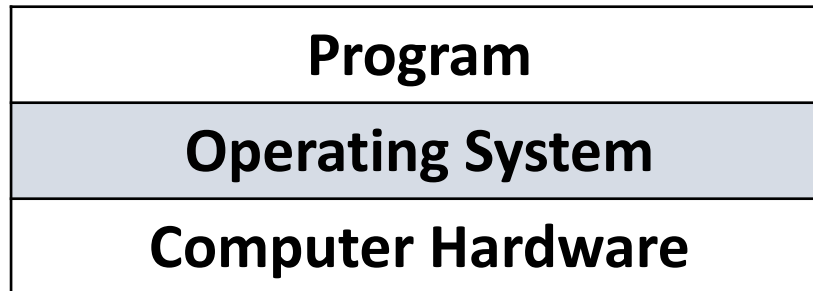


# Operating Systems Intro

11/1/2016

# What is an operating system?

The OS is an interface layer between a user's programs and hardware.



It provides an abstract view of the hardware that makes it easier to use.

How many programs do you think are running on my laptop now?

A. 0-1

B. 2-9

C. 10-99

D. 100-999

E. 1000+

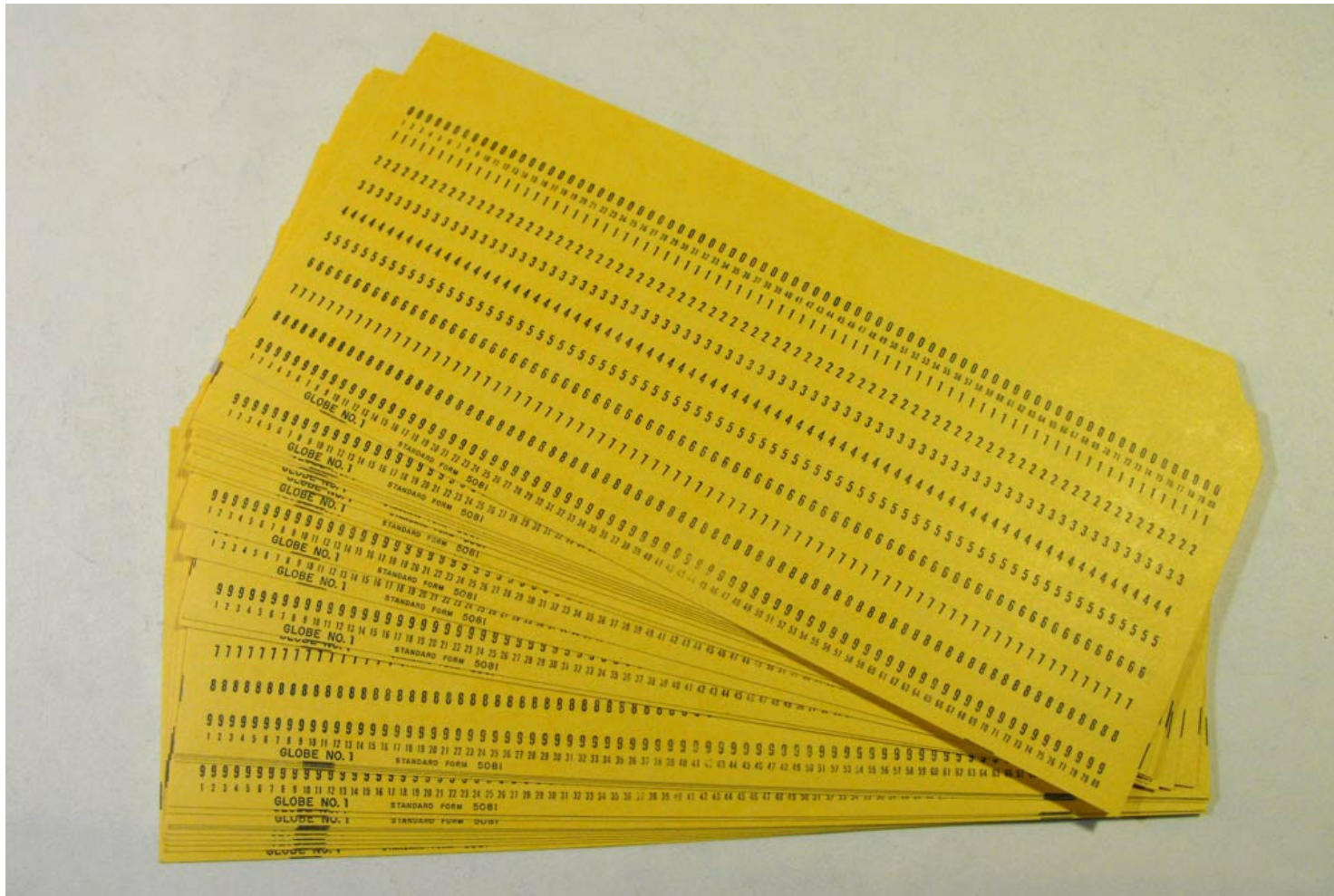
# Sharing resources

All of these programs need to share:

- Memory
- Processor time
- I/O devices

But when you write a program, you don't have to think about the hundreds of other programs that could be running.

# Before Operating Systems



# OS provides virtualization

- Rather than exposing real hardware, introduce a “virtual”, abstract notion of the resource
- Multiple virtual processors
  - By rapidly switching CPU use
- Multiple virtual memories
  - By memory partitioning and re-addressing
- Virtualized devices
  - By simplifying interfaces, and using other resources to enhance function

# We'll focus on the OS 'kernel'

- “Operating system” has many interpretations
  - E.g., all software on machine minus applications
- Our focus is the *kernel*
  - What's necessary for everything else to work
  - Low-level resource control

# The Kernel

- All programs depend on it
  - Loads and runs them
  - Exports system calls to programs
- Works closely with hardware
  - Accesses devices
  - Responds to interrupts
- Allocates basic resources
  - CPU time, memory space
  - Controls I/O devices: display, keyboard, disk, network

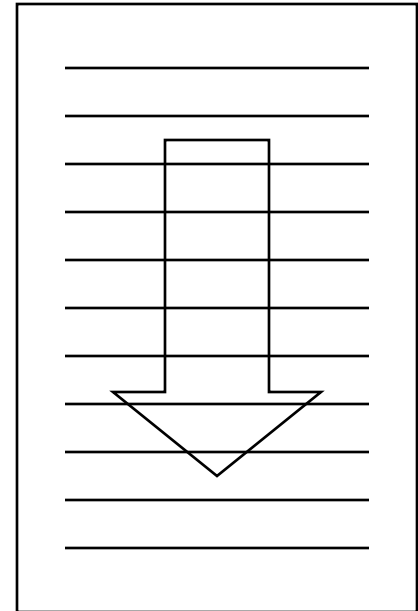


# Kernel provides common functions

- Some functions useful to many programs
  - I/O device control
  - Memory allocation
- Place these functions in central place (kernel)
  - Called by programs (system calls)
  - Or accessed implicitly
- What functions should be included?
  - How many programs need to benefit?
  - Might kernel get too big?

# Main Abstraction: The Process

- Abstraction of a running program
  - “a program in execution”
- Dynamic
  - Has state, changes over time
  - Whereas a program is static
- Basic operations
  - Start/end
  - Suspend/resume



# Basic Resources for Processes

To run, process needs some basic resources:

- CPU
  - Processing cycles (time)
  - To execute instructions
- Memory
  - Bytes or words (space)
  - To maintain state
- Other resources (e.g., I/O)
  - Network, disk, terminal, printer, etc.

# Machine State of a Process

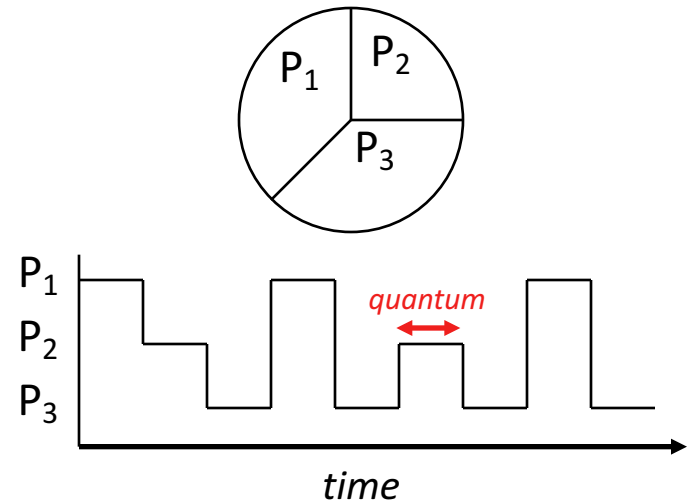
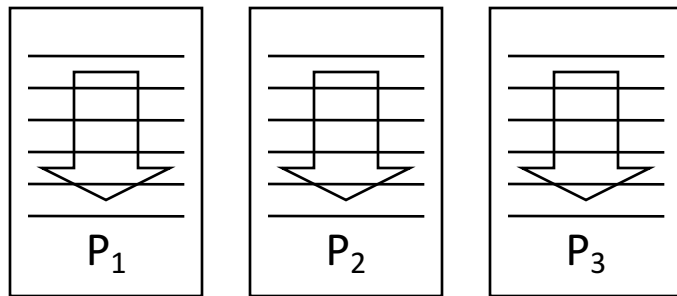
- CPU or processor context
  - PC (program counter)
  - SP (stack pointer)
  - General purpose registers
- Memory
  - Code
  - Global Variables
  - Stack of activation records / frames
  - Other (registers, memory, kernel-related state)

Must keep track of these  
for every running process !

# CPU

- Abstraction goal: make every process think it's running on the CPU all the time.
  - Alternatively: If a process was removed from the CPU and then given it back, it shouldn't be able to tell
- Reality: put a process on CPU, let it run for a short time (~10 ms), switch to another, ... (context switching)

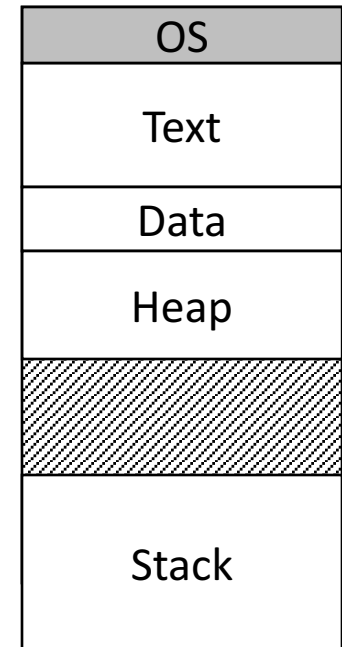
# Timesharing: Sharing the CPU



- Multiple processes, single CPU (uniprocessor)
- Conceptually, each process makes progress over time
- In reality, each periodically gets quantum of CPU time
- Illusion of parallel progress by rapidly switching CPU

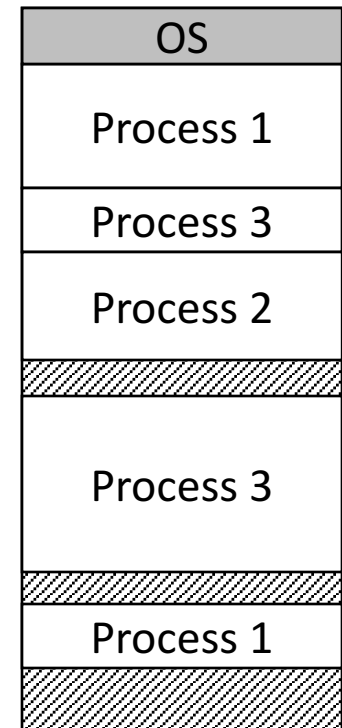
# Memory

- Abstraction goal: make every process think it has the same memory layout.
  - MUCH simpler for compiler if the stack always starts at 0xFFFFFFFF, etc.



# Memory

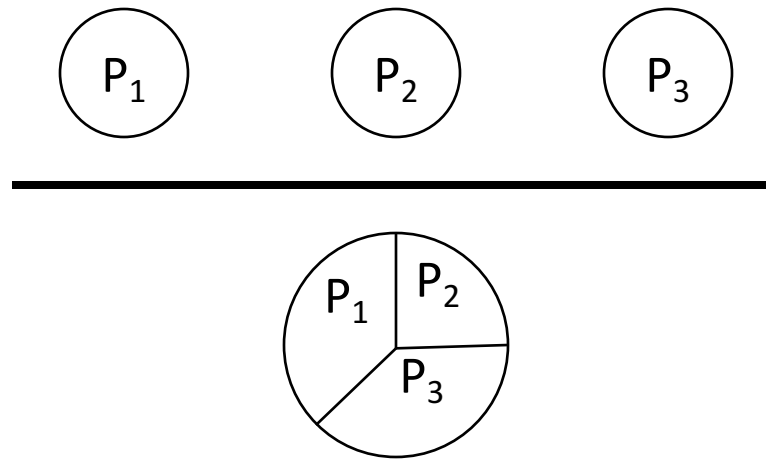
- Abstraction goal: make every process think it has the same memory layout.
  - MUCH simpler for compiler if the stack always starts at 0xFFFFFFFF, etc.
- Reality: there's only so much memory to go around, and no two processes should use the same (physical) memory addresses (unless they're sharing on purpose).



OS (with help from hardware) will keep track of who's using each memory region.



# Virtual Memory: Sharing Storage



- Like CPU cache, memory is a cache for disk.
- Processes never need to know where their memory truly is, OS translates virtual addresses into physical addresses for them.

# The Kernel

- So...how / when should the kernel execute?
- What would you do if you were to build such a thing?

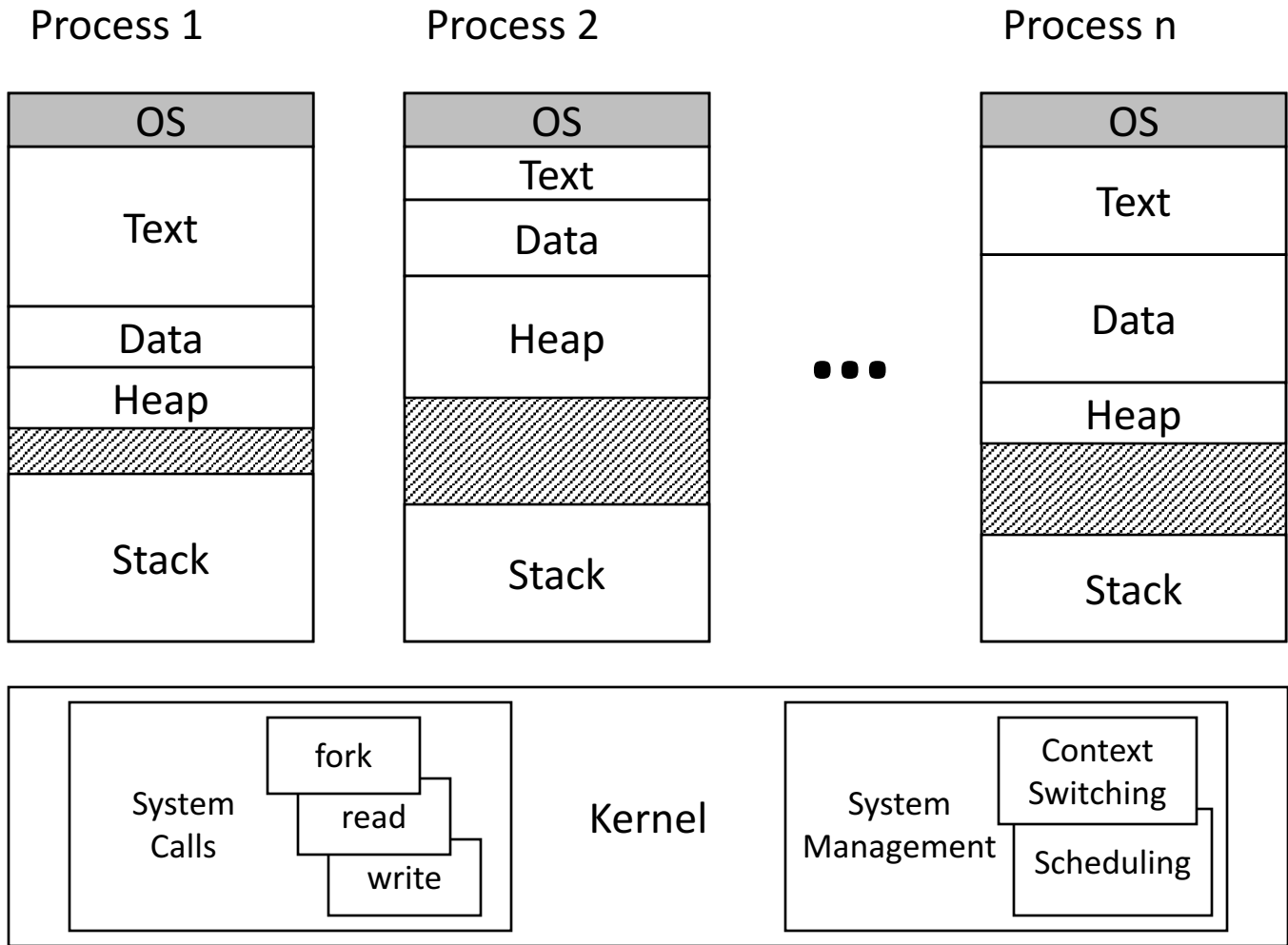
# The operating system kernel...

- A. Executes as a process.
- B. Is always executing, in support of other processes.
- C. Should execute as little as possible.
- D. More than one of the above. (Which ones?)
- E. None of the above.

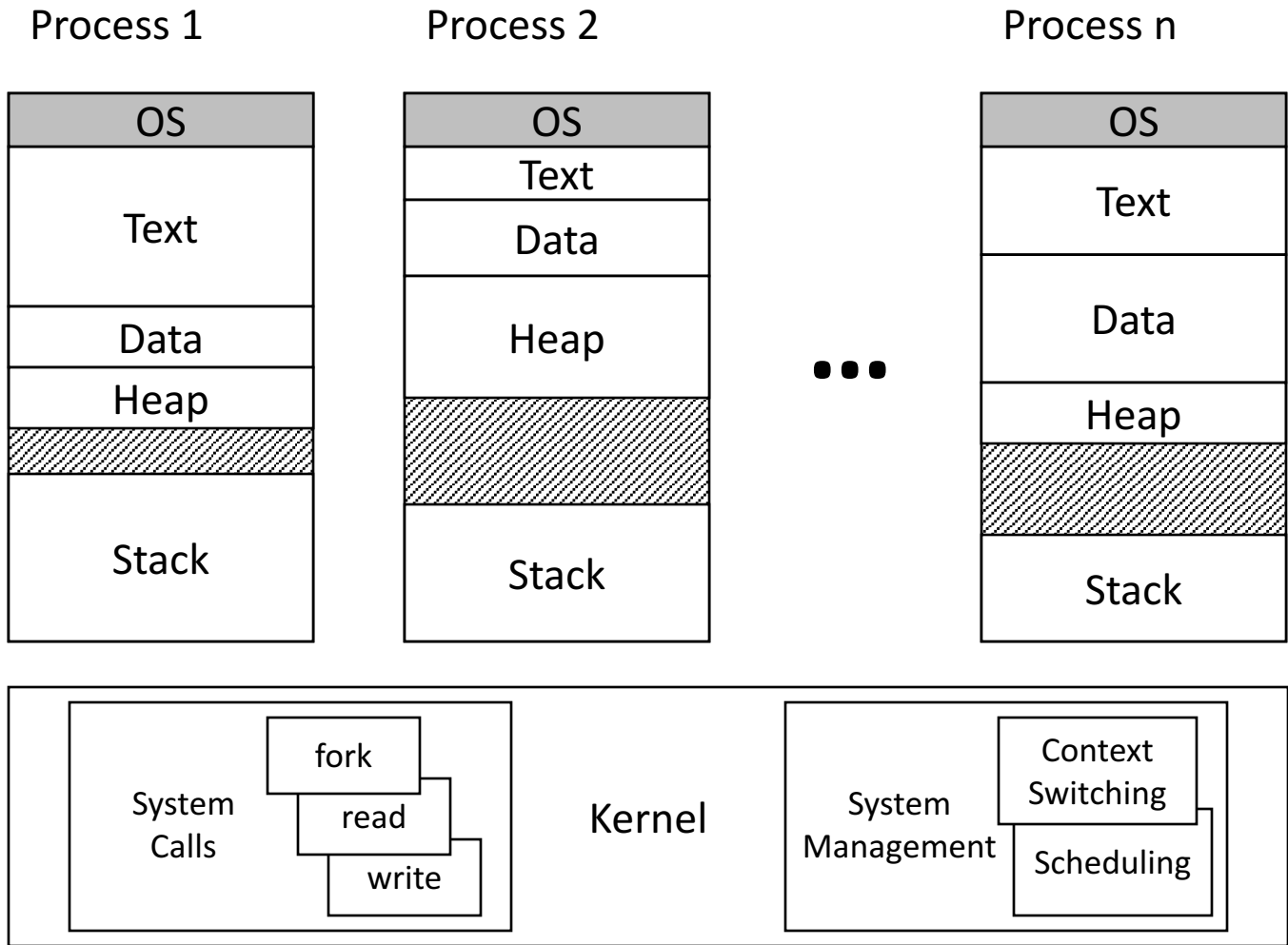
# Process vs. Kernel

- The kernel is the code that supports processes
  - System calls: fork ( ), exit ( ), read ( ), write ( ), ...
  - System management: context switching, scheduling
- When does the kernel run?
  - When system call or hardware interrupt occurs
- Is the kernel a process?
  - No, it supports processes and devices
  - Runs as an extension of process making system call
  - Runs in response to device issuing interrupt

# Process and Kernel Model

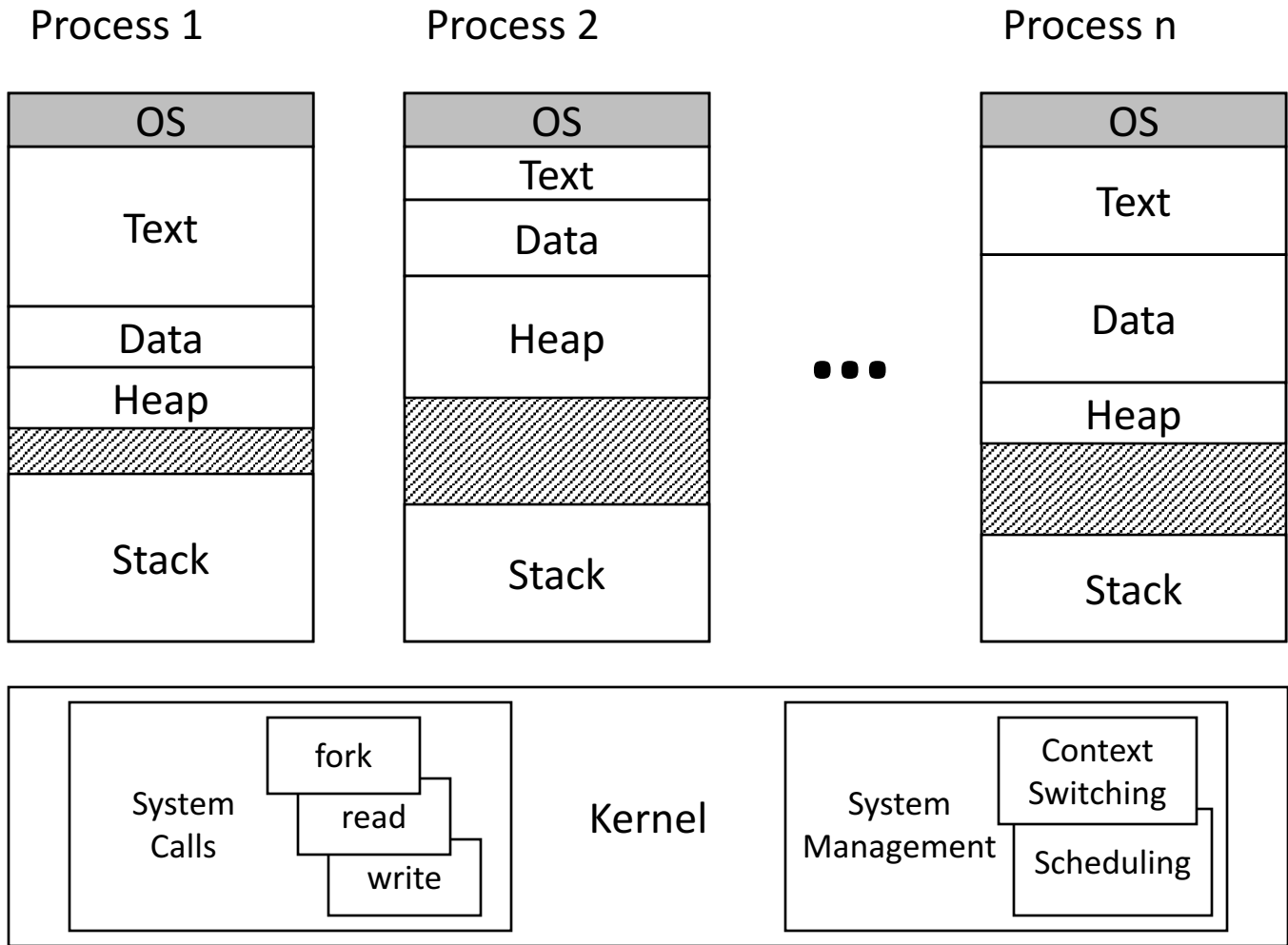


# Process and Kernel Model



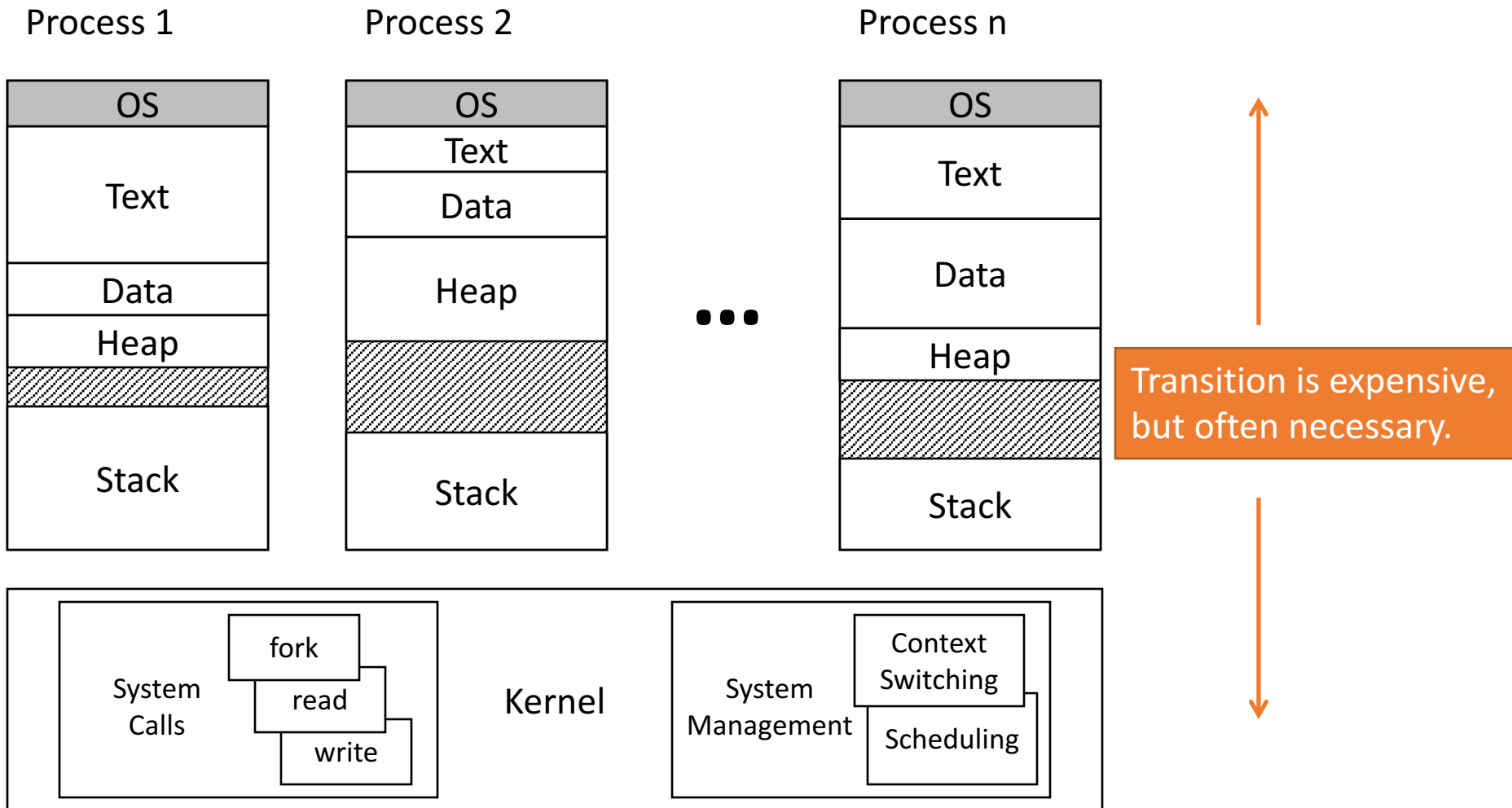
OS has control. It can context switch (and usually does at this point).

# Process and Kernel Model



OS returns control to a process (not usually the same one).

# Process and Kernel Model





# How Does Kernel Get Control

- To switch processes, kernel must get control
- Running process can give up control voluntarily
  - Process makes a blocking system call, e.g., `read ()`
  - To block, call `yield ()` to give up CPU
  - Control goes to kernel, which dispatches new process
- Or, CPU is forcibly taken away: preemption
  - While kernel is running, it sets a timer
  - When timer expires, interrupt is generated
  - Hardware forces control to go to kernel

# Up next...

- How we create/manage processes.
- How we provide the illusion of the same enormous memory space for all processes.