# Hardware, Modularity, and Virtualization
# CS 111
# Operating System Principles
# Peter Reiher

# Outline

- The relationship between hardware and operating systems
  - Processors
  - I/O devices
  - Memory
- Organizing systems via modularity
- Virtualization and operating systems

# Hardware and the Operating System

- One of the major roles of the operating system is to hide details of the hardware
  - Messy and difficult details
  - Specifics of particular pieces of hardware
  - Details that prevent safe operation of the computer
- OS abstractions are built on the hardware, at the bottom
  - Everything ultimately relies on hardware
- A major element of OS design concerns HW

# OS Abstractions and the Hardware

- Many important OS abstractions aren't supported directly by the hardware

- Virtual machines
  - There's one real machine

- Virtual memory
  - There's one set of physical memory
  - And it often isn't as big as even one process thinks it is

- Typical file abstractions

- Many others

- The OS works hard to make up the differences

# Processor Issues

- Execution mode
- Handling exceptions

# Execution Modes

- Modern CPUs can usually execute in two different modes:
  - User mode
  - Supervisor mode

- User mode is to run ordinary programs

- Supervisor mode is for OS use
  - To perform overall control
  - To perform unsafe operations on the behalf of processes

# User Mode

- Allows use of all the "normal" instructions
  - Load and store general registers from/to memory
  - Arithmetic, logical, test, compare, data copying
  - Branches and subroutine calls

- Able to address some subset of memory
  - Controlled by a Memory Management Unit

- <u>Not</u> able to perform privileged operations
  - I/O operations, update the MMU
  - Enable interrupts, enter supervisor mode

# Supervisor Mode

- Allows execution of privileged instructions
  - To perform I/O operations
  - Interrupt enable/disable/return, load PC
  - Instructions to change processor mode
- Can access privileged address spaces
  - Data structures inside the OS
  - Other process's address spaces
  - Can change and create address spaces
- May have alternate registers, alternate stack

# Controlling the Processor Mode

- Typically controlled by the *Processor Status Register* (AKA PS)

- PS also contains condition codes
  - Set by arithmetic/logical operations (0,+,-,ovflo)
  - Tested by conditional branch instructions

- Describes which interrupts are enabled

- May describe which address space to use

- May control other processor features/options
  - Word length, endian-ness, instruction set, ...

# How Do Modes Get Set?

- The computer boots up in supervisor mode
  - Used by bootstrap and OS to initialize the system

- Applications run in user mode
  - OS changes to user mode before running user code
    - User programs cannot do I/O, restricted address space
  - They can't arbitrarily enter supervisor mode
    - Because instructions to change the mode are privileged

- Re-entering supervisor mode is strictly controlled
  - Only in response to traps and interrupts

# So When Do We Go Back To Supervisor Mode?

- In several circumstances

- When a program needs OS services
  - Invokes system call that causes a trap
  - Which returns system to supervisor mode

- When an error occurs
  - Which requires OS to clean up

- When an interrupt occurs
  - Clock interrupts (often set by OS itself)
  - Device interrupts

# Asynchronous Exceptions and Handlers

- Most program errors can be handled "in-line"
  - Overflows may not be errors, noted in condition codes
  - If concerned, program can test for such conditions
- Some errors must interrupt program execution
  - Unable to execute last instruction (e.g., illegal op)
  - Last instruction produced non-results (e.g., divide by zero)
  - Problem unrelated to program (e.g., power failure)
- Most computers use traps to inform OS of problems
  - Define a well specified list of all possible exceptions
  - Provide means for OS to associate handler with each
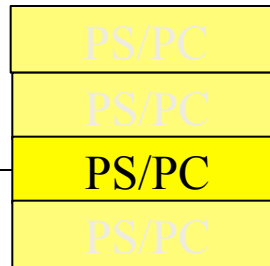
# Control of Supervisor Mode Transitions

- All user-to-supervisor changes via traps/interrupts
  - These happen at unpredictable times
- There is a designated handler for each trap/interrupt
  - Its address is stored in a trap/interrupt vector table managed by the OS
- Ordinary programs can't access these vectors
- The OS controls all supervisor mode transitions
  - By carefully controlling all of the trap/interrupt "gateways"
- Traps/interrupts can happen while in supervisor mode
  - Their handling is similar, but a little easier

# Software Trap Handling

Application Program

| instr ; | instr ; | instr ; | instr ; | instr ; | instr ; |

user mode

supervisor mode

PS/PC
PS/PC
**PS/PC**
PS/PC

TRAP vector table

1st level trap handler
(saves registers and
selects 2nd level handler)

return to
user mode

2nd level handler
(actually deals
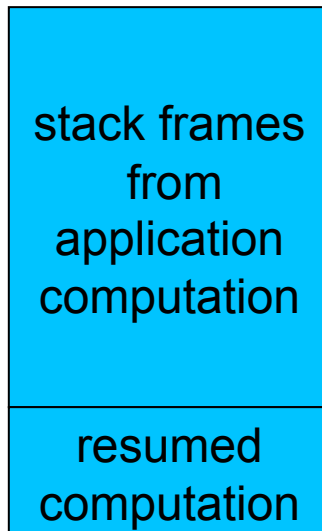with the problem)

# Dealing With the Cause of a Trap

- Some exceptions are handled by the OS
  - For example, page faults, alignment, floating point emulation
  - OS simulates expected behavior and returns

- Some exceptions may be fatal to running task
  - E.g. zero divide, illegal instruction, invalid address
  - OS reflects the failure back to the running process

- Some exceptions may be fatal to the system
  - E.g. power failure, cache parity, stack violation
  - OS cleanly shuts down the affected hardware
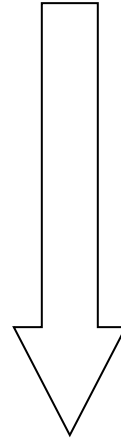
# Returning To User Mode

- Return is opposite of interrupt/trap entry
  - 2nd level handler returns to 1st level handler
  - 1st level handler restores all registers from stack
  - Use privileged return instruction to restore PC/PS
  - Resume user-mode execution after trapped instruction
- Saved registers can be changed before return
  - To set entry point for newly loaded programs
  - To deliver signals to user-mode processes
  - To set return codes from system calls
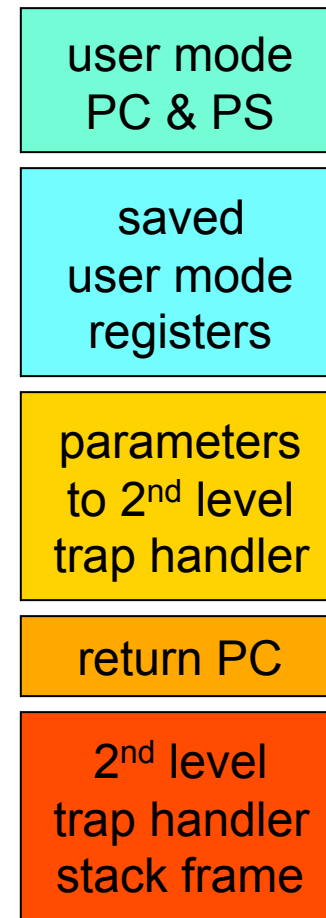
# Stacking and Unstacking a Trap

User-mode Stack

**TRAP!**

Supervisor-mode Stack

| stack frames from application computation |
| resumed computation |

↓

direction
of growth

| user mode PC & PS |
| saved user mode registers |
| parameters to 2nd level trap handler |
| return PC |
| 2nd level trap handler stack frame |

# I/O Architecture

- I/O is:
  - Varied
  - Complex
  - Error prone

- Bad place for the user to be wandering around

- The operating system must make I/O friendlier

- Oriented around handling many different *devices* via *busses* using *device drivers*
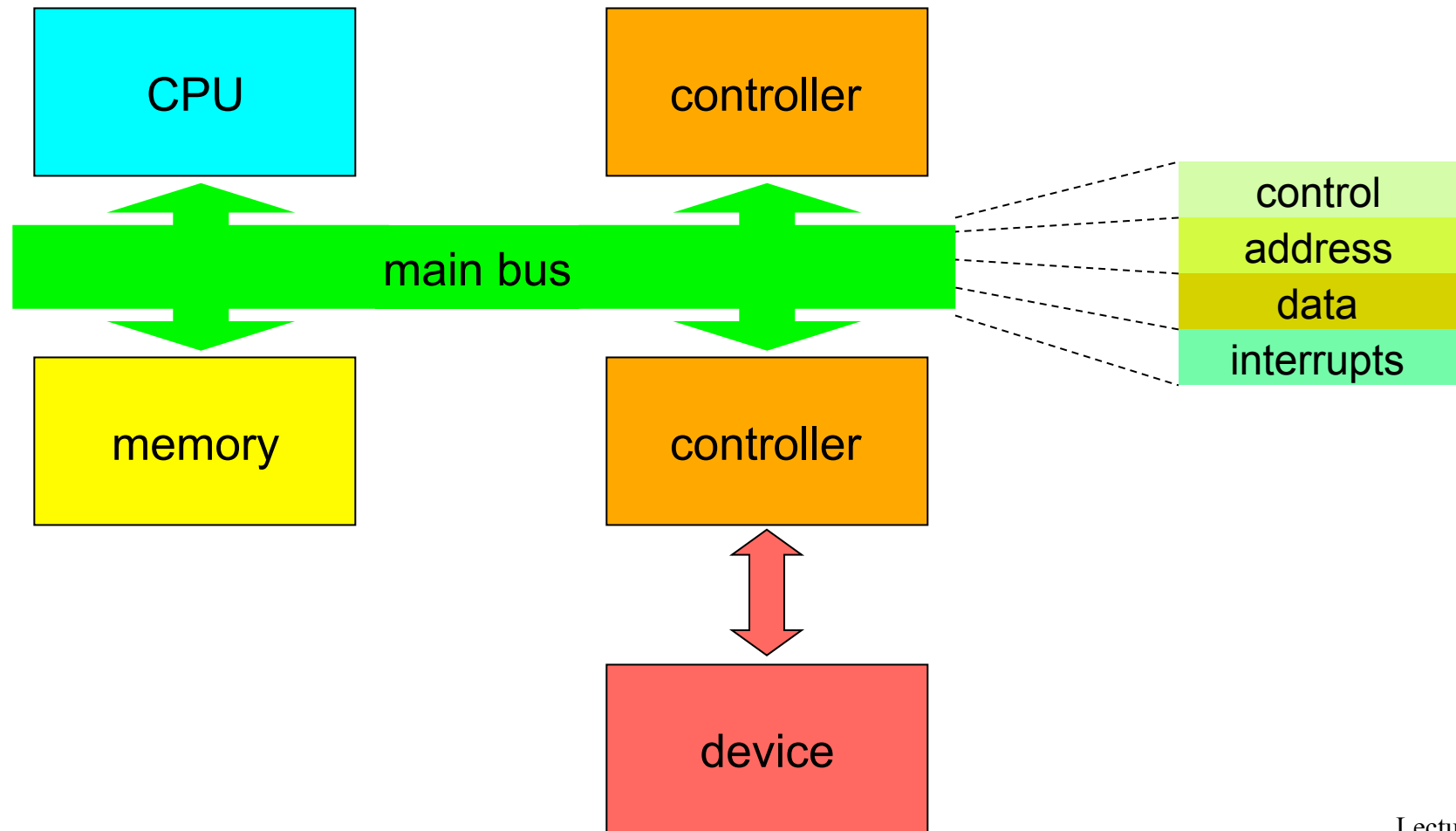
# Sequential vs. Random Access Devices

- Sequential access devices
  - Byte/block N must be read/written before byte/block N+1
  - May be read/write once, or may be rewindable
  - Examples: magnetic tape, printer, keyboard

- Random access devices
  - Possible to directly request any desired byte/block
  - Getting to that byte/block may or may not be instantaneous
  - Examples: memory, magnetic disk, graphics adaptor

- They are used very differently
  - Requiring different handling by the OS

# Busses

- Something has to hook together the components of a computer

  – The CPU, memory, various devices

- Allowing data to flow between them

- That is a *bus*

- A type of communication link abstraction

# A Simple Bus

CPU

controller

memory

controller

main bus

control
address
data
interrupts

device

# Devices and Controllers

- Device controllers connect a device to a bus
  - Communicate control operations to device
  - Relay status information back to the bus, manage DMA, generate device interrupts

- Device controllers export registers to the bus
  - Writing into registers controls device or sends data
  - Reading from registers obtains data/status

- Register access method varies with CPU type
  - May use special instructions (e.g., x86 IN/OUT)
  - May be mapped onto bus just like memory

# Direct Polled I/O

- Method of accessing devices via direct CPU control
  - CPU transfers data to/from device controller registers
  - Transfers are typically one byte or word at a time
  - May be accomplished with normal or I/O instructions
- CPU polls device until it is ready for data transfer
  - Received data is available to be read
  - Previously initiated write operations are completed

+ Very easy to implement (both hardware and software)

– CPU intensive, wastes CPU cycles on I/O control

– Leaves devices idle waiting for CPU when other tasks running

# Direct Memory Access

- Essentially, use the bus without CPU control
  - Move data between memory and device controller
- Bus facilitates data flow in all directions between:
  - CPU, memory, and device controllers
- CPU can be the bus-master
  - Initiating data transfers with memory, device controllers
- But device controllers can also master the bus
  - CPU instructs controller what transfer is desired
  - Device controller does transfer w/o CPU assistance
  - Device controller generates interrupt at end of transfer
- Interrupts tell CPU when DMA is done

# Memory Issues

- Different types of memory handled in different ways

- Cache memory usually handled mostly by hardware
  - Often OS not involved at all

- RAM requires very special handling
  - To be discussed in detail later

- Disks and flash drives treated as devices
  - But often with extra OS support

# Modularity

- Most useful abstractions an OS wants to offer can't be directly realized by hardware

- Modularity is one technique the OS uses to provide better abstractions

- Divide up the overall system you want into well-defined communicating pieces

- Critical issues:
  - Which pieces to treat as modules
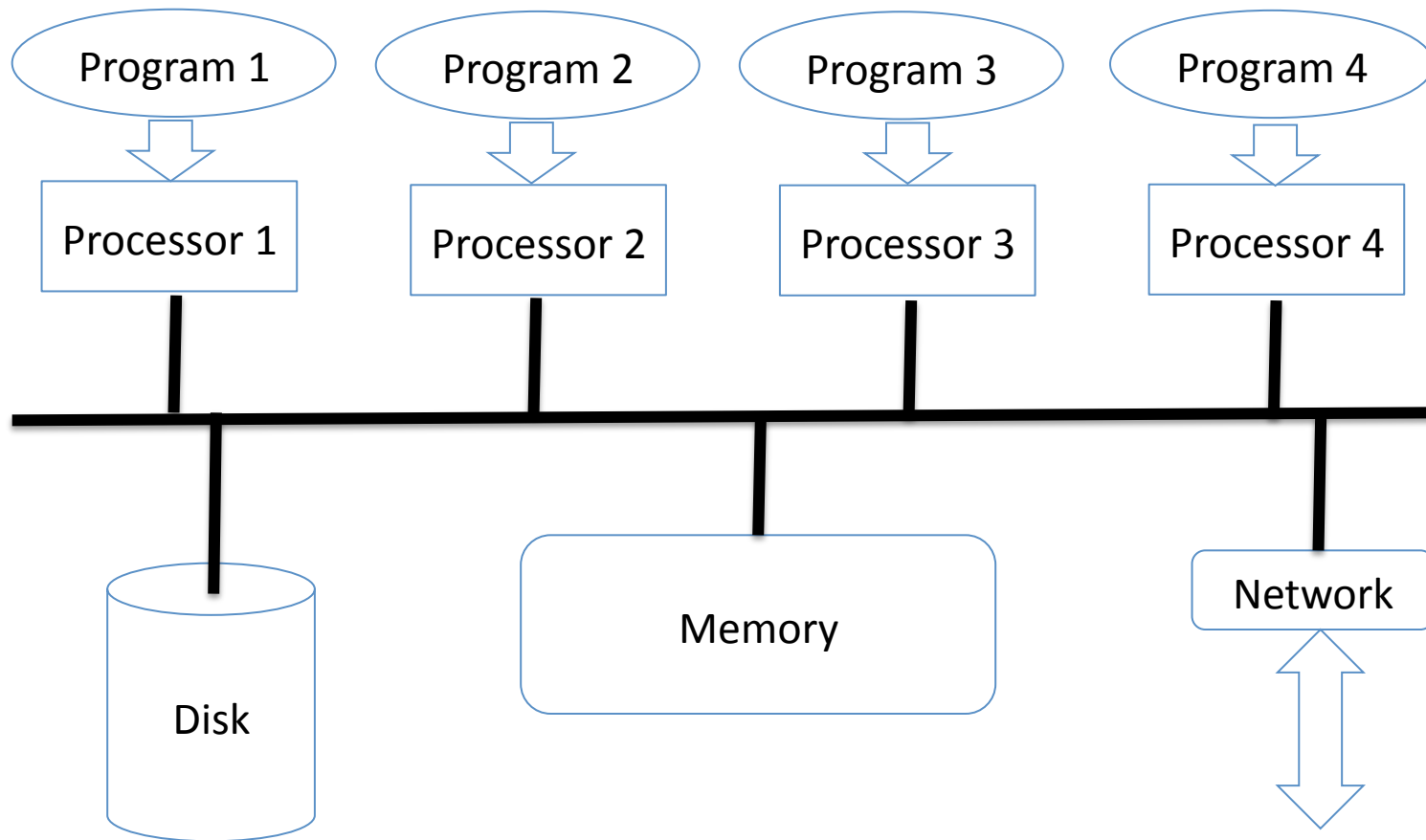  - How to organize the modules
  - Interfaces to modules

# What Does An OS Do?

- At minimum, it enables one to run applications
  - Preferably several on the same machine
  - Preferably several at the same time

- At abstract level, what do we need to do that?
  - Interpreters (to run the code)
  - Memory (to store the code and data)
  - Communications links (to communicate between apps and pieces of the system)

- This suggests the kinds of modules we'll need

# Starting Simple

- We want to run multiple programs
  - Without interference between them
  - Protecting one from the faults of another
- We've got a multicore processor to do so
  - More cores than programs
- We have RAM, a bus, a disk, other simple devices
- What abstractions should we build to ensure that things go well?

# A Simple System



**A machine boundary**

# Exploiting Modularity

- We'll obviously have several SW elements to support the different user programs

- Desirable for each to be modular and self-contained

  – With controlled interactions

- Gives cleaner organization

- Easier to prevent problems from spreading

- Easier to understand what's going on

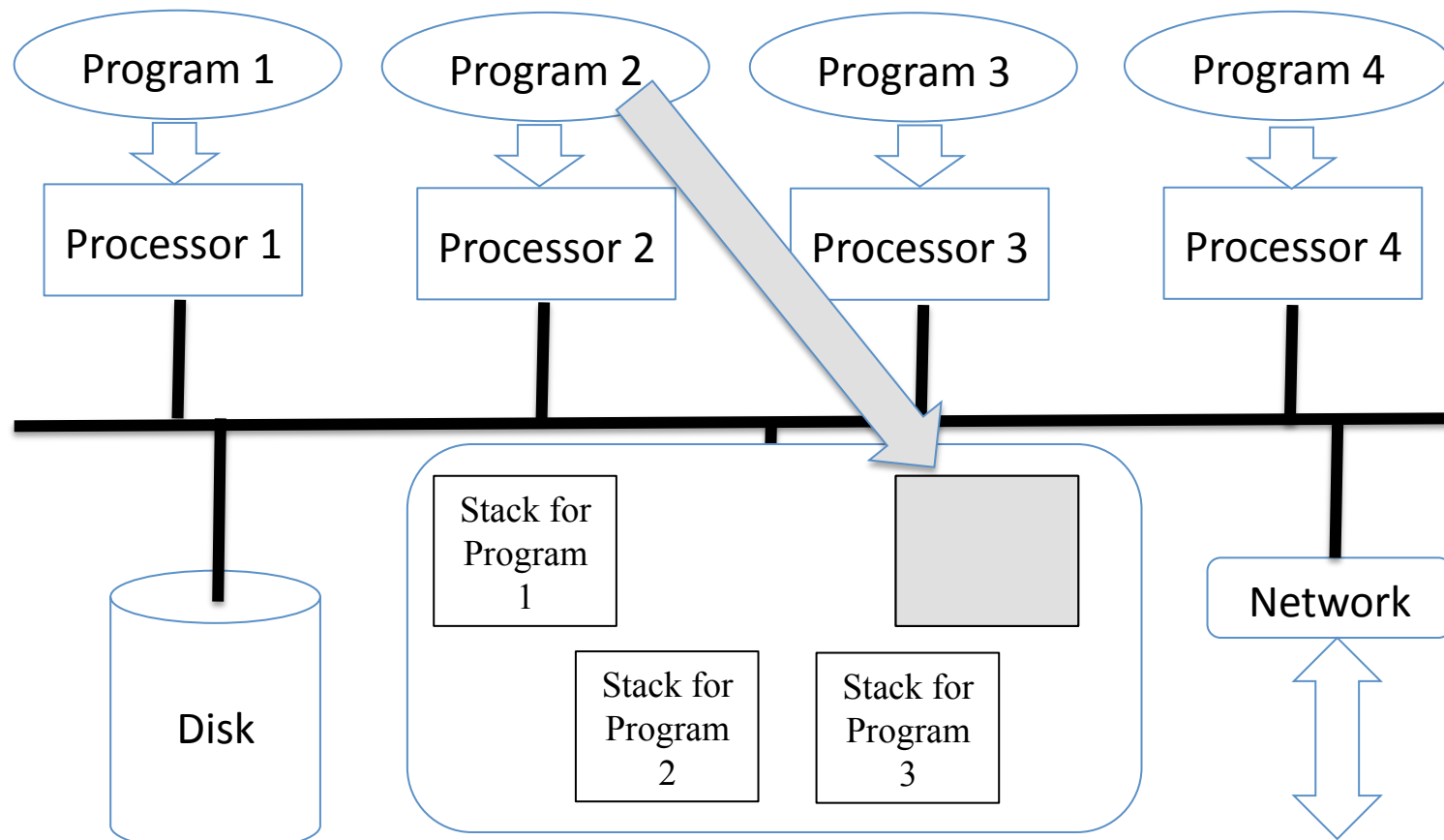- Easier to control each program's behavior

# Subroutine Modularity

- Why not just organize the system as a set of subroutines?
  - All in the same address space
    - A simplifying assumption
    - Allowing easy in-memory communication
- System subroutines call user program subroutines as needed
  - And vice versa
- *Soft modularity*
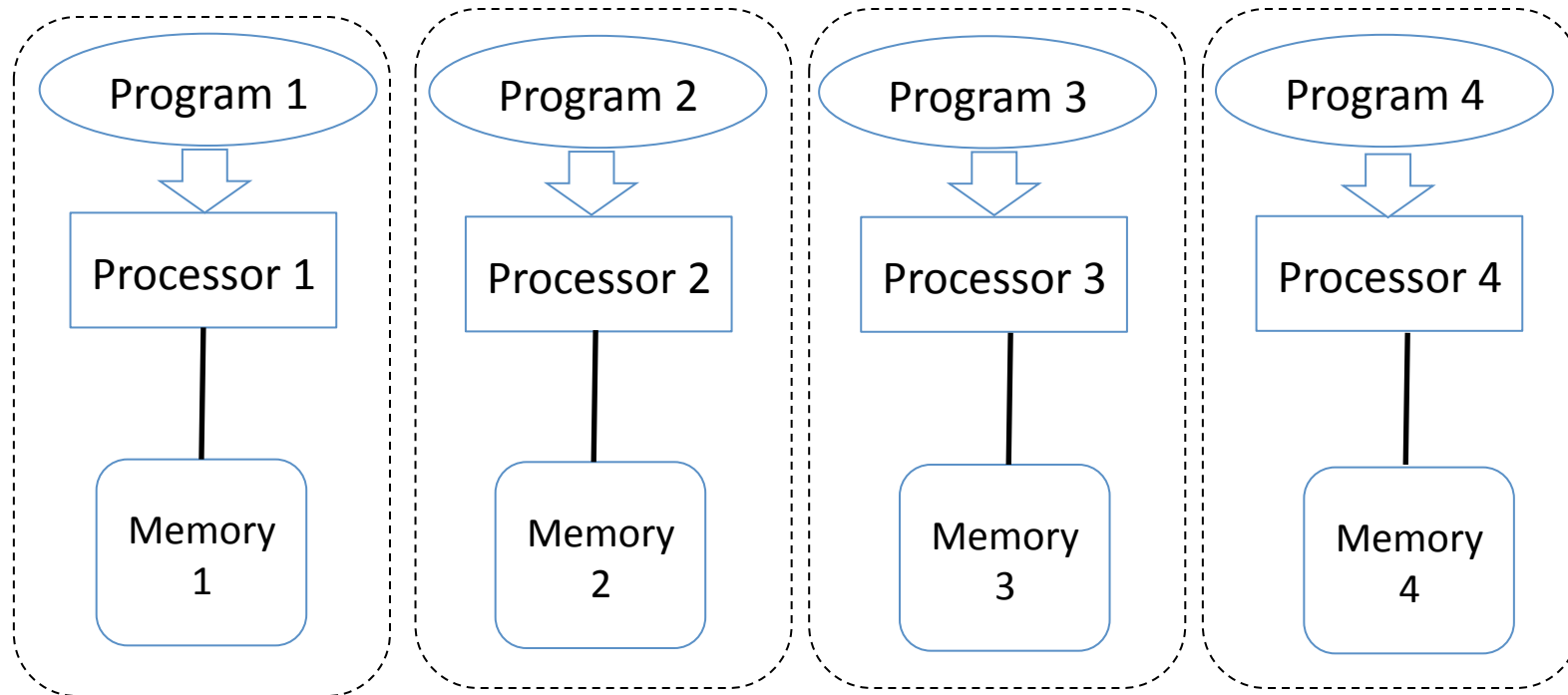
# How Would This Work?

- Each program is a self-contained set of subroutines
  - Subroutines in the program call each other
  - But not subroutines in other programs

- Shared services offered by other subroutines
  - Which any program can call

- Perhaps some "master routine" that calls subroutines in the various programs

- Soft because no OS HW/SW enforces modularity
  - Important resources (like the stack) are shared
  - Only proper program behavior protects one program from the mistakes of another

# Illustrating the Problem

Program 1    Program 2    Program 3    Program 4

Processor 1    Processor 2    Processor 3    Processor 4

Disk

Stack for Program 1

Stack for Program 2

Stack for Program 3

Network

Now Program 4 is in trouble
Even though it did nothing wrong itself

# Hardening the Modularity

| Program 1 | Program 2 | Program 3 | Program 4 |
|---|---|---|---|
| Processor 1 | Processor 2 | Processor 3 | Processor 4 |
| Memory 1 | Memory 2 | Memory 3 | Memory 4 |

**Four separate machines**

**Perhaps in very different places**

**Each program has its own machine**

# System Services In This Model

- Some activities are local to each program
- Other services are intended to be shared
  - Like a file system
- This functionality can be provided by a client/server model
- The system services are provided by the server
- The user programs are clients
- The client sends message to server to get help
- OS uses HW/SW to enforce boundaries

# Benefits of Hard Modularity

- With hard modularity, something beyond good behavior enforces module boundaries

- Here, the physical boundaries of the machine

- A client machine literally cannot touch the memory of the server

  – Or of another client machine

- No error or attack can change that

  – Though flaws in the server can cause problems
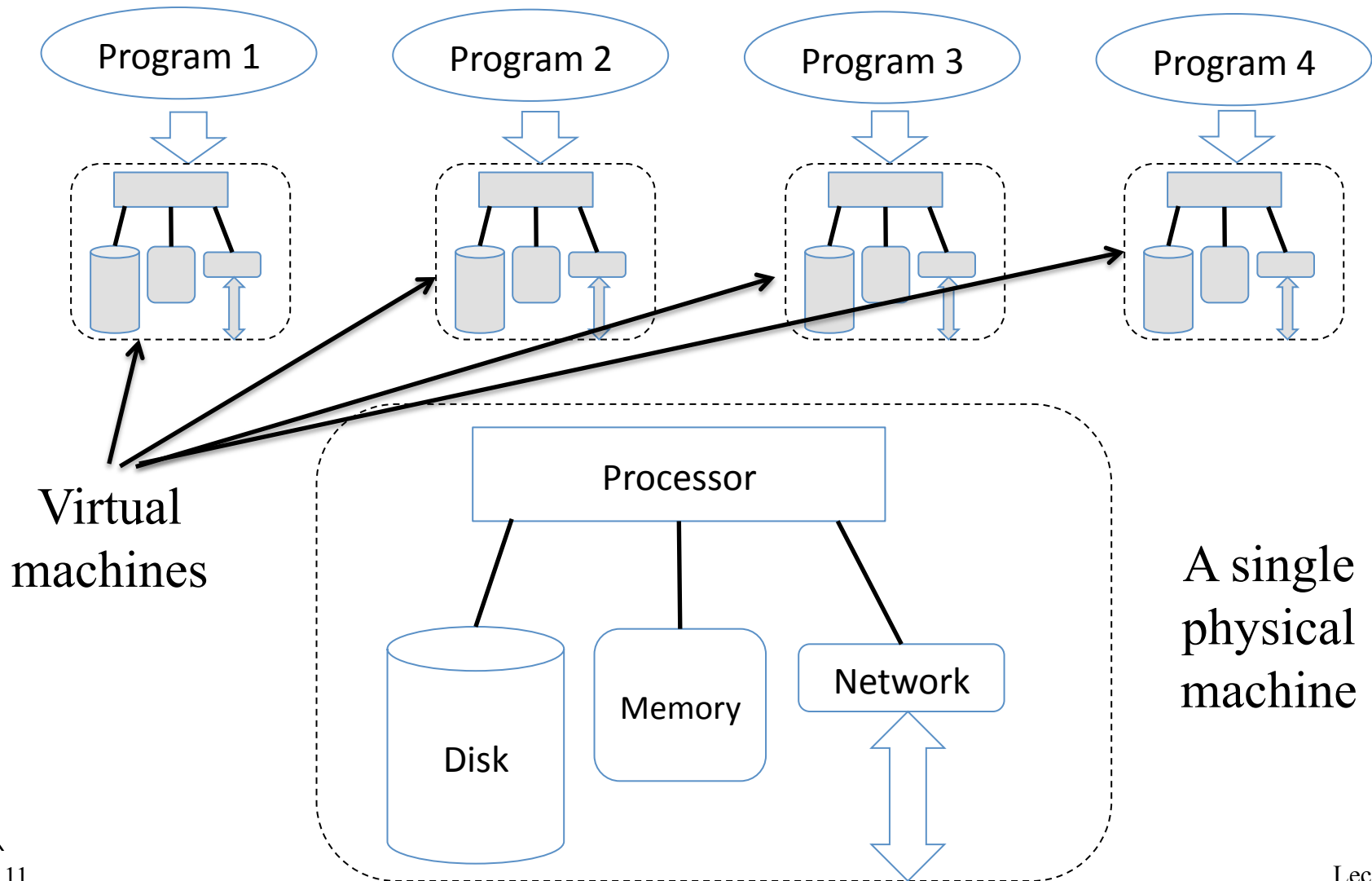
- Provides stronger guarantees all around

# Downsides of Hard Modularity

- The hard boundaries prevent low-cost optimizations

- In client/server organizations, doing anything with another program requires messages
  - Inherently more expensive than memory accesses

- If the boundary sits between components requiring fast interactions, possibly very bad

- Must either give programs pieces of resources or time multiplex use of resources
  - More complexity to do this right

# Virtualization

- Provide the illusion of a complete resource to each program that uses it
  - Hide hard modularity's time/space divisions
- Possible to provide an entire virtual machine per process
- Use shared hardware to instantiate the various virtual devices or machines
- System software (i.e., the operating system) and perhaps special hardware handle it

# The Virtualization Concept

Program 1

Program 2

Program 3

Program 4

Virtual machines

Processor

Disk

Memory

Network

A single physical machine

# The Trick in Virtualization

- All the virtual machines share the same physical hardware

- But each thinks it has its own machine

- Must be sure that one virtual machine doesn't affect behavior of the others
  - Intentionally or accidentally

- With the least possible performance penalty
  - Given that there will be a penalty merely for sharing at all

# Performance and Virtualization

- To achieve good performance, can't run many instructions "virtualized"
  - Most instructions must go directly to the processor
  - Rather than be mapped into multiple instructions via virtualization
- Similarly for access to other HW
  - Can't afford to put lots of virtualization SW in the usual path
- The trick is to virtualize the minimal set of accesses

# Abstractions for Virtualizing Computers

- Some kind of interpreter abstraction

  – A *thread*

- Some kind of communications abstraction

  – *Bounded buffers*

- Some kind of memory abstraction

  – *Virtual memory*

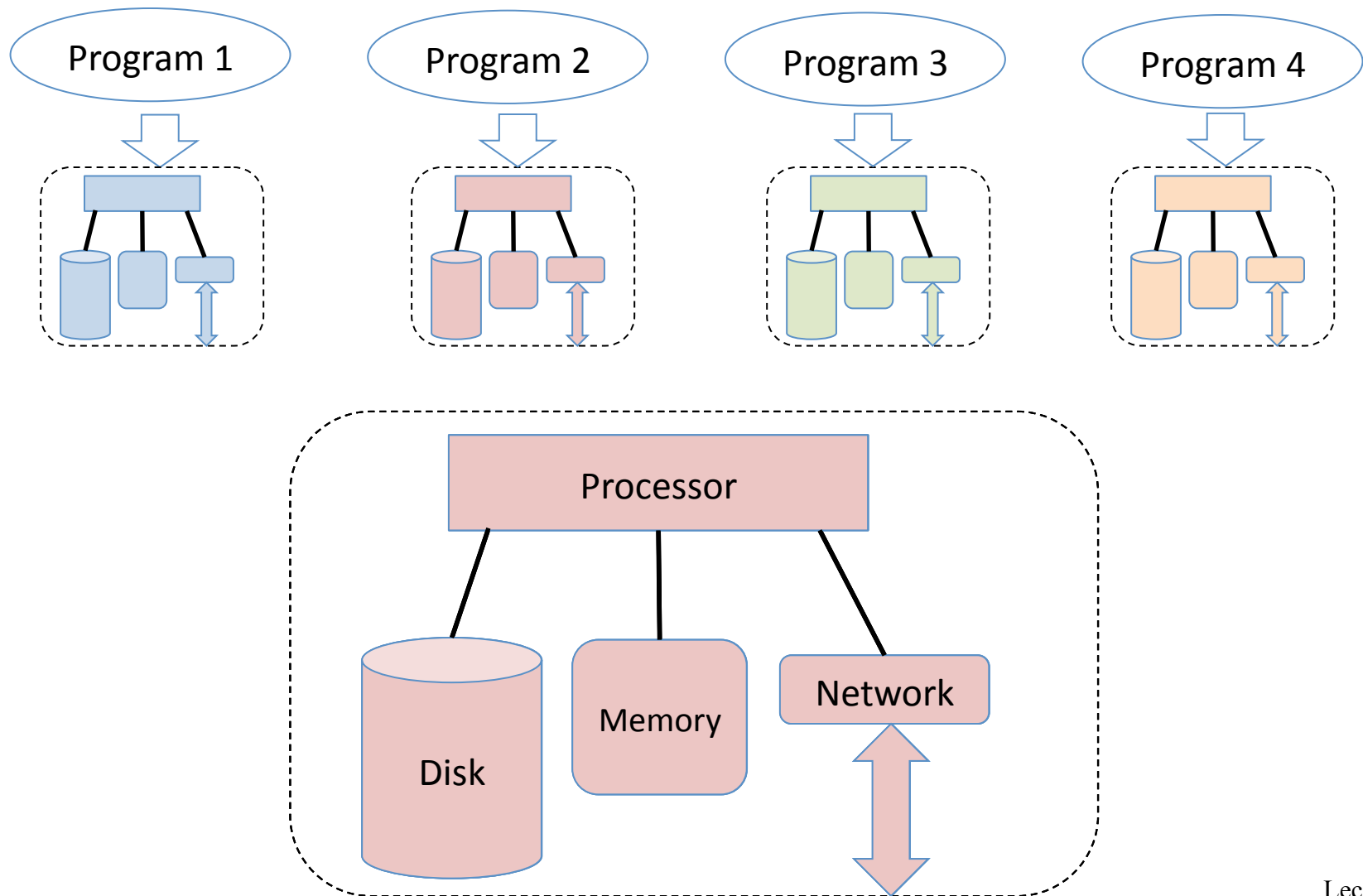- For a virtualized architecture, the operating system provides these kinds of abstractions

# Threads

- Encapsulates the state of a running computation

- So what does it need?

  – Something that describes what computation is to be performed

  – Something that describes where it is in the computation

  – Something that maintains the state of the computation's data

# OS Handling of Threads

- One (or more) threads per running program
- The OS chooses which thread to run
  - To share a processor, the OS must be able to cleanly stop and start threads
- While one thread is using a processor, no other thread should interfere with its use
- To run a thread, OS must:
  - Load its code and data into memory
  - Set up HW control structures (e.g., the PC)
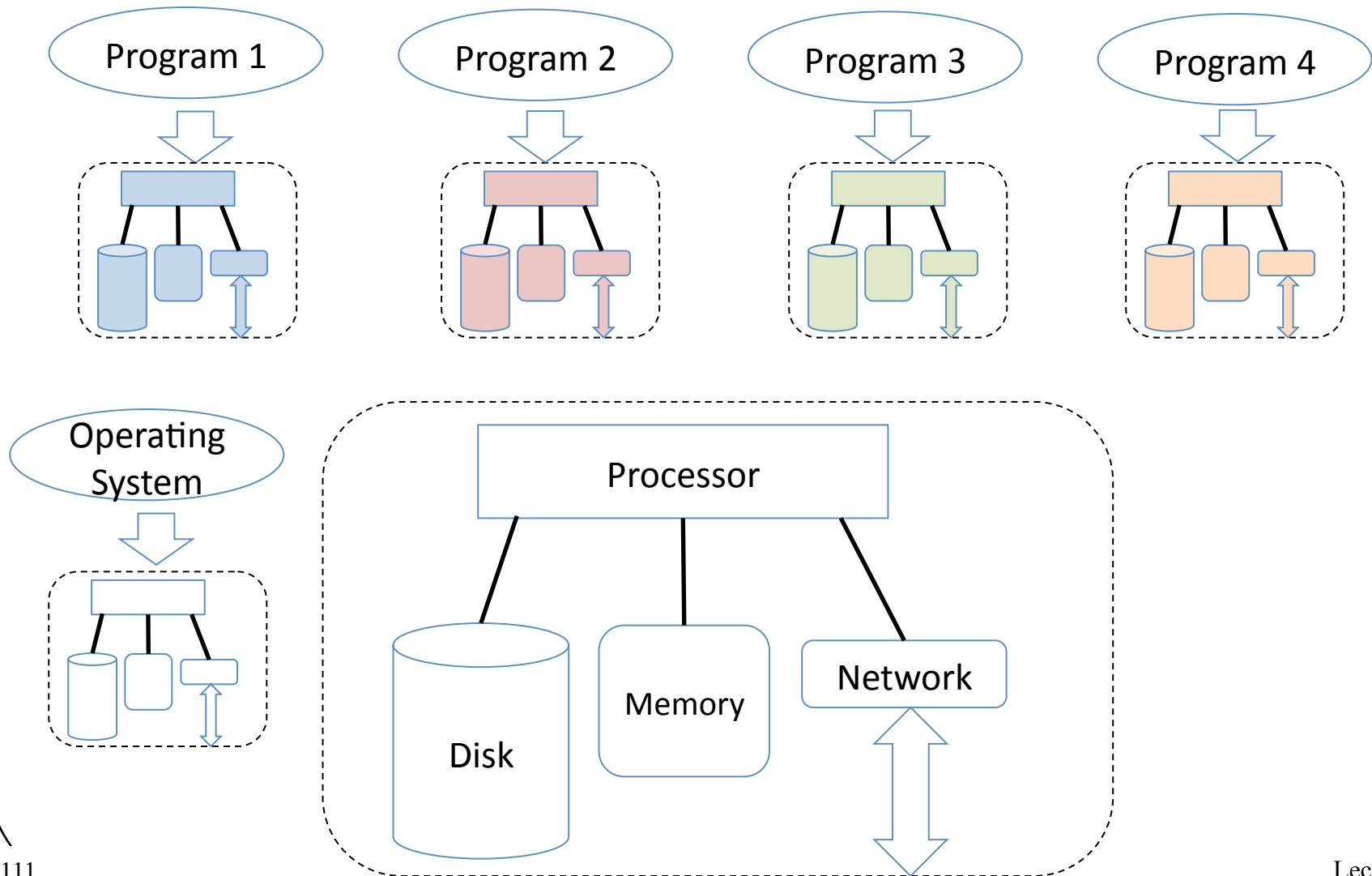  - Transfer control to the thread

# Time Slicing Virtualization

Program 1

Program 2

Program 3

Program 4

Processor

Disk

Memory

Network

# Wait a Minute . . .?

- How does the OS do all that?
- It's just a program itself
  - With its own interpreter, memory, etc.
- It must use the same physical resources as all the other threads
- Basically, the OS itself is a thread
- It creates and manages other threads
- Using privileged supervisor mode to safely and temporarily break virtualization boundaries

# The OS and Virtualization

Program 1

Program 2

Program 3

Program 4

Operating
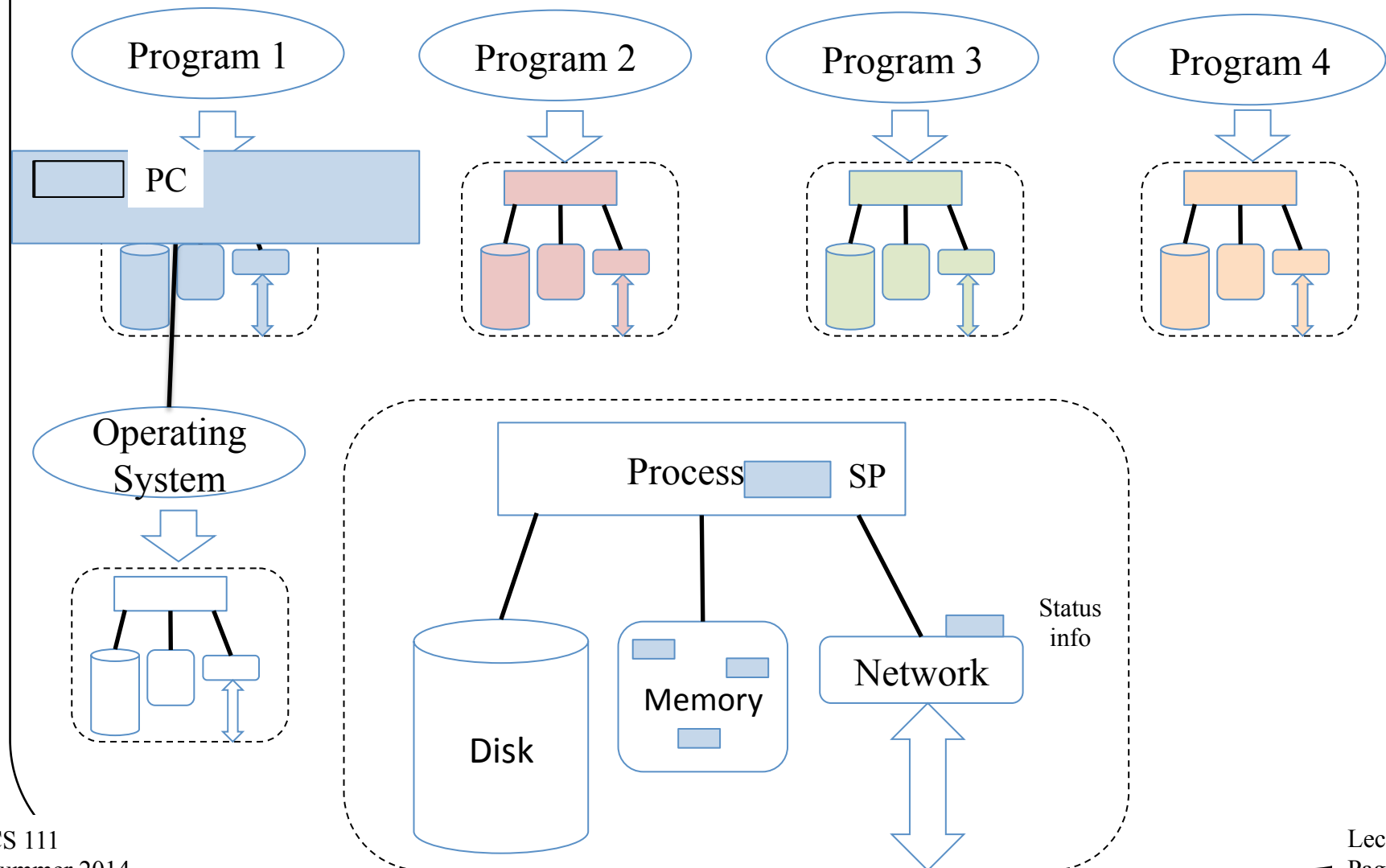System

Processor

Disk

Memory

Network

# Providing Contained Environments

- What must a thread manager control to keep each thread isolated from the others?

- Well, what can each thread do?
  - Run instructions
    - Make sure it can only run its own
  - Access some memory
    - Make sure it can only access its own
  - Communicate to other threads
    - Make sure communication uses a safe abstraction

# What Does This Boil Down To?

- Running threads have access to certain processor registers
  - Program counter, stack pointer, others
  - Thread manager must ensure those are all set correctly

- Running threads have access to some or all pieces of physical memory
  - Thread manager must ensure that a thread can only touch its own physical memory

- Running threads can request services (like communications)
  - Thread manager must provide safe access to those services

# Setting Up a User-Level VM

Program 1

Program 2

Program 3

Program 4

PC

Operating System

Process    SP

Status info

Disk

Memory

Network

# Protecting Threads

- Normal threads usually run in user mode

- Which means they can't touch certain things
  - In particular, each others' stuff

- For certain kinds of resources, that's a problem
  - What if two processes both legitimately need to write to the screen?
  - Do we allow unrestricted writing and hope for the best?
  - Don't allow them to write at all?

- Instead, trap to supervisor mode

# Trapping to Supervisor Mode

- To allow a program safe access to shared resources

- The trap goes to trusted code
  - Not under control of the program

- And performs well-defined actions
  - In ways that are safe

- E.g., program not allowed to write to the screen directly
  - But traps to OS code that writes it safely

# Modularity and Memory

- Clearly, programs must have access to memory
- We need abstractions that give them the required access
  - But with appropriate safety
- What we've really got (typically) is RAM
- RAM is pretty nice
  - But it has few built-in protections
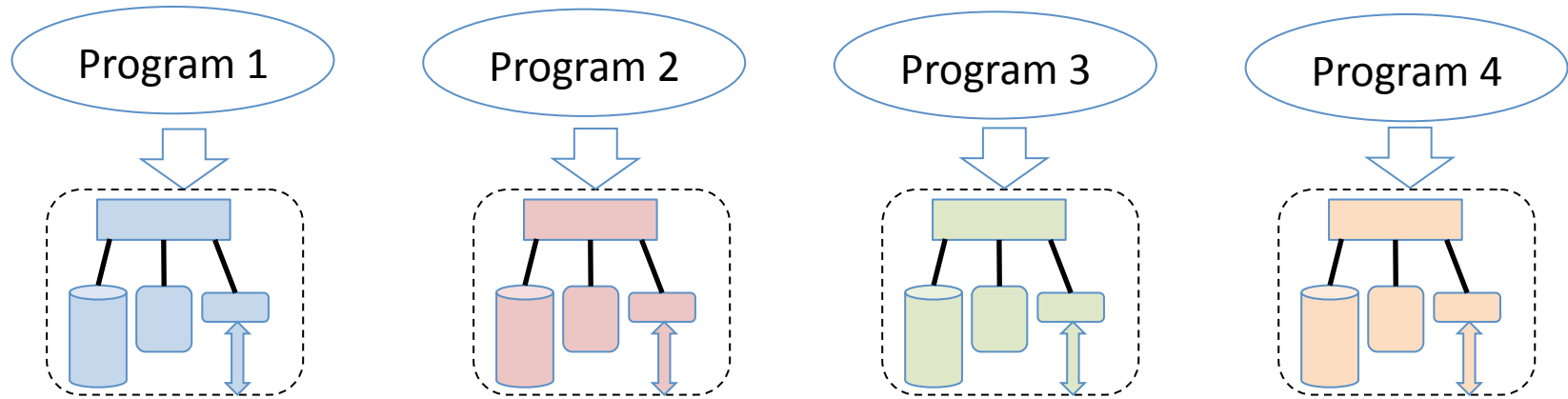- So we want an abstraction that provides RAM with safety

# What's the Safety Issue?

- We have multiple threads running
- Each requires some memory
- Modern architectures typically have one big pool of RAM
- How can we share the same pool of RAM among multiple processes?
  – Giving each what it needs
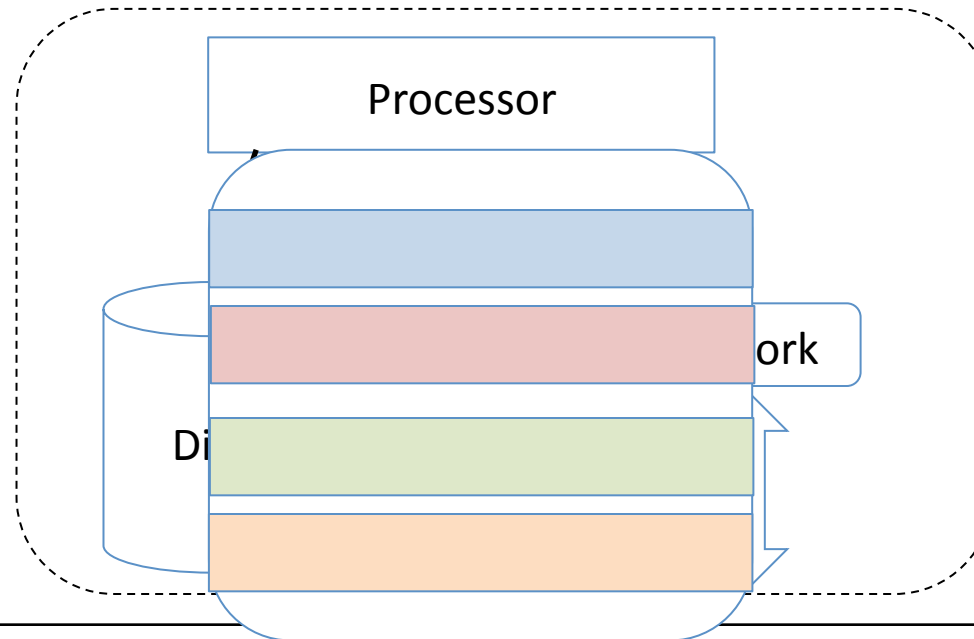  – Not allowing any to harm the others

# Domains

- A simple memory abstraction
- Give each process access to some range of the physical memory
  - Its *domain*
  - Different domain for each process
- Allow process to read/write/execute memory in its domain
- And not touch any memory outside its domain

# Mapping Domains

Program 1    Program 2    Program 3    Program 4

Every process
gets its own
piece of memory

Processor

ork

Di

No process can
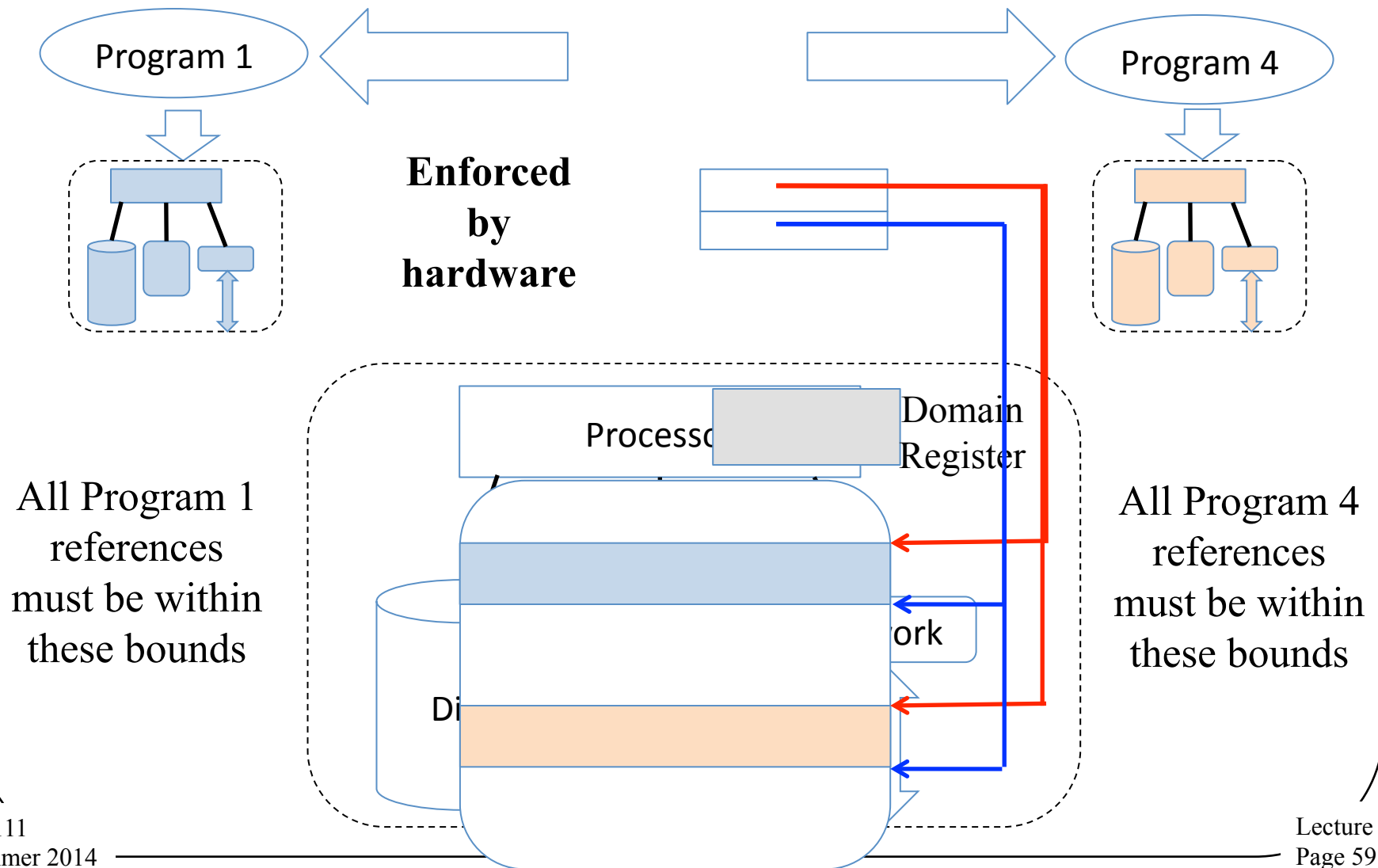interfere with
other processes'
memory

# What Do Domains Require?

- Threads will issue instructions
  - Perhaps using arbitrary memory addresses
- Only honor addresses in the thread's domain
  - Any other address should be caught as an error
- Hard modularity here requires HW support
- E.g., a domain register
  - Specifies the domain associated with the thread currently using the processor
  - By listing the low and high addresses that bound the domain

# The Memory Manager

- Hardware or software that enforces the bounds of the domain register

- When thread reads or writes an address, memory manager checks the domain register

- If within bounds, do the memory operation

- If not, throw an illegal memory reference exception

  – Trapping to supervisor mode

- Only trusted code (i.e., the OS) can change the domain register
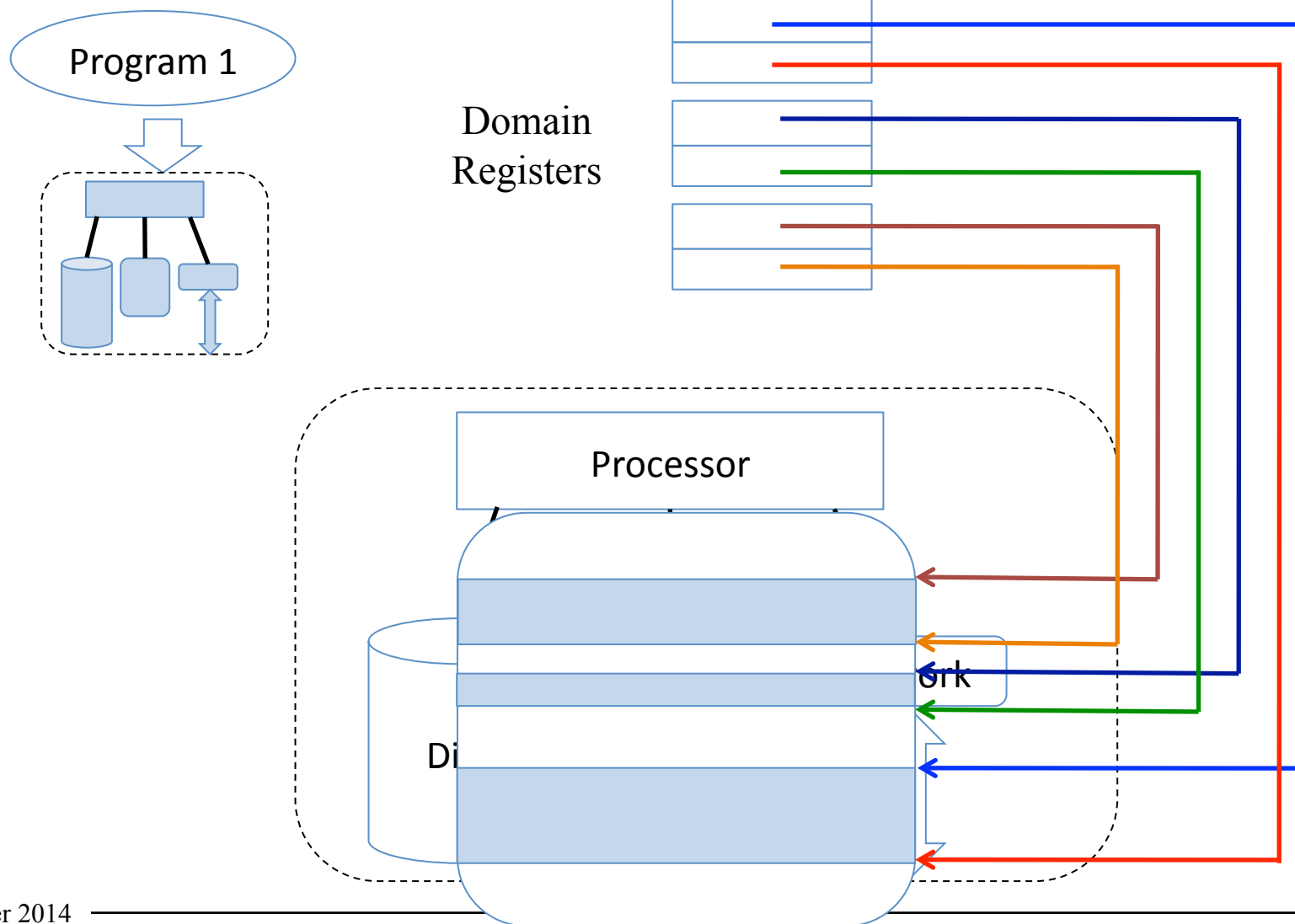
# The Domain Register Concept

Program 1

Program 4

**Enforced by hardware**

Processor

Domain Register

All Program 1 references must be within these bounds

All Program 4 references must be within these bounds

ork

Di

# Multiple Domains

- Limiting a process to a single domain is not too convenient

- The concept is easy to extend
  - Simply allow multiple domains per process

- Obvious way to handle this is with multiple domain registers
  - One per allocated domain

# The Multiple Domain Concept

Program 1

Domain
Registers

Processor

ork

Di

# Handling Multiple Domains

- Programs can request more domains
  - But the OS must set them up

- What does the program get to ask for?
  - A specific range of addresses?
  - Or a domain of a particular size?

- Latter is easier
  - What if requested set of addresses are already used by another program?
  - Memory manager can choose a range of addresses of requested size

# Domains and Access Permissions

- One can typically do three types of things with a memory address
  - Read its contents
  - Write a new value to it
  - Execute an instruction located there

- System can provide useful effects if it does not allow all modes of use to all addresses

- Typically handled on a per-domain basis
  - E.g., read-only domains

- Requires extra bits in domain registers

- And other hardware support

# What If Program Uses a Domain Improperly?

- E.g., it tries to write to a read-only domain

- A *permission error exception*
  - Different than an illegal memory reference exception

- But also handled by a similar mechanism

- Probably want it to be handled by somewhat different code in the OS

- Remember discussion of trap handling in previous lecture?

# Do We Really Need to Switch Processes for OS Services?

- When we trap or make a request for a domain, must we change processes?

  – We lose context doing so

- Instead, run the OS code for the process

  – Which requires changing to supervisor mode

  – Context for process is still available

- But what about safety?

  – Use domain access modes to ensure safety

- We don't do this for all OS services . . .

# Domains in Kernel Mode

- Allow user threads to access certain privileged domains
  - Like code to handle hardware traps
  - Code must be in a user-accessible domain

- But can't allow arbitrary access to those privileged domains

- A supervisor (AKA *kernel*) mode access bit is set on such domains
  - So thread only accesses them when in kernel mode

# How Does a Thread Get to Kernel Mode?

- Can't allow thread to arbitrarily put itself in kernel mode any time
    - Since it might do something unsafe

- Instead, allow entry to kernel mode only in specific ways
    - In particular, only at specific instructions
    - These are called *gates*
    - Typically implemented in hardware using instruction like SVC (supervisor call)