# Chapter 3: Introduction to Processes

CSCI 3753 Operating Systems

William Mortl, MS & Rick Han, PhD

Spring 2016

# Announcements

- Quiz #3 is up early for hackathon kids (deadline is still February 2$^{nd}$
- Programming Assignment #1 is due on February 2$^{nd}$… email me or TA's to use the 1 week extension (20% penalty)
- Computer architecture is not the focal point of this course, but we had to cover in order to talk about OS's... Boot loading started the REAL part of this course

# Recap…

- Booting
    1. POST – Power On Self Test
    2. BIOS – Basic Input Output System loads the 512 byte boot loader
    3. 512 byte boot loader loads secondary stage boot loader (GRUB, LILO)
    4. Second stage boot loader loads the OS kernel

# Bootstrapping the OS in PCs

- Multi-stage procedure:
  1. Power On Self Test (POST) from ROM
     - Check hardware, e.g. CPU and memory, to make sure it's OK
  2. BIOS (Basic Input/Output System) looks for a device to boot from…
     - May be prioritized to look for a USB flash drive or a CD/DVD-ROM drive before a hard disk drive
     - Can also boot from network
  3. BIOS finds a hard disk drive to boot from
     - Looks at Master Boot Record (MBR) in sector 0 of disk
     - Only 512 bytes long (Intel systems), contains primitive code for later stage loading and a partition table listing an active partition, or the location of the bootloader
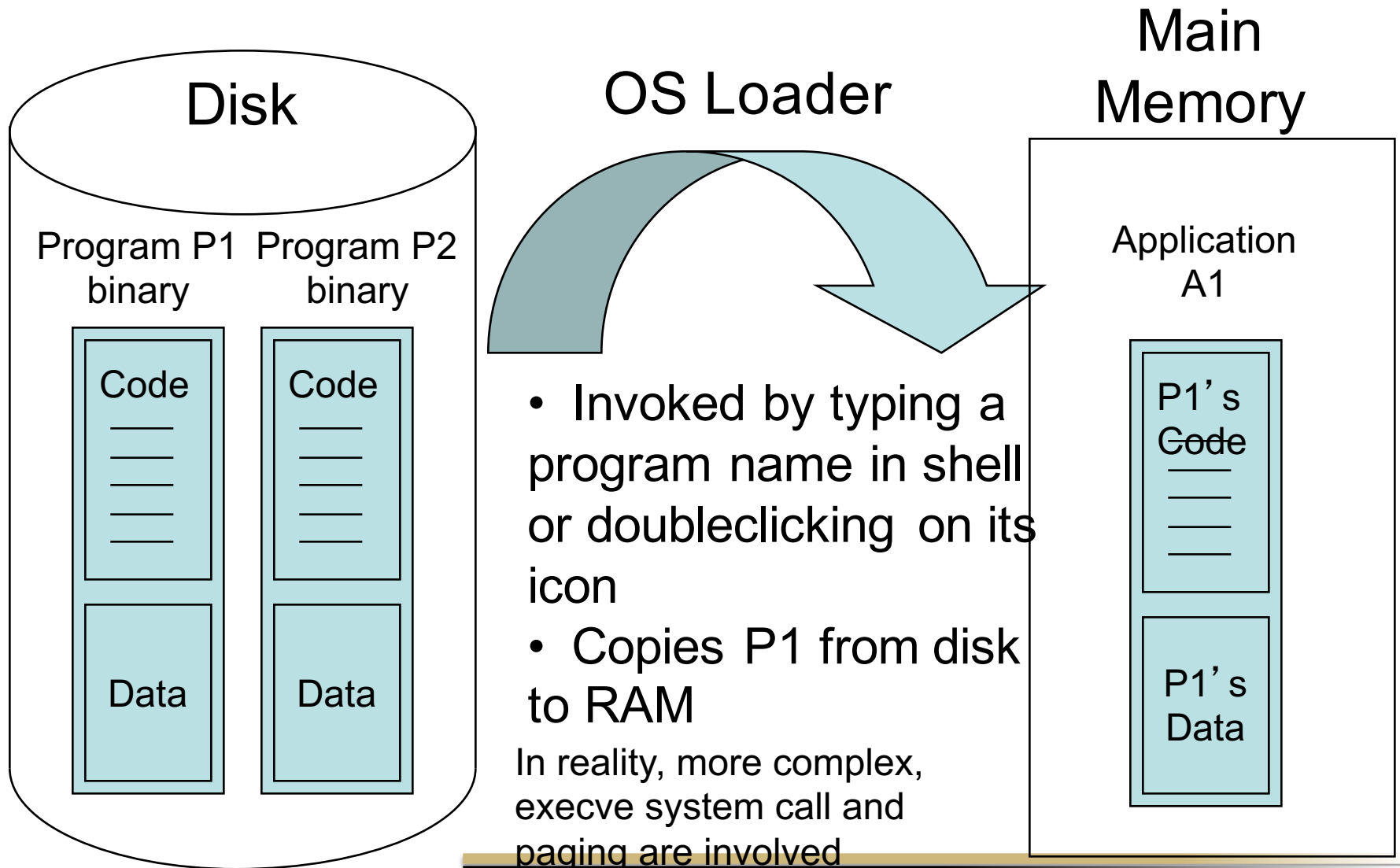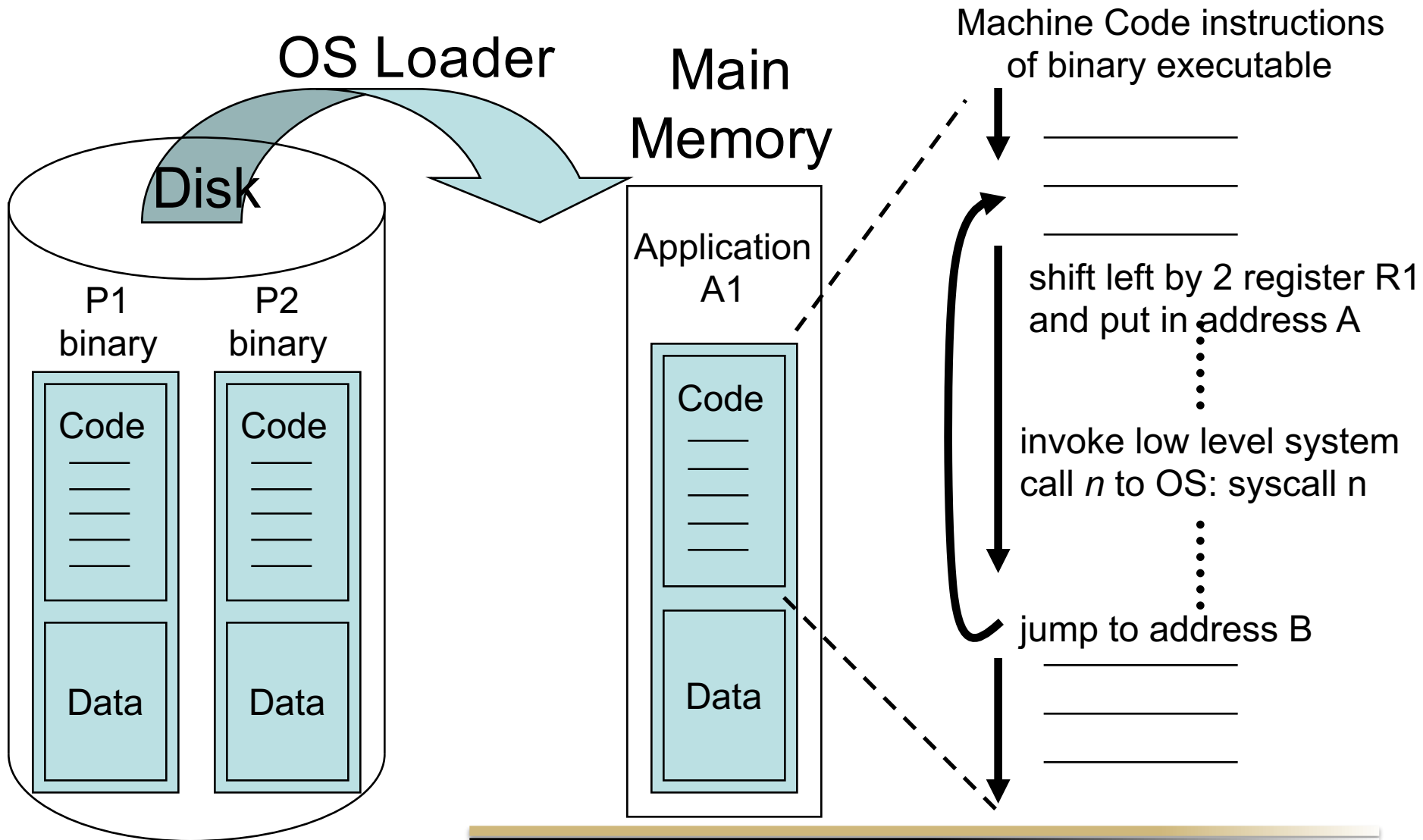
# Bootstrapping the OS in PCs

- Multi-stage procedure: (continued)
  4. Primitive loader then loads the secondary stage bootloader
     - Examples of this bootloader include LILO (Linux Loader), and GRUB (Grand Unified Bootloader)
     - Can select among multiple OS's (on different partitions) – i.e. dual booting
     - Once OS is selected, the bootloader goes to that OS's partition, finds the boot sector, and starts loading the OS's kernel

University of Colorado **Boulder**

# Loading a Program into Memory

**Disk**

**OS Loader**

**Main Memory**

Program P1 binary

Program P2 binary

Code
————
————
————
————
————

Code
————
————
————
————
————

Data

Data

- Invoked by typing a program name in shell or doubleclicking on its icon
- Copies P1 from disk to RAM

In reality, more complex, execve system call and paging are involved

Application A1

P1's ~~Code~~
————
————
————

P1's Data

# Loading and Executing a Program

OS Loader

Main Memory

Machine Code instructions of binary executable

Disk

P1 binary

P2 binary

Code

Code

Data

Data

Application A1

Code

Data

shift left by 2 register R1 and put in address A

invoke low level system call *n* to OS: syscall n

jump to address B

# Loading and Executing a Program



OS Loader

Disk

P1 binary

P2 binary

Code

Code

Data

Data

Main Memory

Application A1

P1's Code

P1's Data

Fetch Code and Data
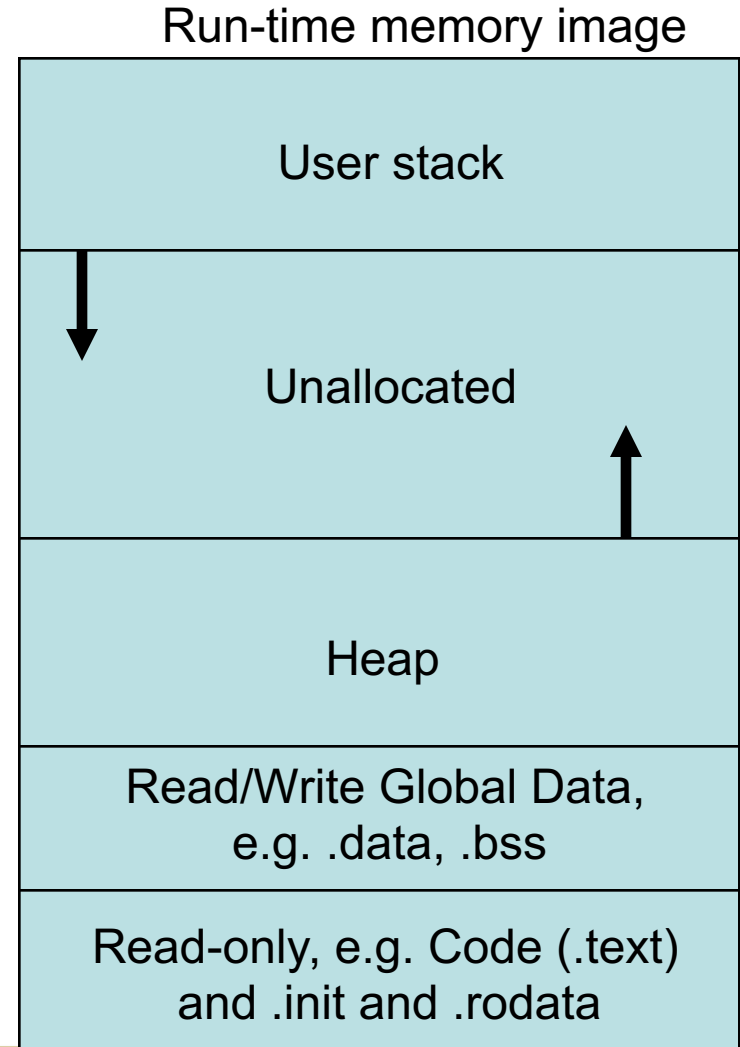
CPU

Program Counter (PC)

Registers

ALU

Write Data

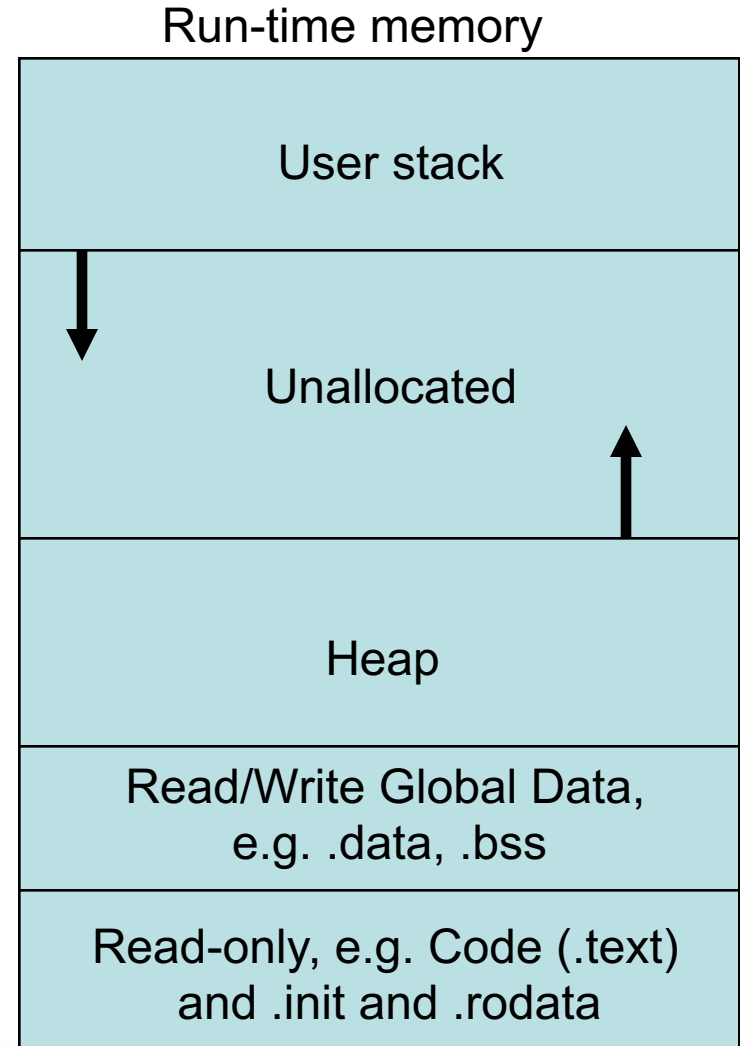University of Colorado **Boulder**

# Loading Executable Object Files

- When a program is loaded into RAM, it becomes an actively executing application

- The OS allocates a stack and heap to the app in addition to code and global data.
  - A call stack is for local variables
  - A heap is for dynamic variables, e.g. malloc()
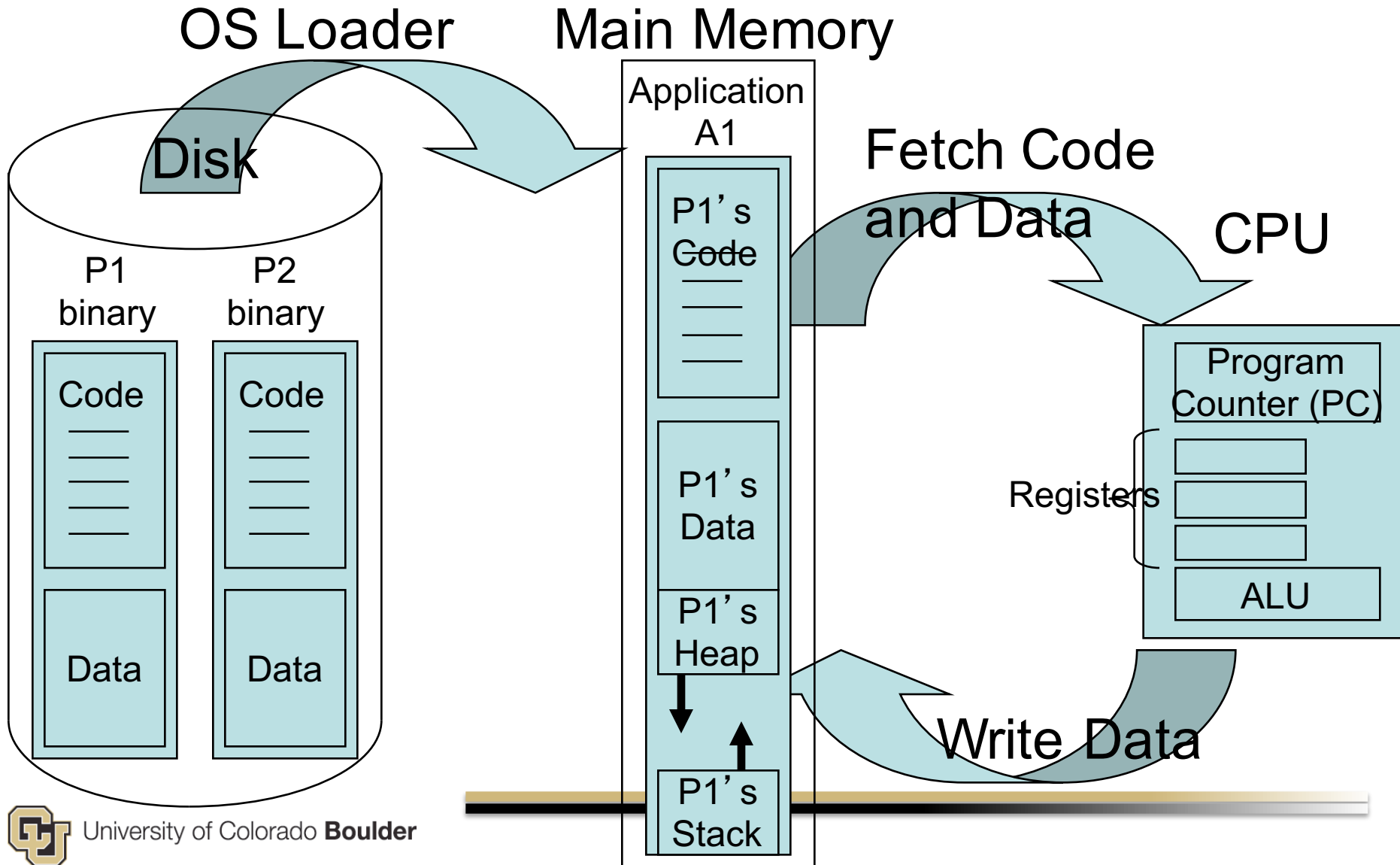  - Usually, stack grows downward from high memory, heap grows upward from low

Run-time memory image

| User stack |
|---|
| Unallocated |
| Heap |
| Read/Write Global Data, e.g. .data, .bss |
| Read-only, e.g. Code (.text) and .init and .rodata |

, but this is architecture-specific

# Running Executable Object Files

Run-time memory

- Stack contains local variables
  - As main() calls function f1, we allocate f1's local variables on the stack
  - If f1 calls f2, we allocate f2's variables on the stack below f1's, thereby growing the stack, etc…
  - When f2 is done, we deallocate f2's local variables, popping them off the stack, and return to f1
- Stack dynamically expands and contracts as program runs and different levels of nested functions are called
- Heap contains run-time variables/buffers
  - Obtained from malloc()
  - Program should free() the malloc'ed memory
- Heap can also expand and contract during program execution

| Run-time memory |
|---|
| User stack |
| Unallocated |
| Heap |
| Read/Write Global Data, e.g. .data, .bss |
| Read-only, e.g. Code (.text) and .init and .rodata |

University of Colorado **Boulder**

# Loading and Executing a Program – a more complete picture

OS Loader

Main Memory

Disk

P1 binary

P2 binary

Code

Code

Data

Data

Application A1

P1's Code

P1's Data

P1's Heap

P1's Stack

Fetch Code and Data

CPU

Program Counter (PC)

Registers

ALU

Write Data

# Multiple Applications + OS

**Main Memory**

**CPU Execution**

**Application A1**

Code

Data

Heap

Stack

**Application A2**

Code

Data

Heap

Stack

**Operating System**

Program Counter (PC)

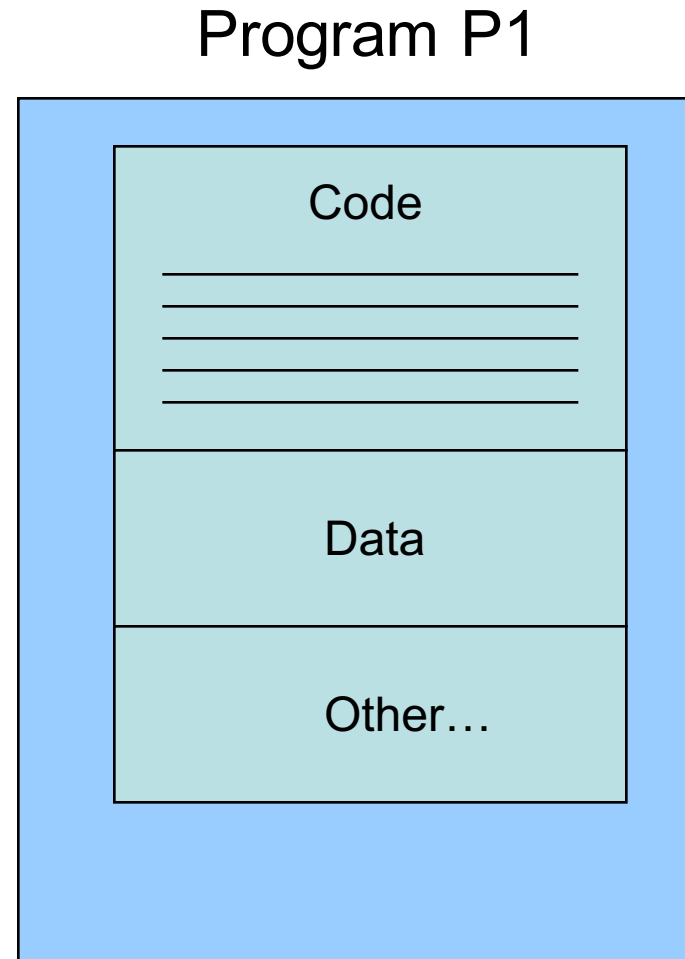Instruction Register

Stack Pointer
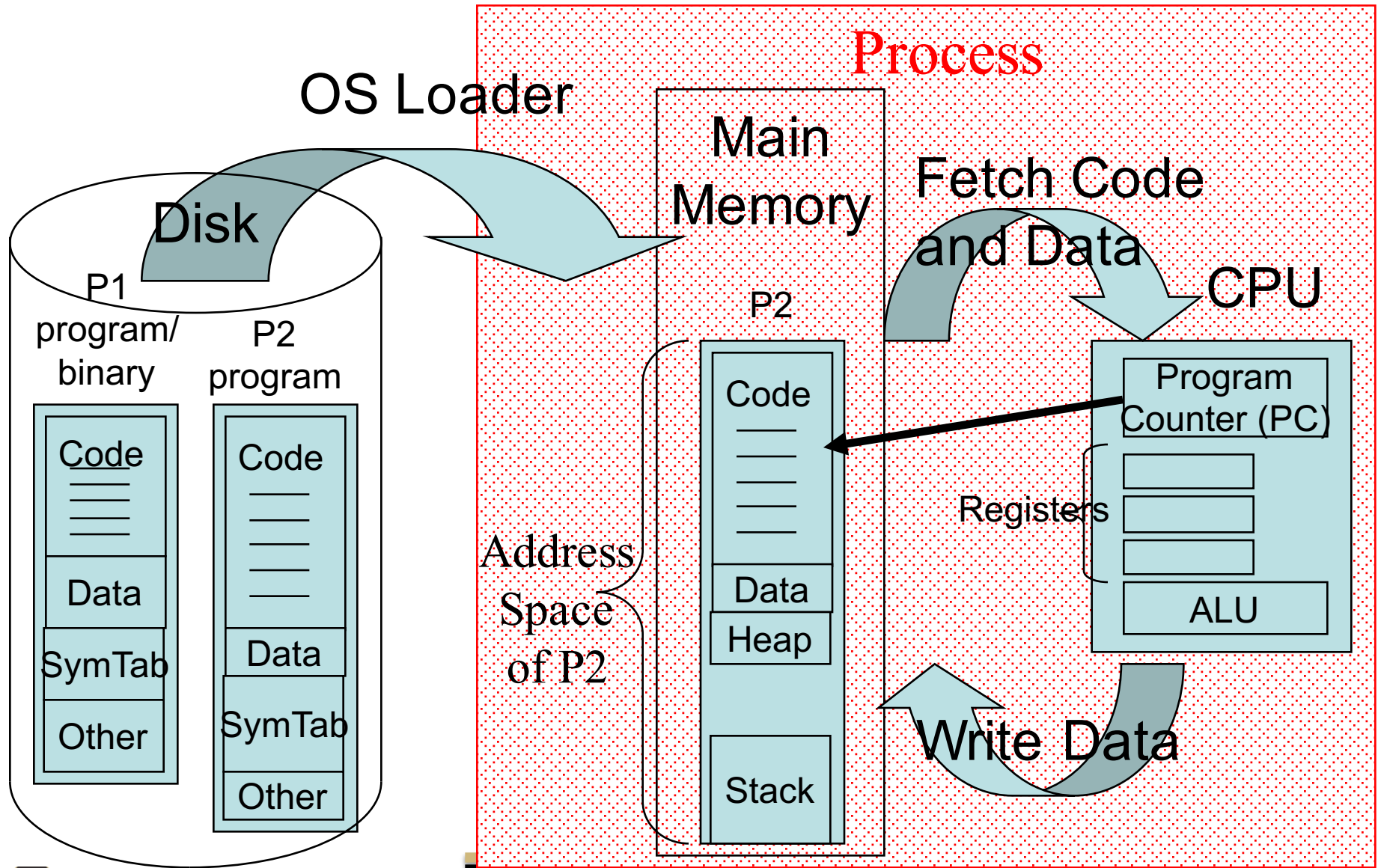
ALU

University of Colorado **Boulder**

# Chapter 3: What is a Process?

- A software *program* consist of a sequence of code instructions and data stored on disk
  - A program is a *passive* entity
- A *process* is a program *actively executing* from main memory within its *own address space*
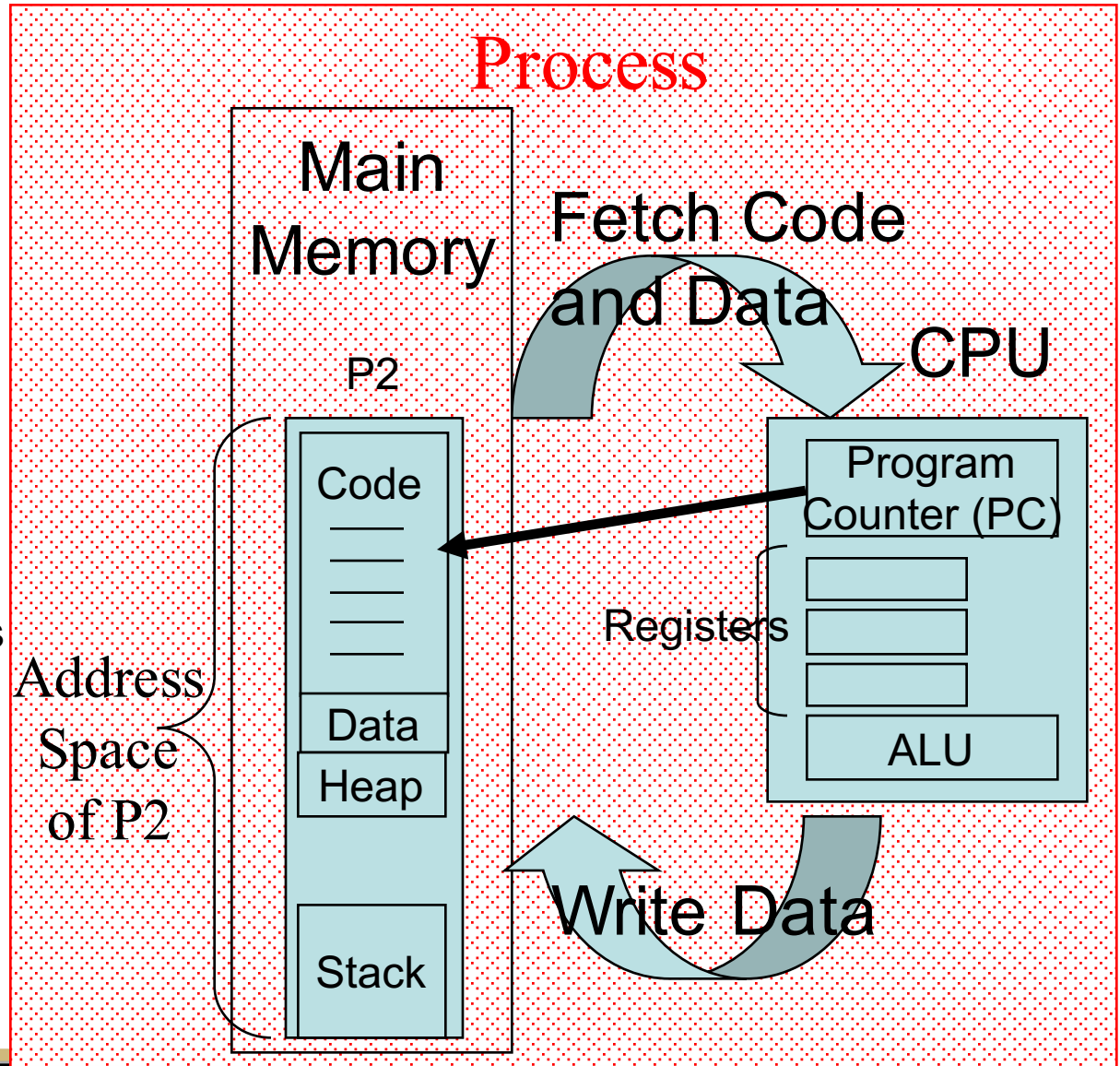
Program P1

| Code |
| --- |
| Data |
| Other… |

# What Is a Process? (2)



Process

OS Loader

Disk

P1 program/ binary

Code
Data
SymTab
Other

P2 program

Code
Data
SymTab
Other

Main Memory

P2

Address Space of P2

Code
Data
Heap
Stack

Fetch Code and Data

Write Data

CPU

Program Counter (PC)

Registers

ALU
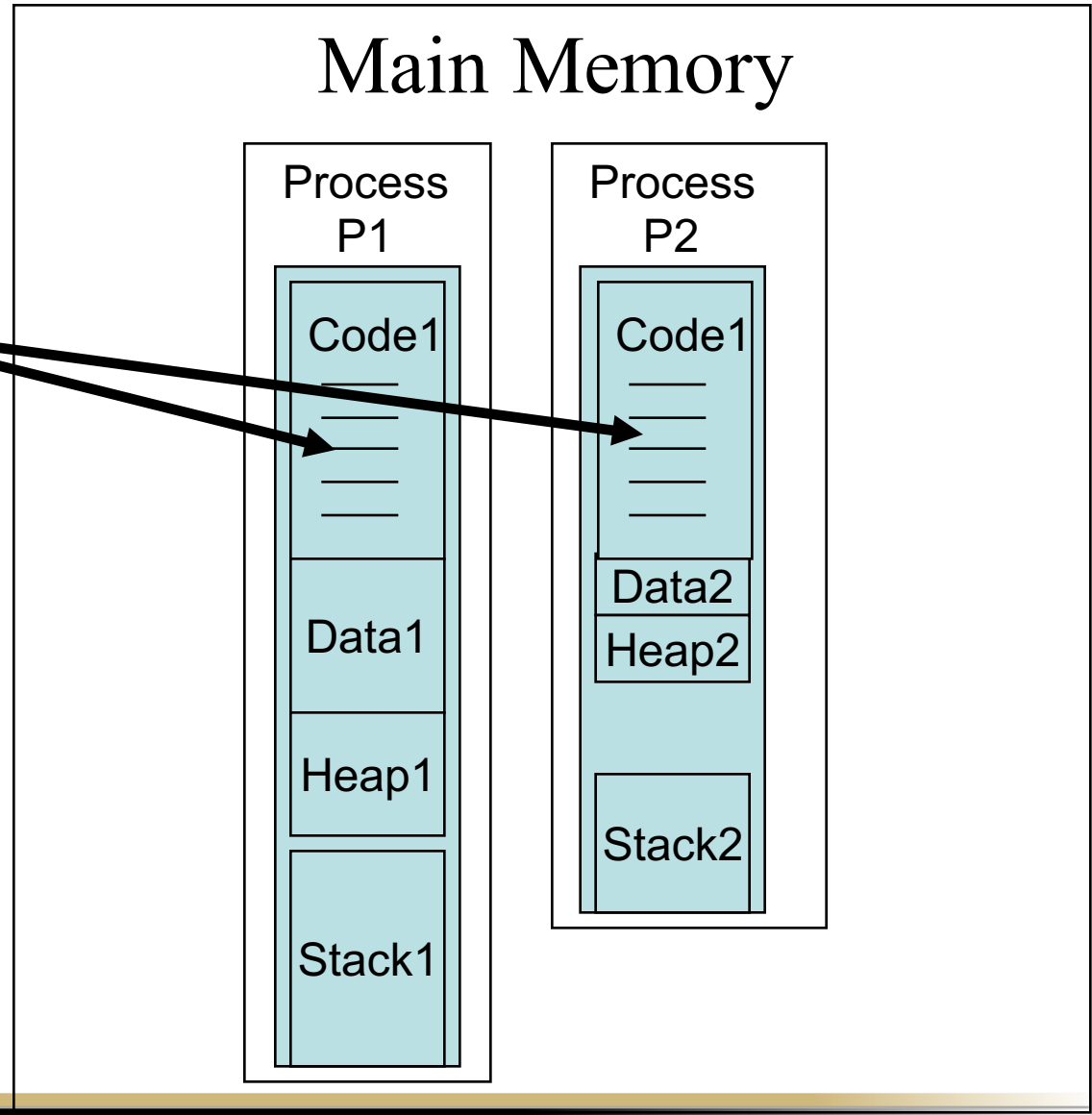
University of Colorado **Boulder**

# What is a Process? (3)

- A *process* is a program *actively executing* from main memory

  – has a Program Counter (PC) and execution state associated with it

    • CPU registers keep state

    • OS keeps process state in memory

    • it's alive!

  – Owns its own *address space*

    • a limited set of (virtual) addresses that can be accessed by the executing code



Process

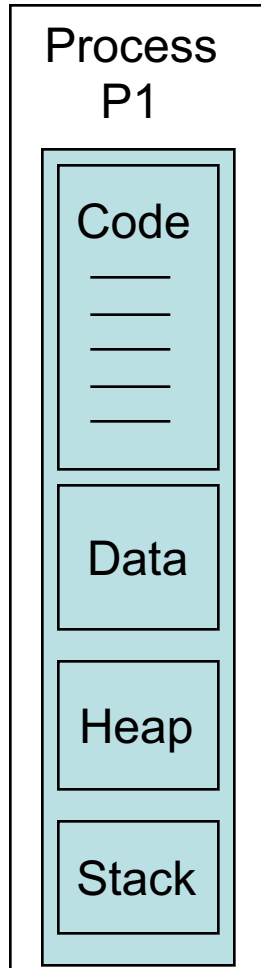Main Memory

Fetch Code and Data

CPU

P2

Code

Data

Heap

Stack

Address Space of P2

Program Counter (PC)

Registers

ALU

Write Data

# What is a Process? (4)

– 2 processes may execute the same program code, but they are considered *separate execution* sequences
  - e.g. two shell terminals

## Main Memory

### Process P1

Code1

Data1

Heap1

Stack1

### Process P2

Code1

Data2

Heap2

Stack2

# A Process Executes in its Own Address Space

## Main Memory

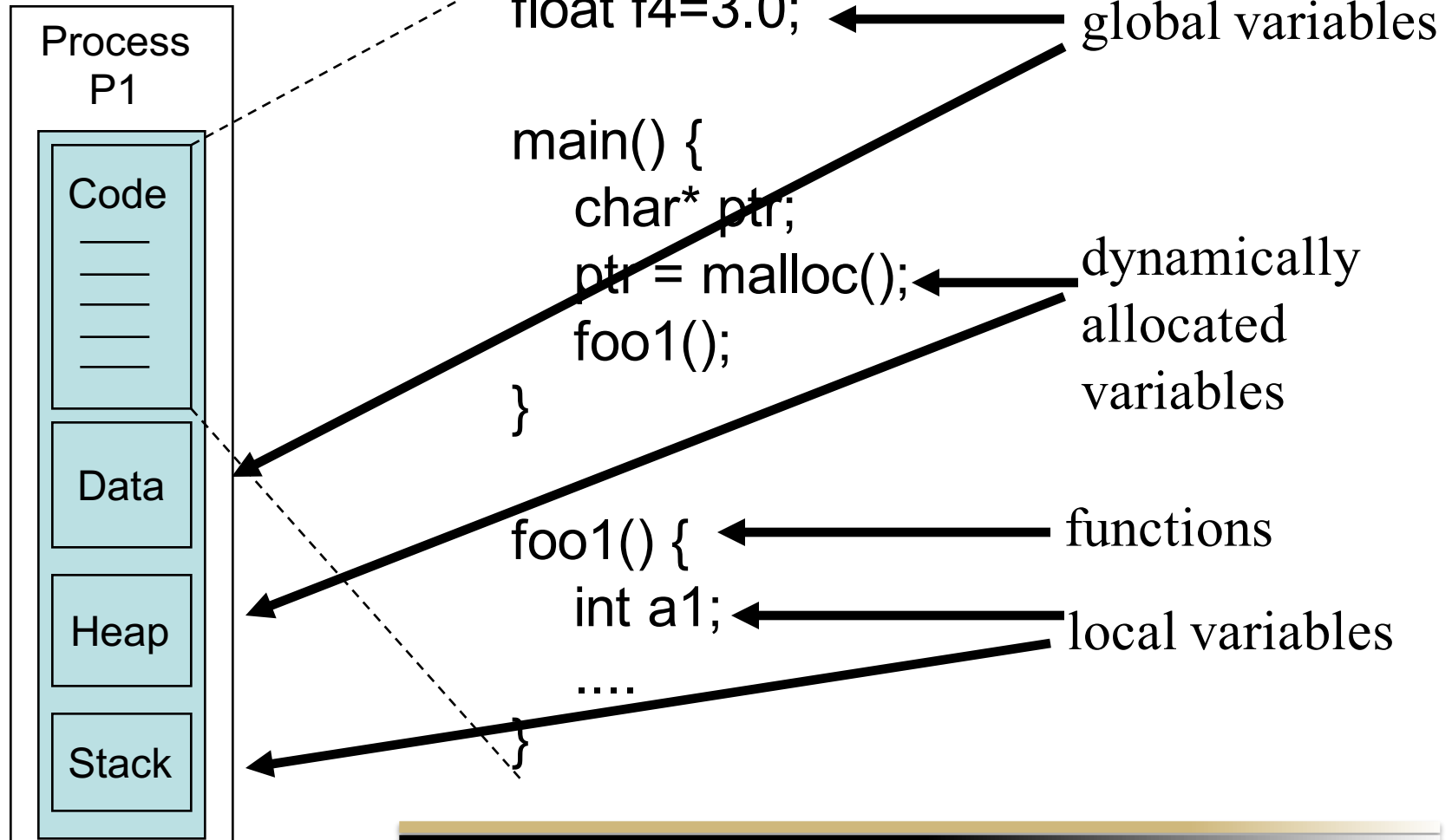| Process P1 |
|:---:|
| Code |
| ——— |
| ——— |
| ——— |
| Data |
| Heap |
| Stack |

- OS tries to provide the illusion or *abstraction* to the process that it executes on its own abstract machine
  - in its own subset of RAM, i.e. its own address space – achieved using virtual memory paging
  - on its own subset (time slice) of the CPU – achieved by preemptive multitasking
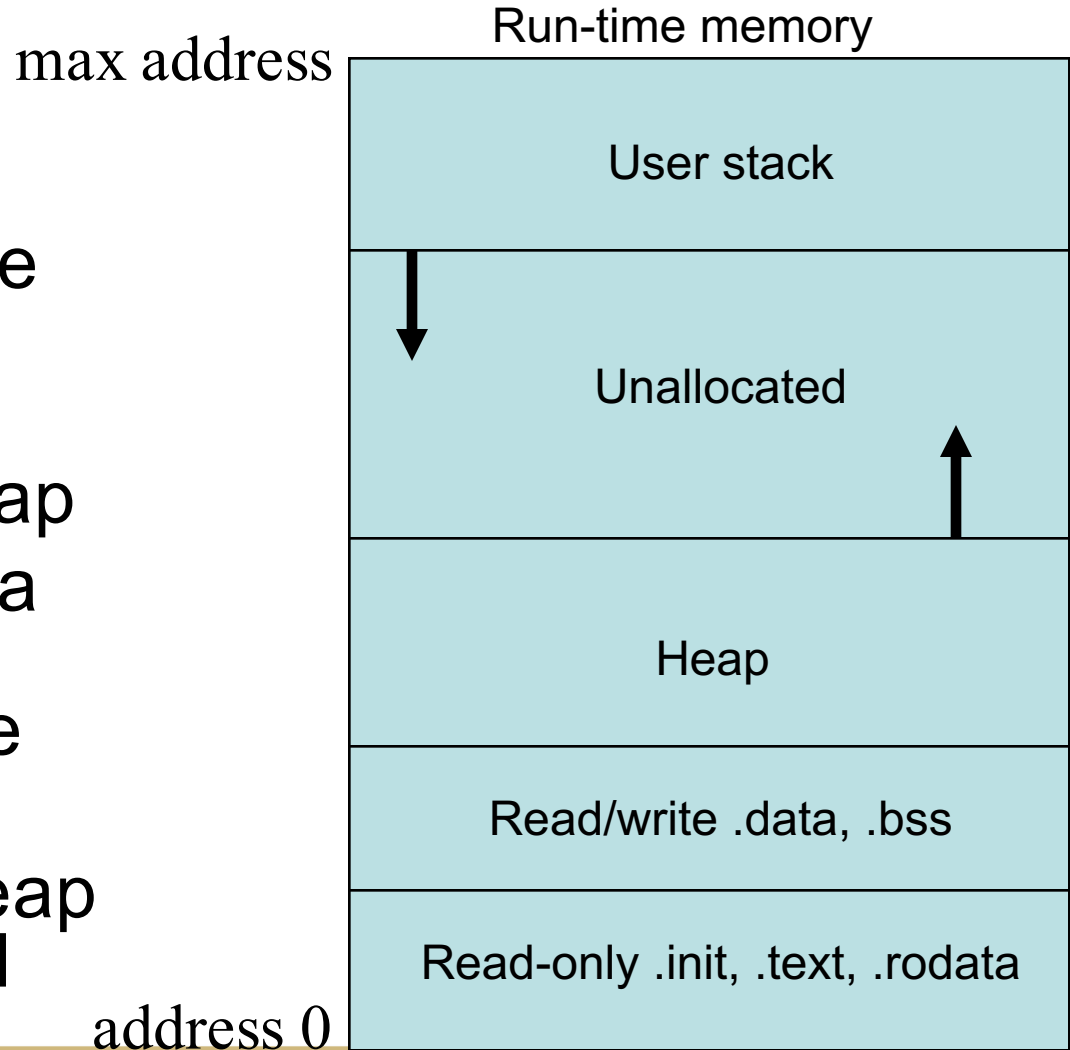
University of Colorado **Boulder**

# How is a Process Structured in Memory?

## Main Memory

Process P1

| Code |
|------|
| — |
| — |
| — |

Data

Heap

Stack

```
float f4=3.0;          ← global variables

main() {
    char* ptr;
    ptr = malloc();     ← dynamically
    foo1();                allocated
}                          variables


foo1() {               ← functions
    int a1;             ← local variables
    ....
}
```

# How is a Process Structured in Memory?

- Run-time memory image
- Essentially code, data, stack, and heap
- Code and data loaded from executable file
- Stack grows downward, heap grows upward

Run-time memory

max address

| |
|---|
| User stack |
| Unallocated |
| Heap |
| Read/write .data, .bss |
| Read-only .init, .text, .rodata |

address 0

# Applications and Processes

- Application = $\Sigma$ Processes$_i$

  - e.g. a server could be split into multiple processes, each one dedicated to a specific task (UI, computation, communication, etc.)

  - The Application's various processes talk to each other using Inter-Process Communication (IPC). We'll see various forms of IPC later.

University of Colorado **Boulder**