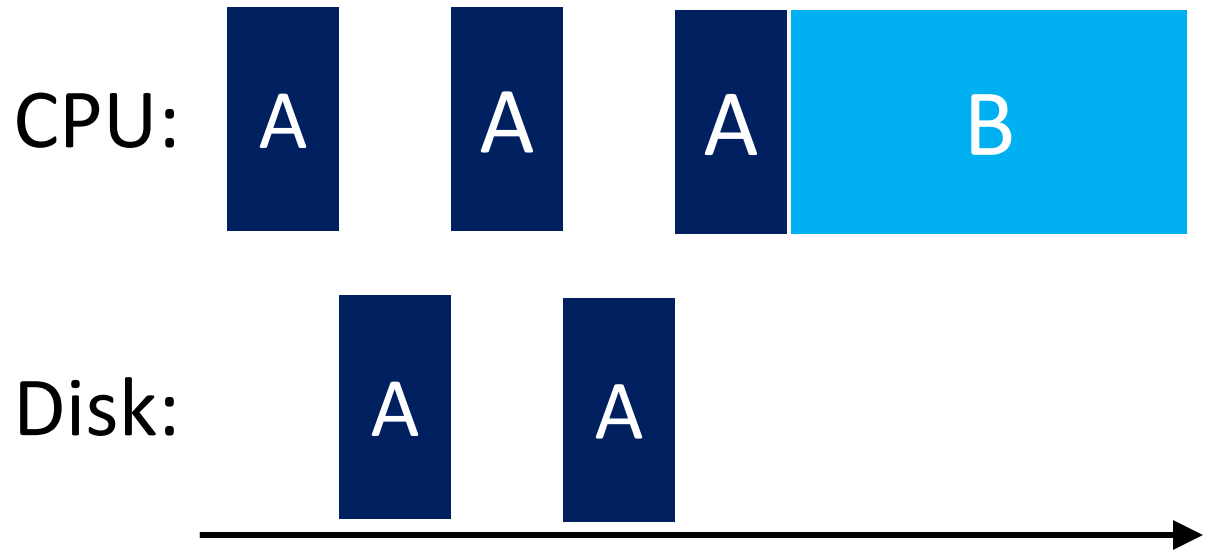
(a) Not interleaved



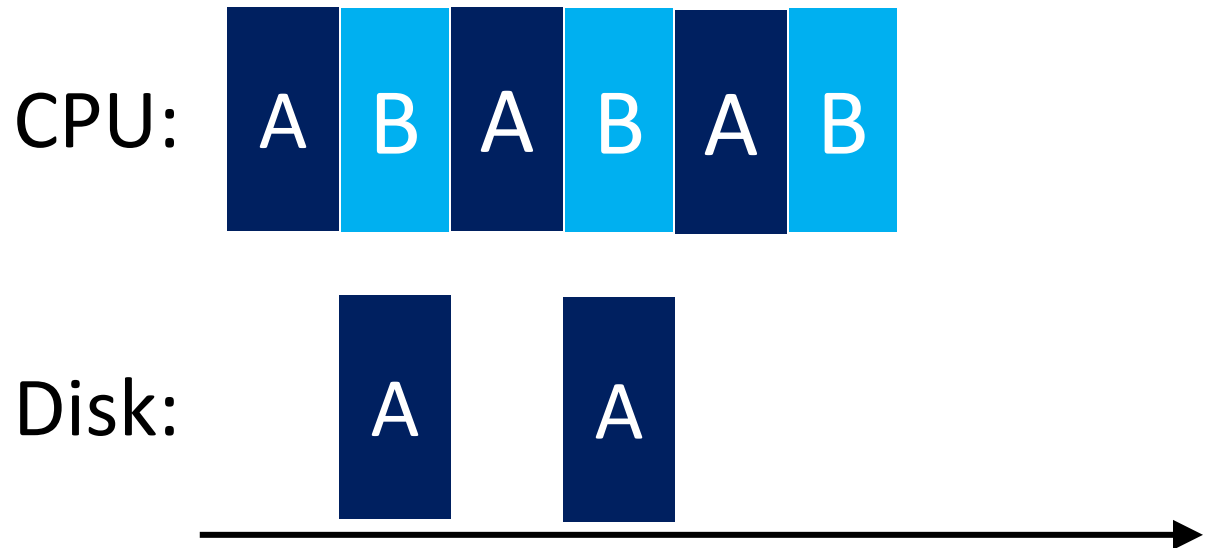
(b) Interleaved



(a) Not interleaved



(b) Interleaved



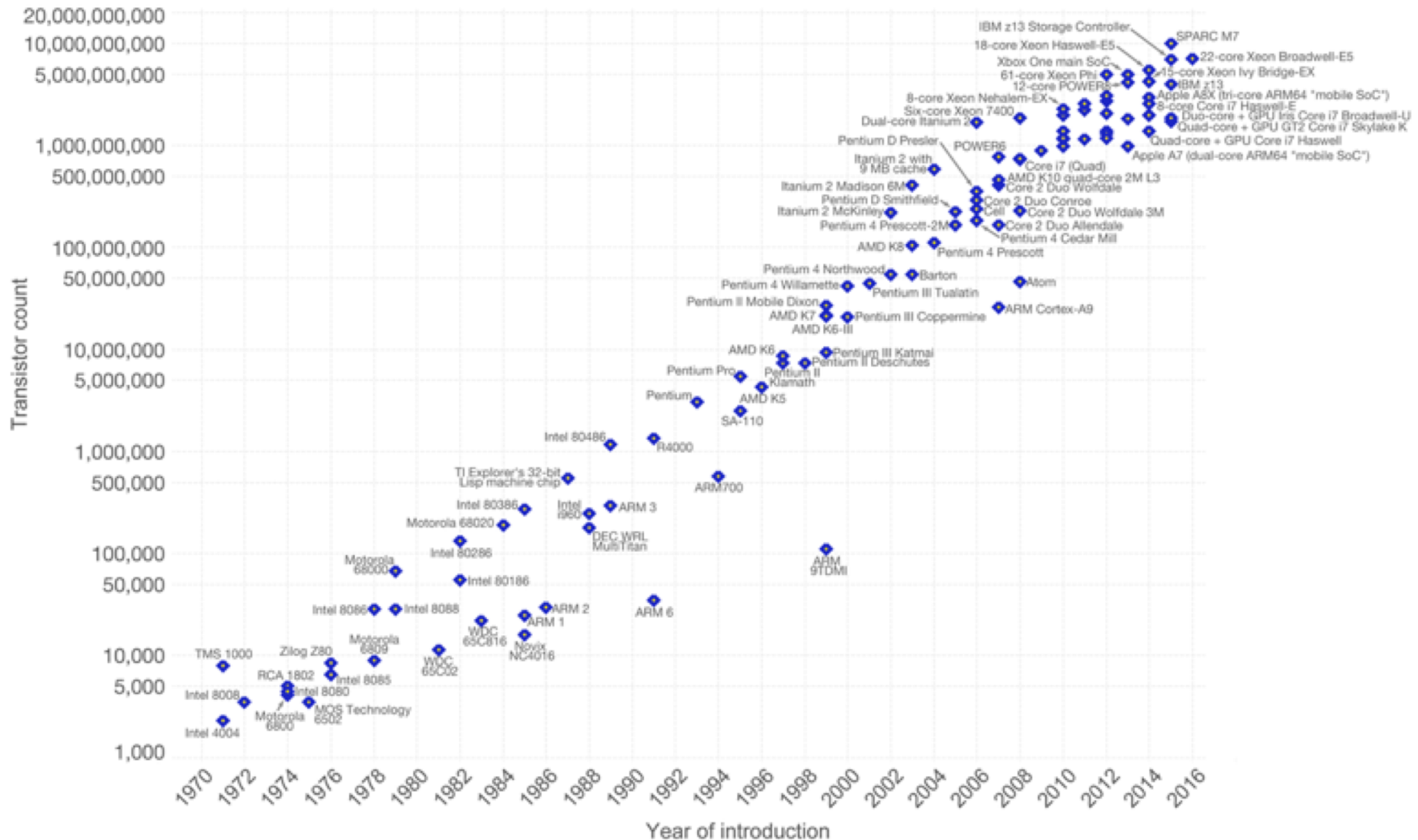
What if there is only one process?

Moore's law: # transistors doubles every ~2 years

Moore's Law – The number of transistors on integrated circuit chips (1971-2016)



Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are strongly linked to Moore's law.



Data source: Wikipedia (https://en.wikipedia.org/wiki/Transistor_count)

The data visualization is available at [OurWorldinData.org](https://www.ourworldindata.org). There you find more visualizations and research on this topic.

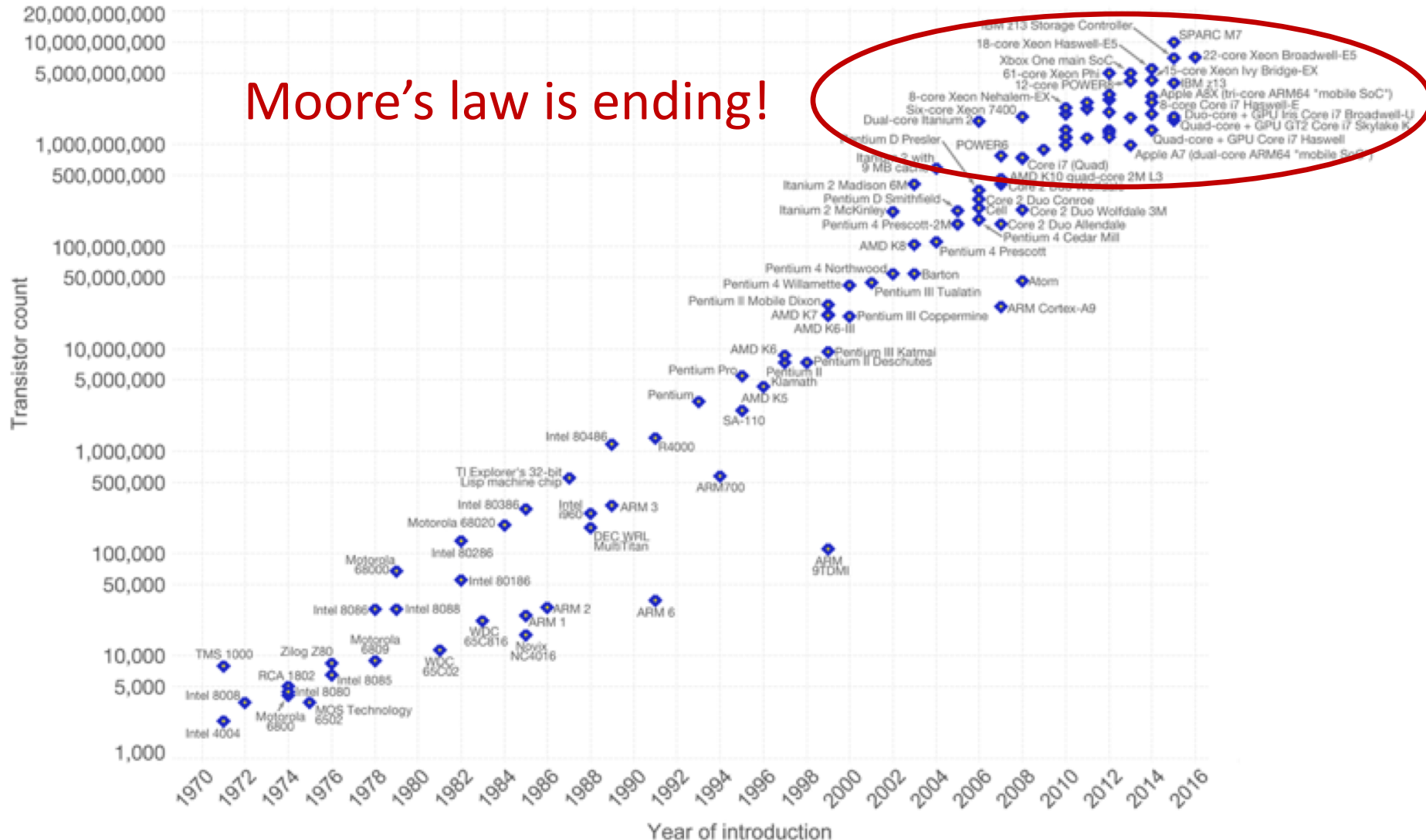
Licensed under [CC-BY-SA](https://creativecommons.org/licenses/by-sa/4.0/) by the author Max Roser.

Moore's law: # transistors doubles every ~2 years

Moore's Law – The number of transistors on integrated circuit chips (1971-2016)



Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are strongly linked to Moore's law.



Data source: Wikipedia (https://en.wikipedia.org/wiki/Transistor_count)

The data visualization is available at [OurWorldinData.org](https://www.ourworldindata.org). There you find more visualizations and research on this topic.

Licensed under [CC-BY-SA](https://creativecommons.org/licenses/by-sa/4.0/) by the author Max Roser.

CPU Trends – What Moore's Law Implies...

- The future
 - Same CPU speed
 - More cores (to scale-up)
- Faster programs => concurrent execution
- **Goal:** Write applications that fully utilize many CPU cores...

Goal

- Write applications that fully utilize many CPUs...

Strategy 1

- Build applications from many communication processes
 - Like Chrome (process per tab)
 - Communicate via `pipe()` or similar
- Pros/cons?

Strategy 1

- Build applications from many communication processes
 - Like Chrome (process per tab)
 - Communicate via `pipe()` or similar
- Pros/cons? – That we've talked about in previous slides
 - Pros: Don't need new abstractions!
 - Cons:
 - Cumbersome programming using IPC
 - Copying overheads
 - Expensive context switching

Strategy 2

- New abstraction: the **thread**

Introducing Thread Abstraction

- New abstraction: the **thread**
- Threads are just **like processes**, but threads **share the address space**

Thread

- A process, as defined so far, has only **one thread of execution**
- **Idea:** Allow multiple threads of **concurrently running** execution within the same process environment, **to a large degree independent** of each other
 - Each thread may be **executing different code** at the same time

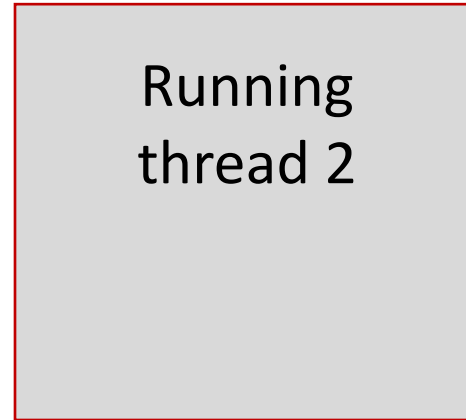
Process vs. Thread

- Multiple threads within a process will share
 - The address space
 - Open files (file descriptors)
 - Other resources
- Thread
 - Efficient and fast resource sharing
 - Efficient utilization of many CPU cores with only one process
 - Less context switching overheads

CPU 1



CPU 2



CPU 1

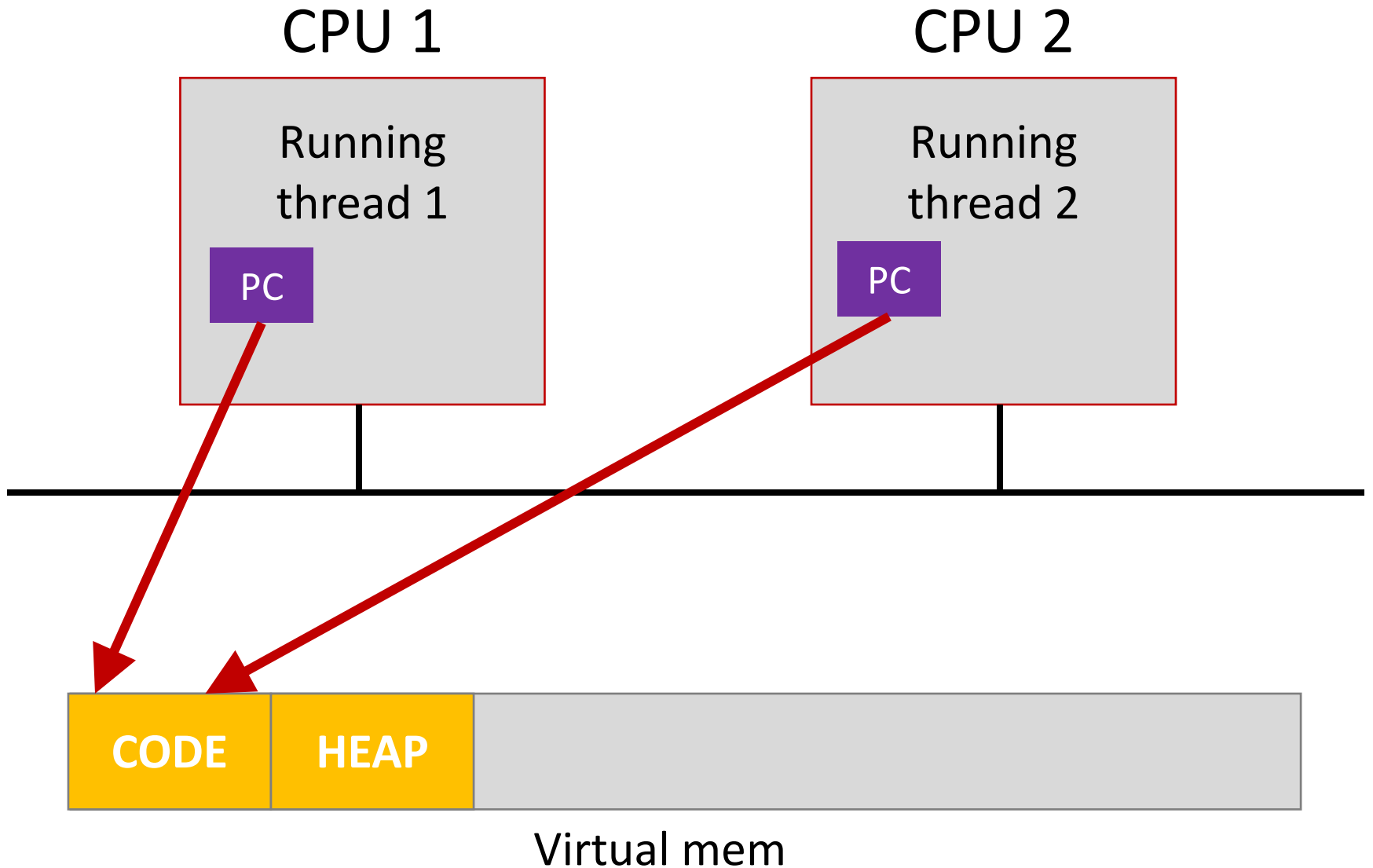
Running
thread 1

PC

CPU 2

Running
thread 2

PC



CPU 1

Running
thread 1

PC

CPU 2

Running
thread 2

PC

Each thread may be executing
different code at the same time

CODE

HEAP

Virtual mem

CPU 1

Running
thread 1

PC

CPU 2

Running
thread 2

PC



Virtual mem

CPU 1

Running
thread 1

PC

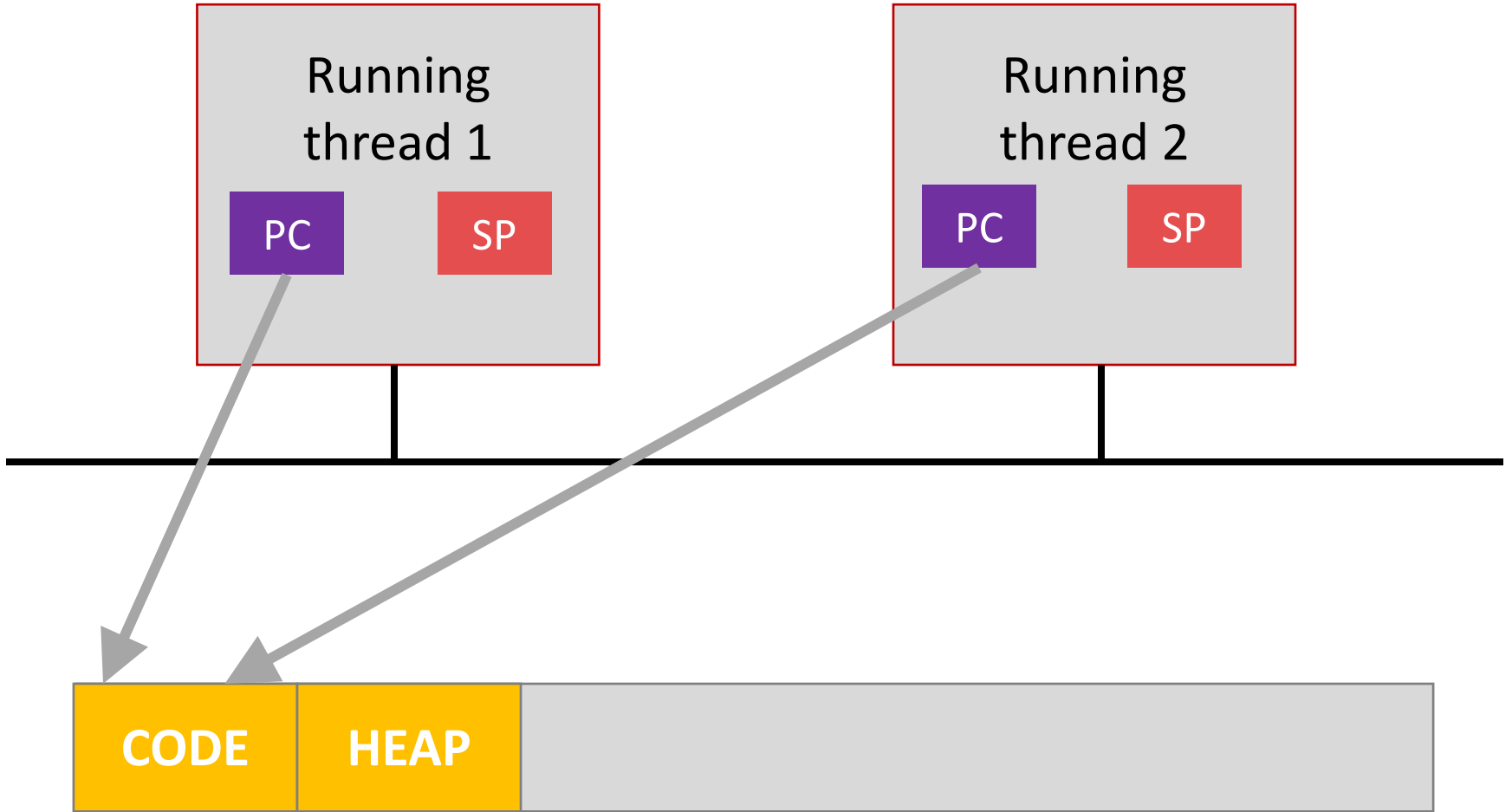
SP

CPU 2

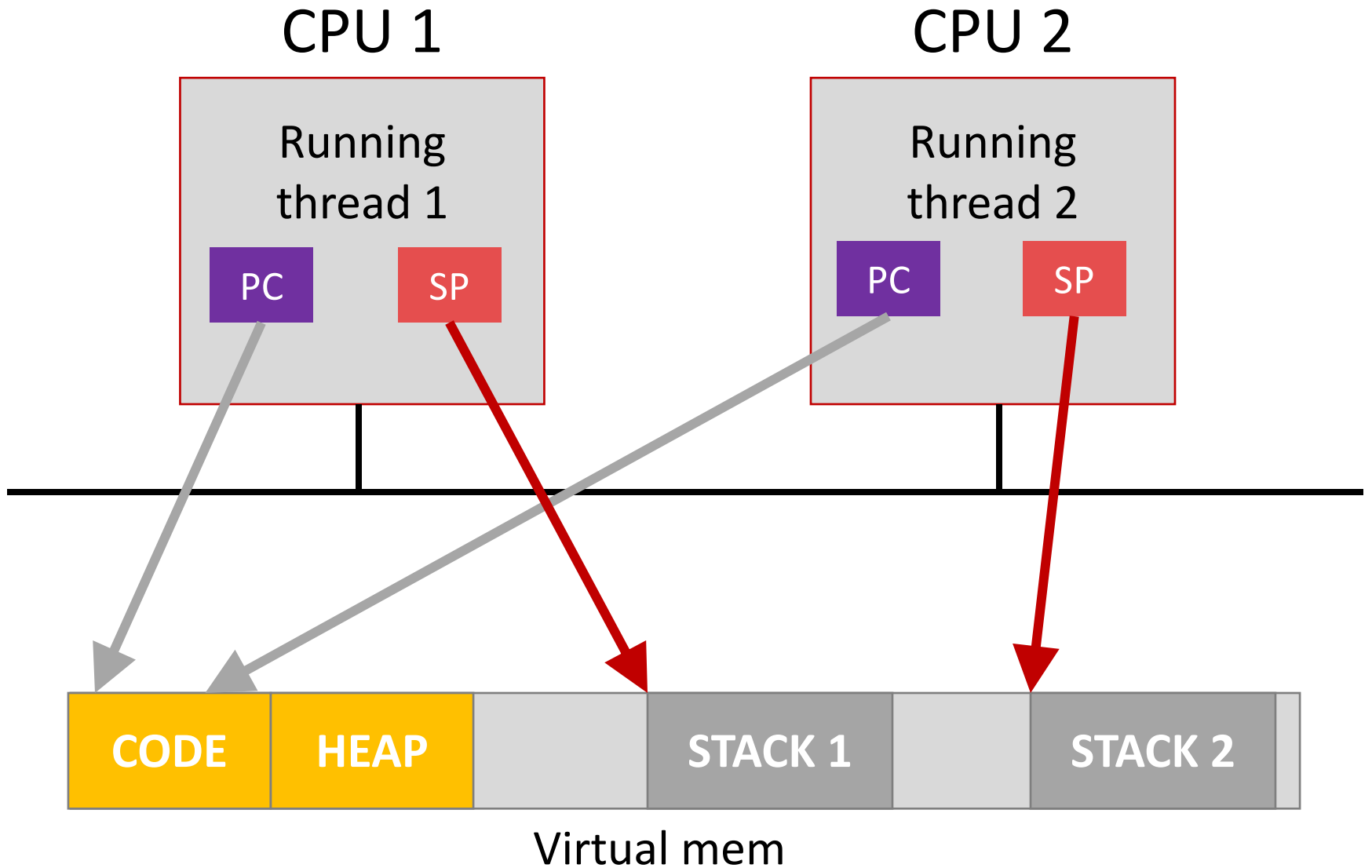
Running
thread 2

PC

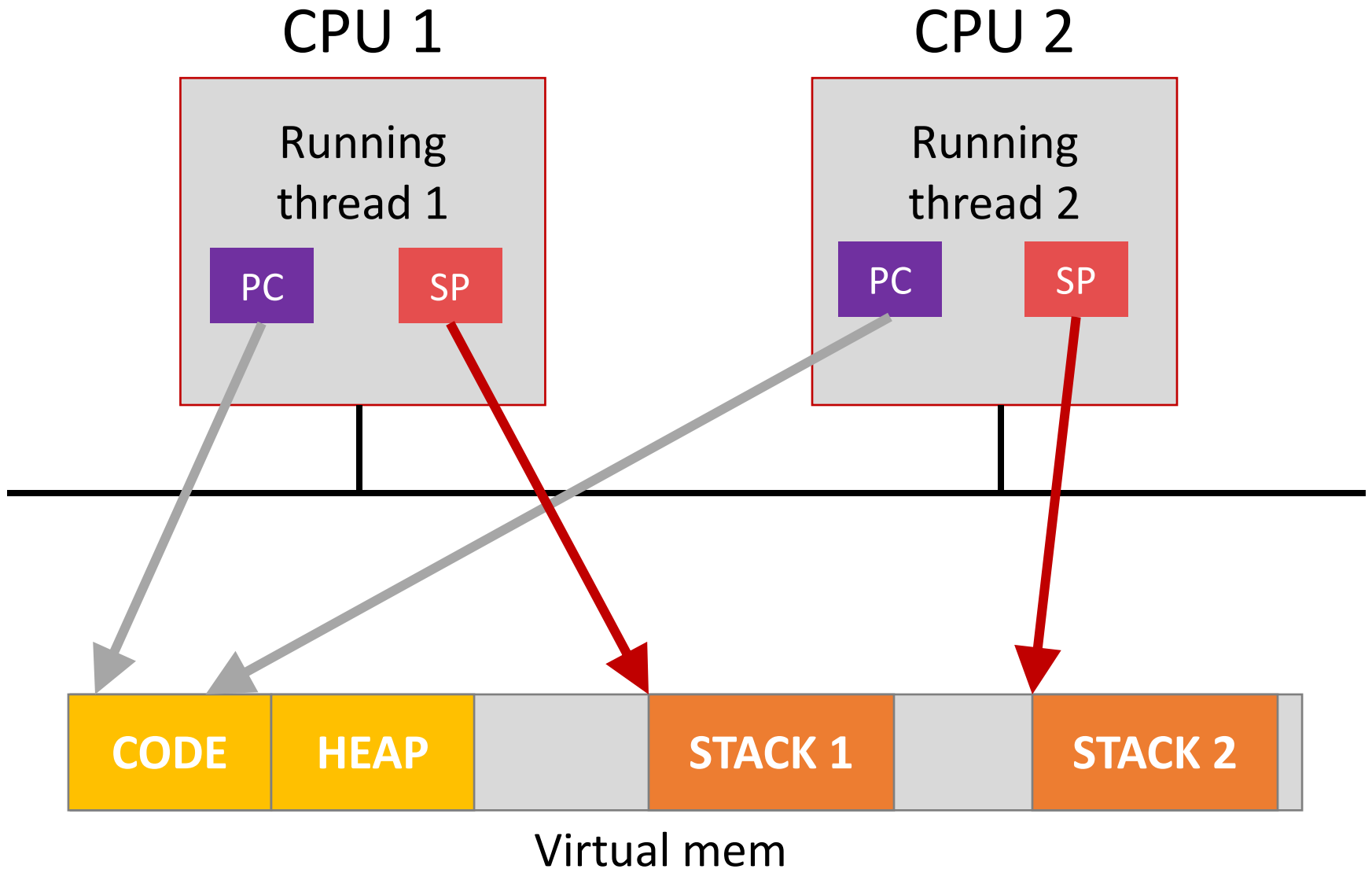
SP

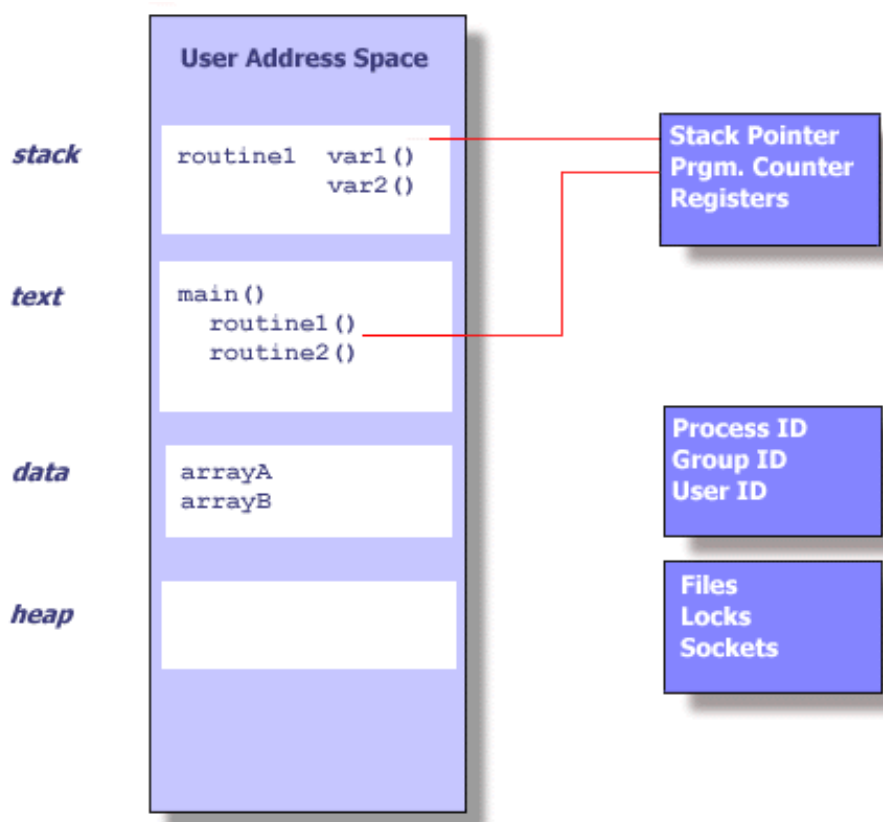


Virtual mem

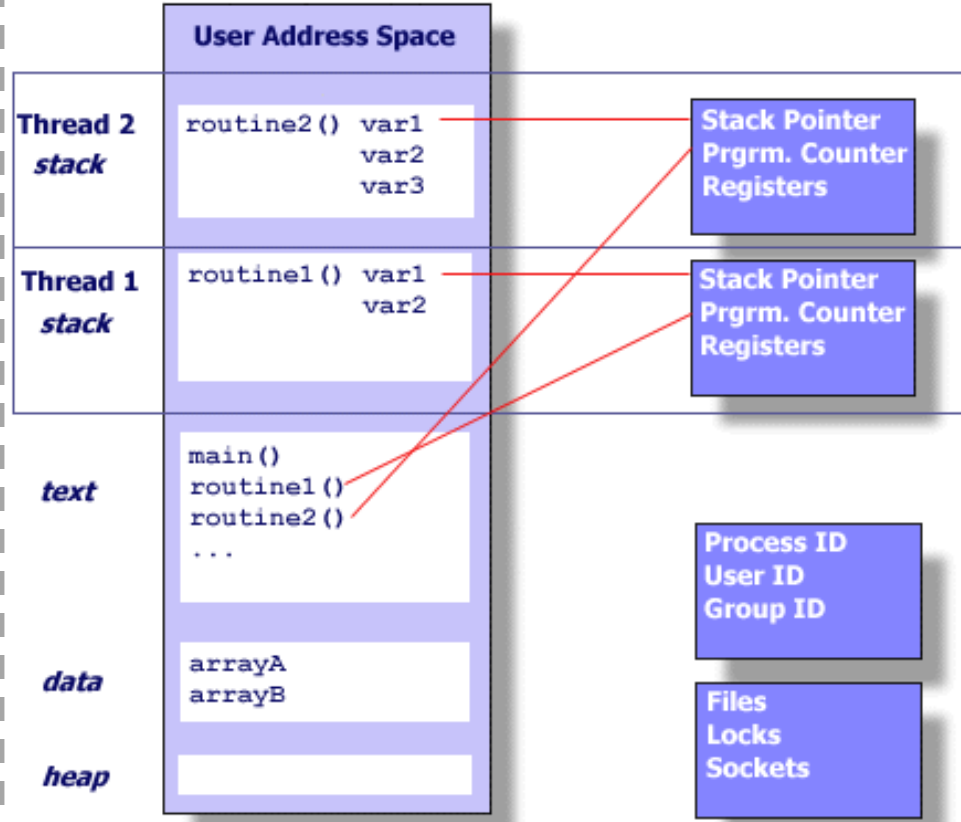


Thread executing **different functions** need **different stacks**



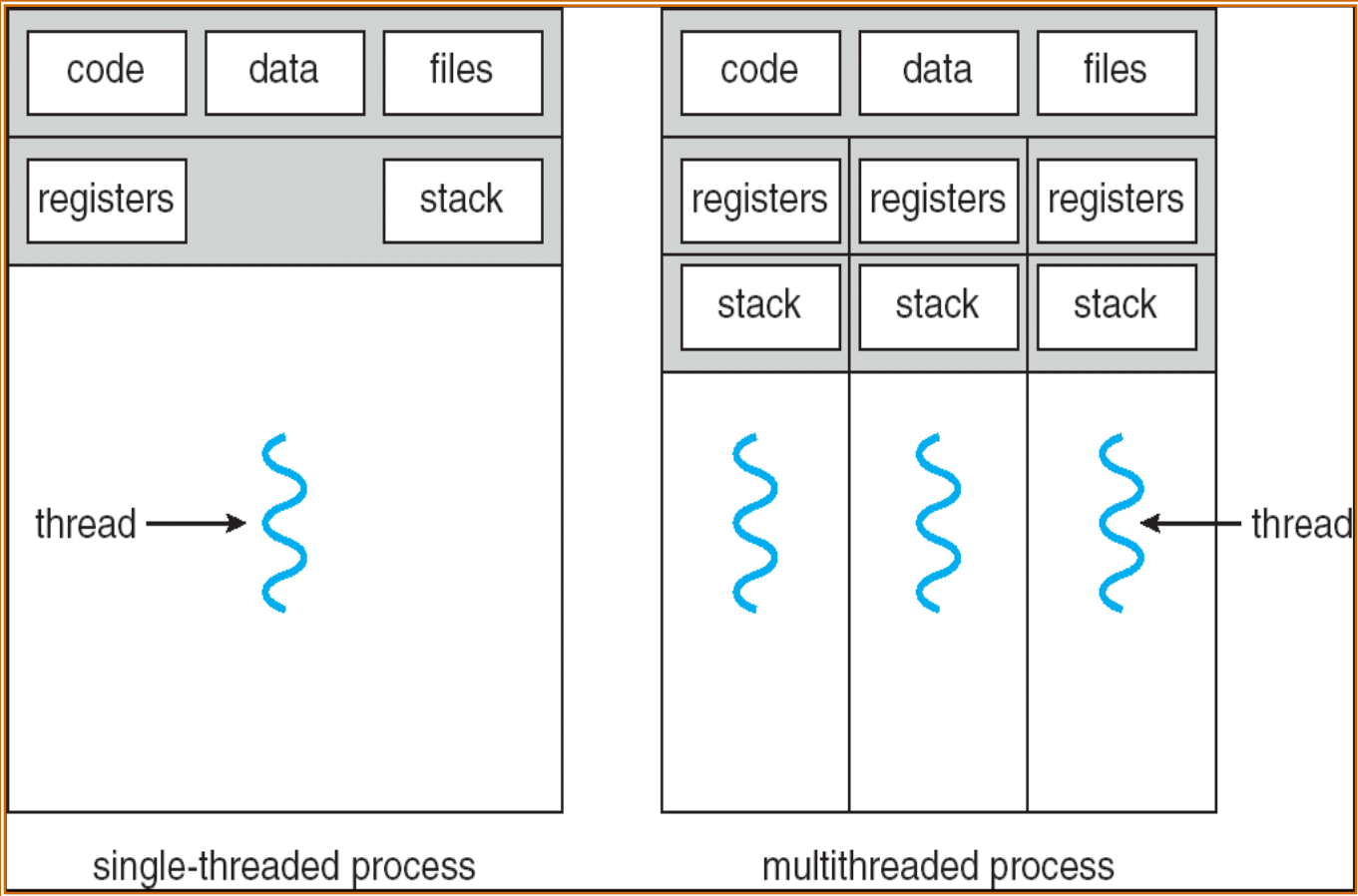


Linux process



Threads within a Linux process

Single- vs. Multi-threaded Process



Using Threads

- Processes usually start with a single thread
- Usually, library procedures are invoked to manage threads
 - `thread_create`: typically specifies the name of the procedure for the new thread to run
 - `thread_exit`
 - `thread_join`: blocks the calling thread until another (specific) thread has exited
 - `thread_yield`: voluntarily gives up the CPU to let another thread run

Pthread

- A POSIX standard (IEEE 1003.1c) API for thread creation and synchronization
- API specifies behavior of the thread library, implementation is up to development of the library
- Common in UNIX (e.g., Linux) OSes

Pthread APIs

Thread Call	Description
<code>pthread_create</code>	Create a new thread in the caller's address space
<code>pthread_exit</code>	Terminate the calling thread
<code>pthread_join</code>	Wait for a thread to terminate
<code>pthread_mutex_init</code>	Create a new mutex
<code>pthread_mutex_destroy</code>	Destroy a mutex
<code>pthread_mutex_lock</code>	Lock a mutex
<code>pthread_mutex_unlock</code>	Unlock a mutex
<code>pthread_cond_init</code>	Create a condition variable
<code>pthread_cond_destroy</code>	Destroy a condition variable
<code>pthread_cond_wait</code>	Wait on a condition variable
<code>pthread_cond_signal</code>	Release one thread waiting on a condition variable

Pthread APIs

Thread Call	Description	
<code>pthread_create</code>	Create a new thread in the caller's address space	Thread creation
<code>pthread_exit</code>	Terminate the calling thread	
<code>pthread_join</code>	Wait for a thread to terminate	
<code>pthread_mutex_init</code>	Create a new mutex	Thread lock
<code>pthread_mutex_destroy</code>	Destroy a mutex	
<code>pthread_mutex_lock</code>	Lock a mutex	
<code>pthread_mutex_unlock</code>	Unlock a mutex	
<code>pthread_cond_init</code>	Create a condition variable	Thread CV
<code>pthread_cond_destroy</code>	Destroy a condition variable	
<code>pthread_cond_wait</code>	Wait on a condition variable	
<code>pthread_cond_signal</code>	Release one thread waiting on a condition variable	

Example of Using Pthread

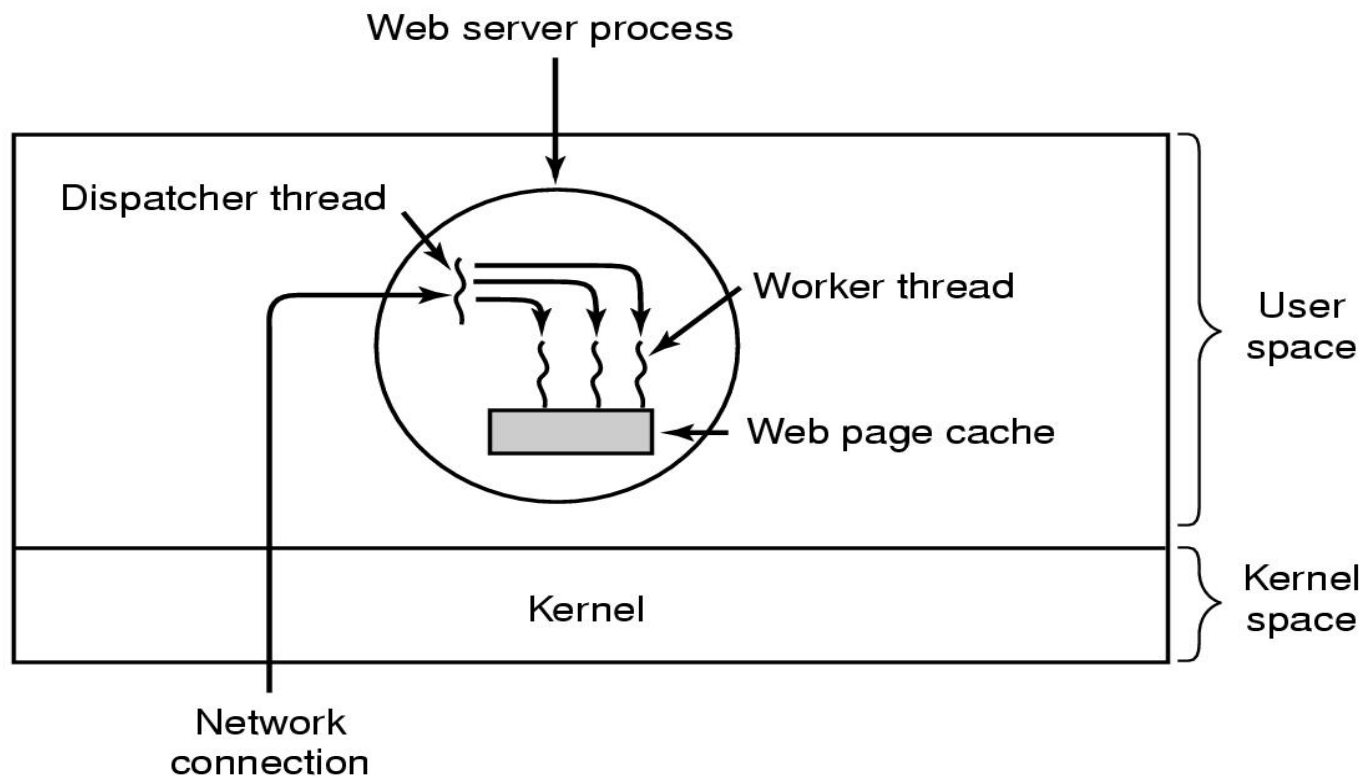
```
1  #include <stdio.h>
2  #include <assert.h>
3  #include <pthread.h>
4
5  void *mythread(void *arg) {
6      printf("%s\n", (char *) arg);
7      return NULL;
8  }
9
10 int
11 main(int argc, char *argv[]) {
12     pthread_t p1, p2;
13     int rc;
14     printf("main: begin\n");
15     rc = pthread_create(&p1, NULL, mythread, "A"); assert(rc == 0);
16     rc = pthread_create(&p2, NULL, mythread, "B"); assert(rc == 0);
17     // join waits for the threads to finish
18     rc = pthread_join(p1, NULL); assert(rc == 0);
19     rc = pthread_join(p2, NULL); assert(rc == 0);
20     printf("main: end\n");
21     return 0;
22 }
```

Demo: Basic Threads

- Fork the demo code repo at:
<https://github.com/tddg/demo-ostep-code>
- In today's lecture, we showed the demo in dir:
thread-api

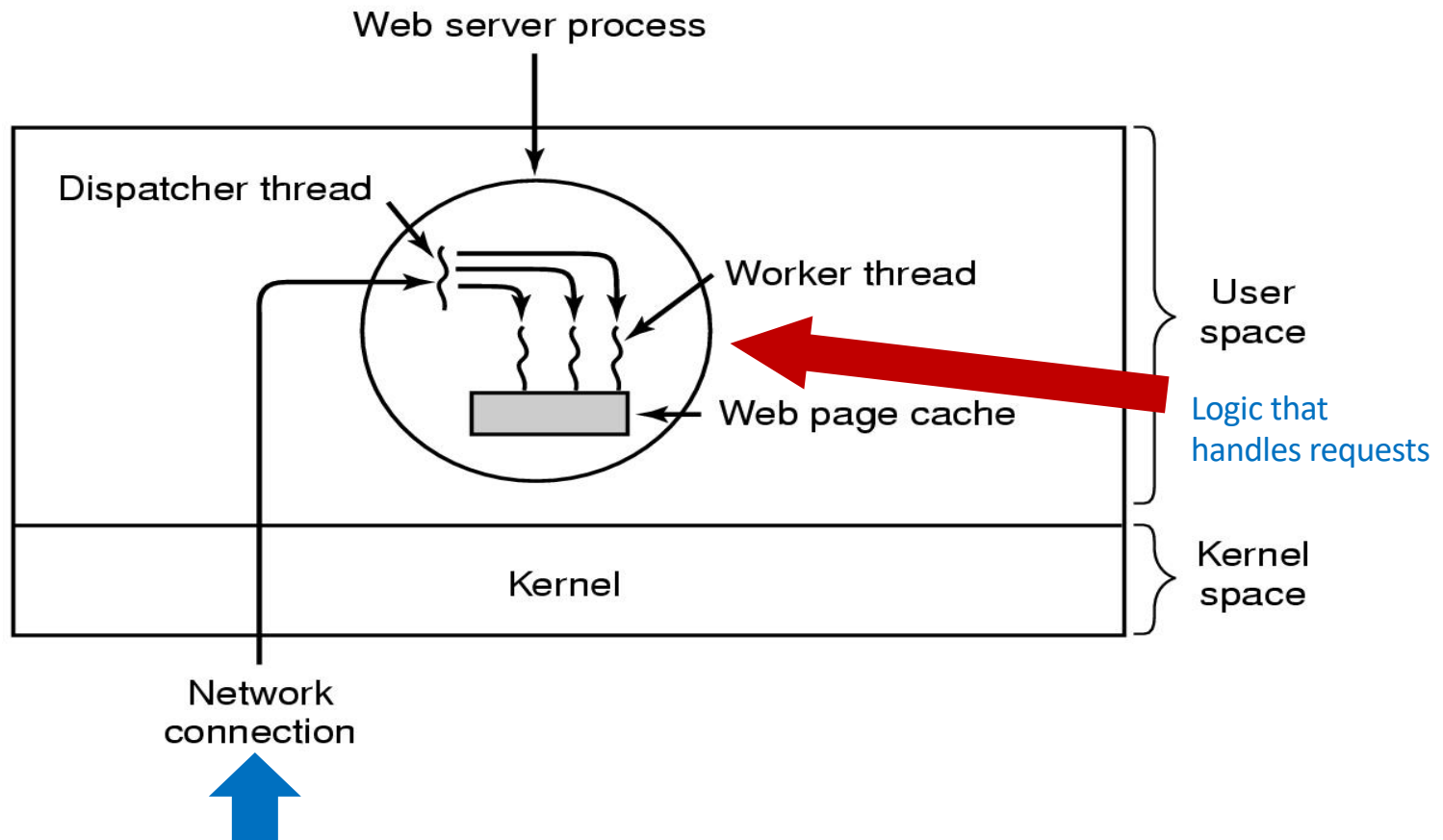
Example Multithreaded Applications

A multithreaded web server



Example Multithreaded Applications

A multithreaded web server



Code Sketch

```
while (TRUE) {  
    get_next_request(&buf);  
    handoff_work(&buf);  
}
```

(a) Dispatcher thread

```
while (TRUE) {  
    wait_for_work(&buf);  
    check_cache(&buf; &page);  
    if (not_in_cache)  
read_from_disk(&buf, &page);  
    return_page(&page);  
}
```

(b) Worker thread

Benefits of Multi-threading

- **Resource sharing**
 - Sharing the address space and other resources may result in high degree of cooperation
- **Economy**
 - Creating/managing processes much more time consuming than managing threads: e.g., context switch
- **Better utilization of multicore architectures**
 - Threads are doing job concurrently (in parallel)
 - Multithreading an interactive application may allow a program to continue running even if part of it is blocked or performing a lengthy operation

Real-world Example: Memcached

- Memcached — A high-performance memory-based caching system
 - 14k lines of C source code
 - <https://memcached.org/>
- A typical multithreaded server implementation
 - `Pthread + libevent`
 - A dispatcher thread dispatches newly coming connections to the worker threads in a round-robin manner
 - Event-driven: Each worker thread is responsible for serving requests from the established connections



Multithreading vs. Multi-processes

- Real-world debate
 - Multithreading vs. Multi-processes
 - Memcached vs. Redis
- Redis — A single-threaded memory-based data store
 - <https://redis.io/>



Memcached



redis

Wish List for Redis...

<http://goo.gl/N9UTKD>

Wish List For Redis

- Explicit memory management.
- **Deployable (Lua) Scripts.** Talked about near the start.
- **Multi-threading.** Would make cluster management easier. Twitter has a lot of “tall boxes,” where a host has 100+ GB of memory and a lot of CPUs. To use the full capabilities of a server a lot of Redis instances need to be started on a physical machine. With multi-threading fewer instances would need to be started which is much easier to manage.