# CS162
## Operating Systems and Systems Programming
## Lecture 16

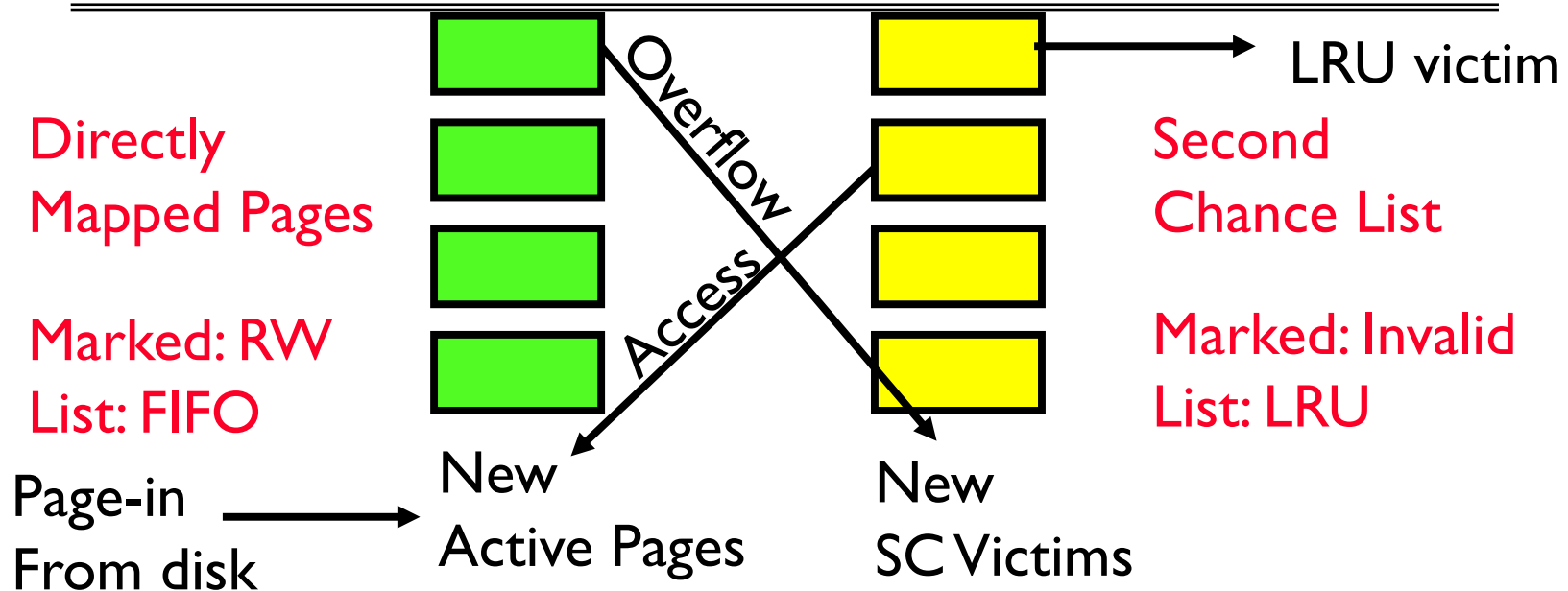## Finish caching,
## General I/O

October 20th, 2018

Ion Stoica

http://cs162.eecs.Berkeley.edu

# Second-Chance List Algorithm (VAX/VMS)



**Directly Mapped Pages**

Overflow

LRU victim

**Second Chance List**

**Marked: RW**
**List: FIFO**

Access

**Marked: Invalid**
**List: LRU**

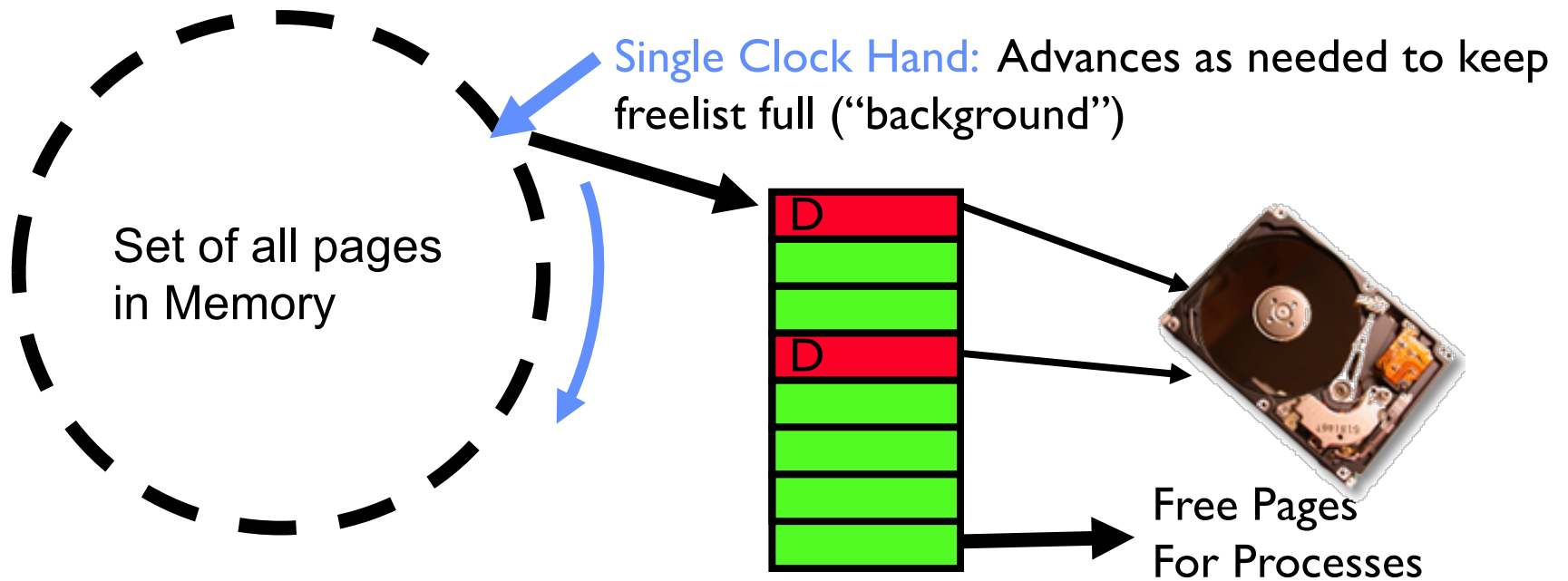Page-in From disk

New Active Pages

New SC Victims

- Split memory in two: Active list (RW), SC list (Invalid)
- Access pages in Active list at full speed
- Otherwise, Page Fault
  - Always move overflow page from end of Active list to front of Second-chance list (SC) and mark invalid
  - Desired Page On SC List: move to front of Active list, mark RW
  - Not on SC list: page in to front of Active list, mark RW; page out LRU victim at end of SC list

# Second-Chance List Algorithm (con't)

- How many pages for second chance list?
  - If 0 $\Rightarrow$ FIFO
  - If all $\Rightarrow$ LRU, but page fault on every page reference
- Pick intermediate value. Result is:
  - Pro: Few disk accesses (page only goes to disk if unused for a long time)
  - Con: Increased overhead trapping to OS (software / hardware tradeoff)
- With page translation, we can adapt to any kind of access the program makes
  - Later, we will show how to use page translation / protection to share memory between threads on separated machines
- Question: why didn't VAX include "use" bit?
  - Strecker (architect) asked OS people, they said they didn't need it, so didn't implement it
  - He later got blamed, but VAX did OK anyway

# Free List

Single Clock Hand: Advances as needed to keep freelist full ("background")

Set of all pages in Memory

D

D

Free Pages For Processes

- Keep set of free pages ready for use in demand paging
  - Freelist filled in background by Clock algorithm or other technique ("Pageout demon")
  - Dirty pages start copying back to disk when enter list
- Like VAX second-chance list
  - If page needed before reused, just return to active set
- Advantage: faster for page fault
  - Can always use page (or pages) immediately on fault

# Demand Paging (more details)

- Does software-loaded TLB need use bit?
  Two Options:
  - Hardware sets use bit in TLB; when TLB entry is replaced, software copies use bit back to page table
  - Software manages TLB entries as FIFO list; everything not in TLB is Second-Chance list, managed as strict LRU

- Core Map
  - Page tables map virtual page $\rightarrow$ physical page
  - Do we need a reverse mapping (i.e. physical page $\rightarrow$ virtual page)?
    » Yes. Clock algorithm runs through page frames. If sharing, then multiple virtual-pages per physical page
    » Can't push page out to disk without invalidating all PTEs

# Allocation of Page Frames (Memory Pages)

- How do we allocate memory among different processes?
  - Does every process get the same fraction of memory?  Different fractions?
  - Should we completely swap some processes out of memory?
- Each process needs *minimum* number of pages
  - Want to make sure that all processes that are loaded into memory can make forward progress
  - Example:  IBM 370 – 6 pages to handle SS MOVE instruction:
    - » instruction is 6 bytes, might span 2 pages
    - » 2 pages to handle *from*
    - » 2 pages to handle *to*

- Possible Replacement Scopes:

  - Global replacement – process selects replacement frame from set of all frames; one process can take a frame from another

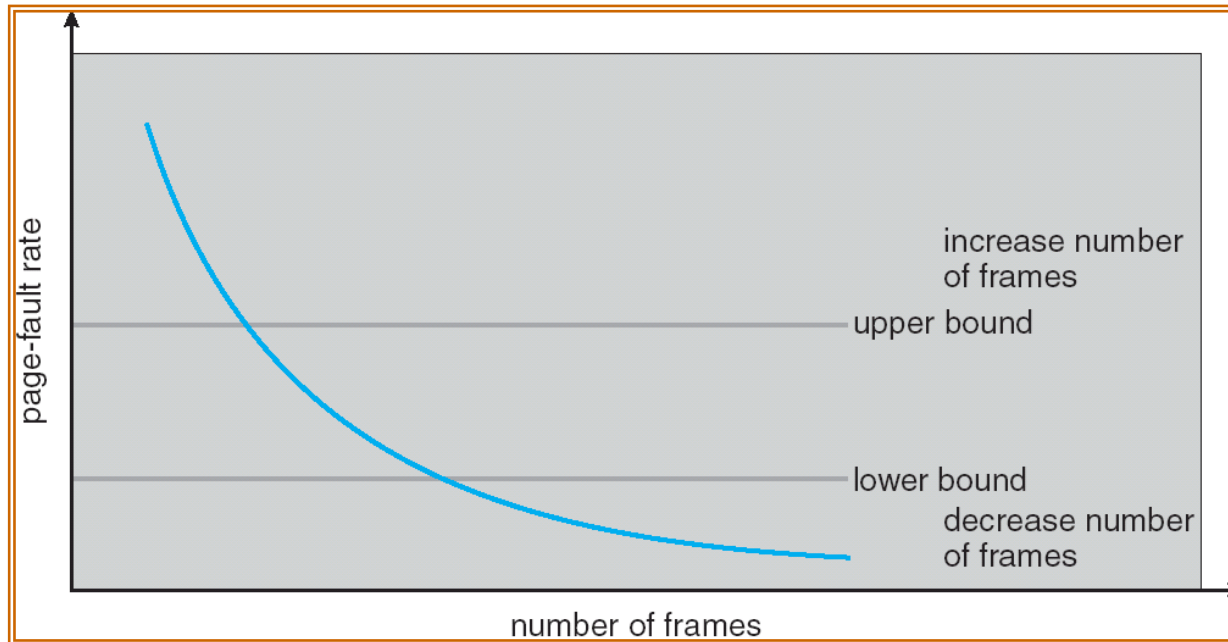  - Local replacement – each process selects from only its own set of allocated frames

# Fixed/Priority Allocation

- Equal allocation (Fixed Scheme):
  - Every process gets same amount of memory
  - Example: 100 frames, 5 processes → process gets 20 frames
- Proportional allocation (Fixed Scheme)
  - Allocate according to the size of process
  - Computation proceeds as follows:
    - $s_i$ = size of process $p_i$ and $S = \Sigma s_i$
    - $m$ = total number of frames

    $$a_i = \text{allocation for } p_i = \frac{s_i}{S} \times m$$

- Priority Allocation:
  - Proportional scheme using priorities rather than size
    - » Same type of computation as previous scheme
  - Possible behavior: If process $p_i$ generates a page fault, select for replacement a frame from a process with lower priority number
- Perhaps we should use an adaptive scheme instead???
  - What if some application just needs more memory?
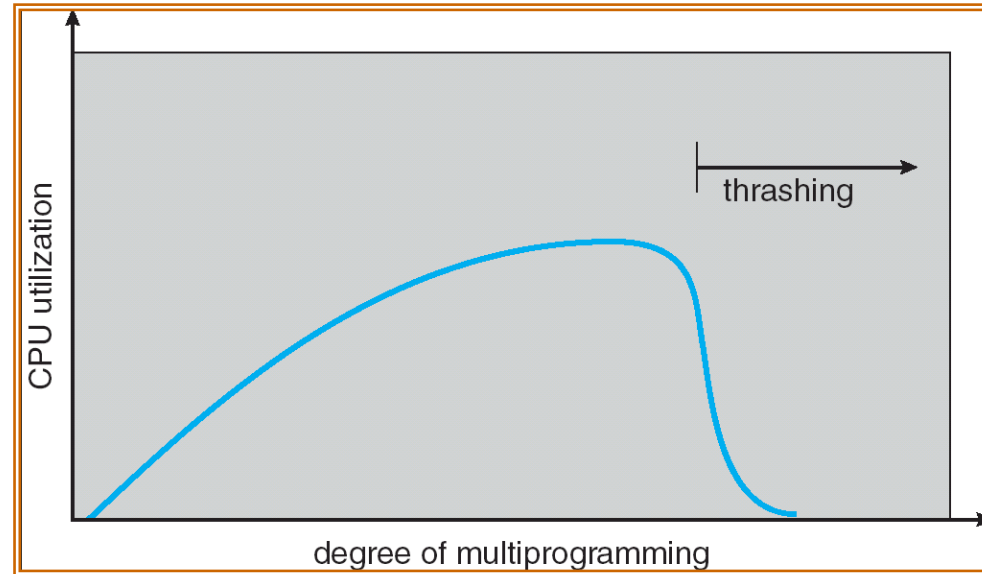
# Page-Fault Frequency Allocation

- Can we reduce Capacity misses by dynamically changing the number of pages/application?



- Establish "acceptable" page-fault rate
  - If actual rate too low, process loses frame
  - If actual rate too high, process gains frame
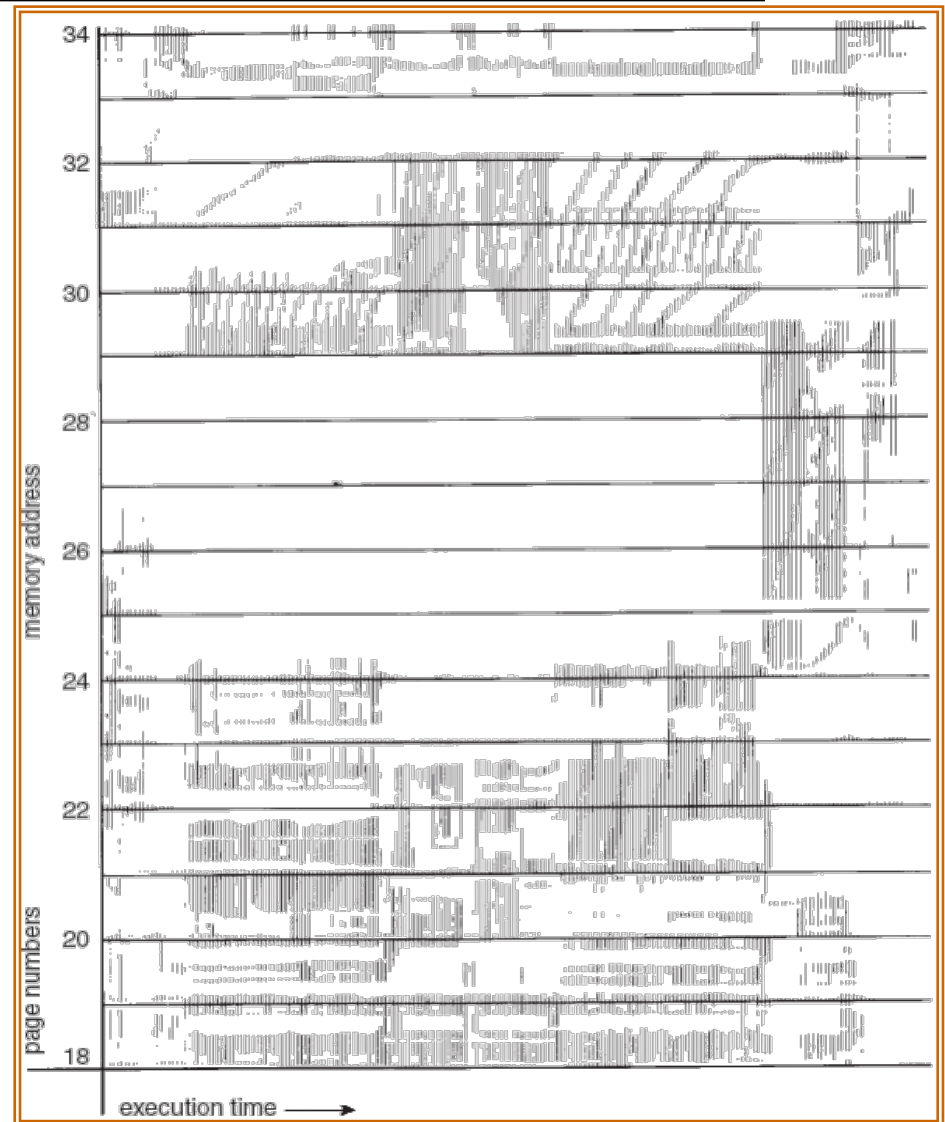- Question: What if we just don't have enough memory?

# Thrashing



A graph of CPU utilization (y-axis) versus degree of multiprogramming (x-axis). The curve rises, peaks, and then drops sharply in the region labeled "thrashing".
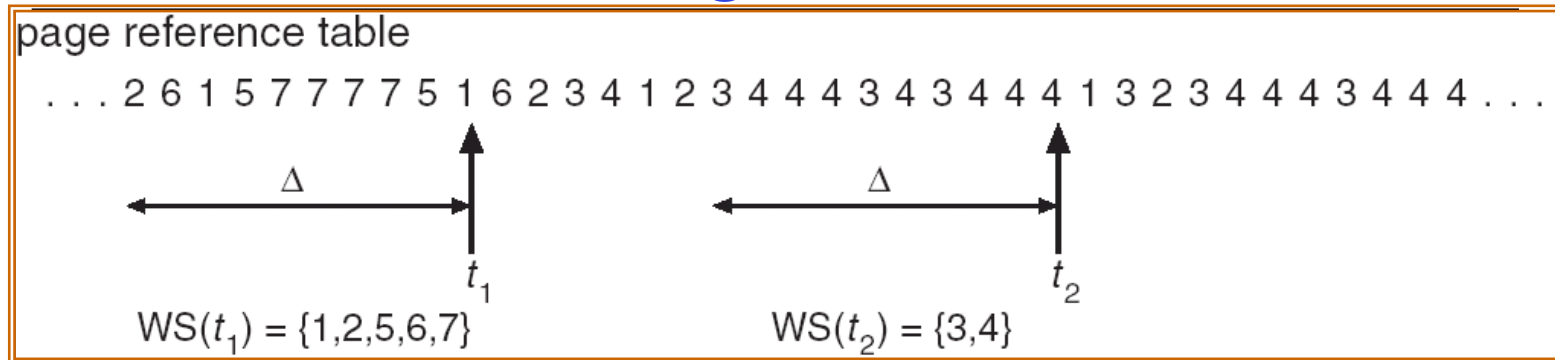
- If a process does not have "enough" pages, the page-fault rate is very high. This leads to:
  - low CPU utilization
  - operating system spends most of its time swapping to disk
- Thrashing ≡ a process is busy swapping pages in and out
- Questions:
  - How do we detect Thrashing?
  - What is best response to Thrashing?

# Locality In A Memory-Reference Pattern

- Program Memory Access Patterns have temporal and spatial locality
  - Group of Pages accessed along a given time slice called the "Working Set"
  - Working Set defines minimum number of pages needed for process to behave "well"
- Not enough memory for Working Set $\Rightarrow$ Thrashing
  - Better to swap out process?

# Working-Set Model

```
page reference table
    . . . 2 6 1 5 7 7 7 7 5 1 6 2 3 4 1 2 3 4 4 4 3 4 3 4 4 4 1 3 2 3 4 4 4 3 4 4 4 . . .
```

$WS(t_1) = \{1,2,5,6,7\}$        $WS(t_2) = \{3,4\}$

- $\Delta \equiv$ working-set window $\equiv$ fixed number of page references
  - Example: 10,000 instructions
- $WS_i$ (working set of Process $P_i$) = total set of pages referenced in the most recent $\Delta$ (varies in time)
  - if $\Delta$ too small will not encompass entire locality
  - if $\Delta$ too large will encompass several localities
  - if $\Delta = \infty \Rightarrow$ will encompass entire program
- $D = \Sigma|WS_i| \equiv$ total demand frames
- if $D > m \Rightarrow$ Thrashing
  - Policy: if $D > m$, then suspend/swap out processes
  - This can improve overall system behavior by a lot!

# What about Compulsory Misses?

- Recall that compulsory misses are misses that occur the first time that a page is seen
  - Pages that are touched for the first time
  - Pages that are touched after process is swapped out/swapped back in
- Clustering:
  - On a page-fault, bring in multiple pages "around" the faulting page
  - Since efficiency of disk reads increases with sequential reads, makes sense to read several sequential pages
- Working Set Tracking:
  - Use algorithm to try to track working set of application
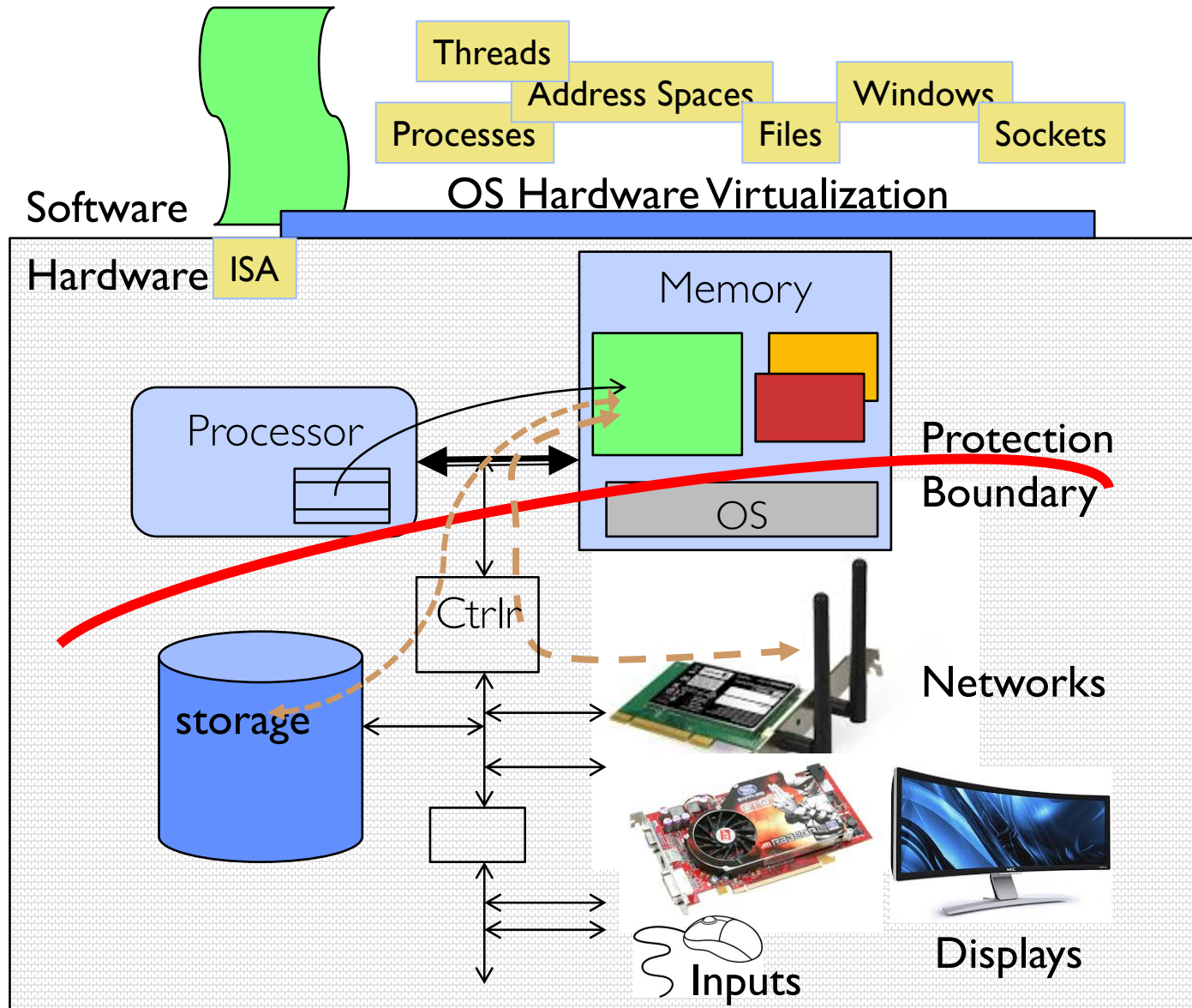  - When swapping process back in, swap in working set

# GENERAL I/O

# The Requirements of I/O

- So far in this course:
  - We have learned how to manage CPU and memory

- What about I/O?
  - Without I/O, computers are useless (disembodied brains?)
  - But… thousands of devices, each slightly different
    - » How can we standardize the interfaces to these devices?
  - Devices unreliable: media failures and transmission errors
    - » How can we make them reliable???
  - Devices unpredictable and/or slow
    - » How can we manage them if we don't know what they will do or how they will perform?

# OS Basics: I/O

Threads

Address Spaces

Windows

Processes

Files

Sockets

Software

OS Hardware Virtualization

Hardware

ISA

Memory

Processor

OS

Protection
Boundary

Ctrlr

storage

Networks

Displays

Inputs

# In a Picture



- I/O devices you recognize are supported by I/O Controllers
- Processors accesses them by reading and writing IO registers as if they were memory
  - Write commands and arguments, read status and results

# Operational Parameters for I/O

- Data granularity: Byte vs. Block
  - Some devices provide single byte at a time (e.g., keyboard)
  - Others provide whole blocks (e.g., disks, networks, etc.)

- Access pattern: Sequential vs. Random
  - Some devices must be accessed sequentially (e.g., tape)
  - Others can be accessed "randomly" (e.g., disk, cd, etc.)
    - » Fixed overhead to start transfers
  - Some devices require continual monitoring
  - Others generate interrupts when they need service

- Transfer Mechanism: Programmed IO and DMA

# Kernel Device Structure

| The System Call Interface | | | | |
|---|---|---|---|---|
| Process Management | Memory Management | Filesystems | Device Control | Networking |
| Concurrency, multitasking | Virtual memory | Files and dirs: the VFS | TTYs and device access | Connectivity |
| Architecture Dependent Code | Memory Manager | File System Types<br>⬛⬛⬛⬛<br><br>Block Devices<br>⬛⬛⬛⬛ | Device Control | Network Subsystem<br><br>IF drivers<br>⬛⬛⬛⬛ |

CS162 © UCB Fall 2018

# The Goal of the I/O Subsystem

- Provide Uniform Interfaces, Despite Wide Range of Different Devices

  - This code works on many different devices:

    ```
    FILE fd = fopen("/dev/something", "rw");
    for (int i = 0; i < 10; i++) {
       fprintf(fd, "Count %d\n", i);
    }
    close(fd);
    ```

  - Why? Because code that controls devices ("device driver") implements standard interface

- We will try to get a flavor for what is involved in actually controlling devices in rest of lecture

  - Can only scratch surface!

# Want Standard Interfaces to Devices

- Block Devices: e.g. disk drives, tape drives, DVD-ROM
  - Access blocks of data
  - Commands include `open()`, `read()`, `write()`, `seek()`
  - Raw I/O or file-system access
  - Memory-mapped file access possible
- Character Devices: e.g. keyboards, mice, serial ports, some USB devices
  - Single characters at a time
  - Commands include `get()`, `put()`
  - Libraries layered on top allow line editing
- Network Devices: e.g. Ethernet, Wireless, Bluetooth
  - Different enough from block/character to have own interface
  - Unix and Windows include socket interface
    » Separates network protocol from network operation
    » Includes `select()` functionality
  - Usage: pipes, FIFOs, streams, queues, mailboxes

# How Does User Deal with Timing?

- **Blocking Interface: "Wait"**
  - When request data (e.g. `read()` system call), put process to sleep until data is ready
  - When write data (e.g. `write()` system call), put process to sleep until device is ready for data

- **Non-blocking Interface: "Don't Wait"**
  - Returns quickly from read or write request with count of bytes successfully transferred
  - Read may return nothing, write may write nothing

- **Asynchronous Interface: "Tell Me Later"**
  - When request data, take pointer to user's buffer, return immediately; later kernel fills buffer and notifies user
  - When send data, take pointer to user's buffer, return immediately; later kernel takes data and notifies user

# Administrivia

- Midterm 2 coming up on <span style="color:red">Mon 10/29 5:00-6:30PM</span>
  - All topics up to and including Lecture 17
    - » Focus will be on Lectures 11 – 17 and associated readings
    - » Projects 1 and 2
    - » Homework 0 – 2
  - Closed book
  - 2 pages hand-written notes both sides

- Please fill in the survey:
  https://www.surveymonkey.com/r/7YNS59R

# BREAK

# Chip-scale Features of 2015 x86 (Sky Lake)

- Significant pieces:
  - Four OOO cores with deeper buffers
    - » New Intel MPX (Memory Protection Extensions)
    - » New Intel SGX (Software Guard Extensions)
    - » Issue up to 6 μ-ops/cycle
  - Integrated GPU, System Agent (Mem, Fast I/O)
  - Large shared L3 cache with on-chip ring bus
    - » 2 MB/core instead of 1.5 MB/core
    - » High-BW access to L3 Cache
- Integrated I/O
  - Integrated memory controller (IMC)
    - » Two independent channels of DDR3L/DDR4 DRAM
  - High-speed PCI-Express (for Graphics cards)
  - Direct Media Interface (DMI) Connection to PCH (Platform Control Hub)
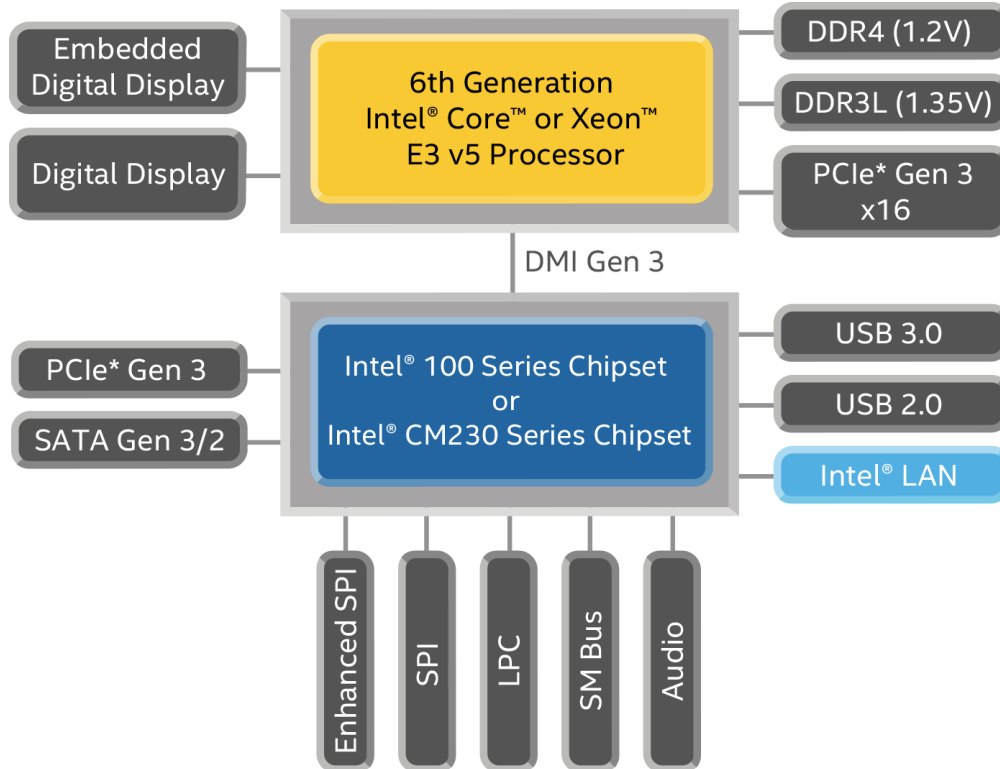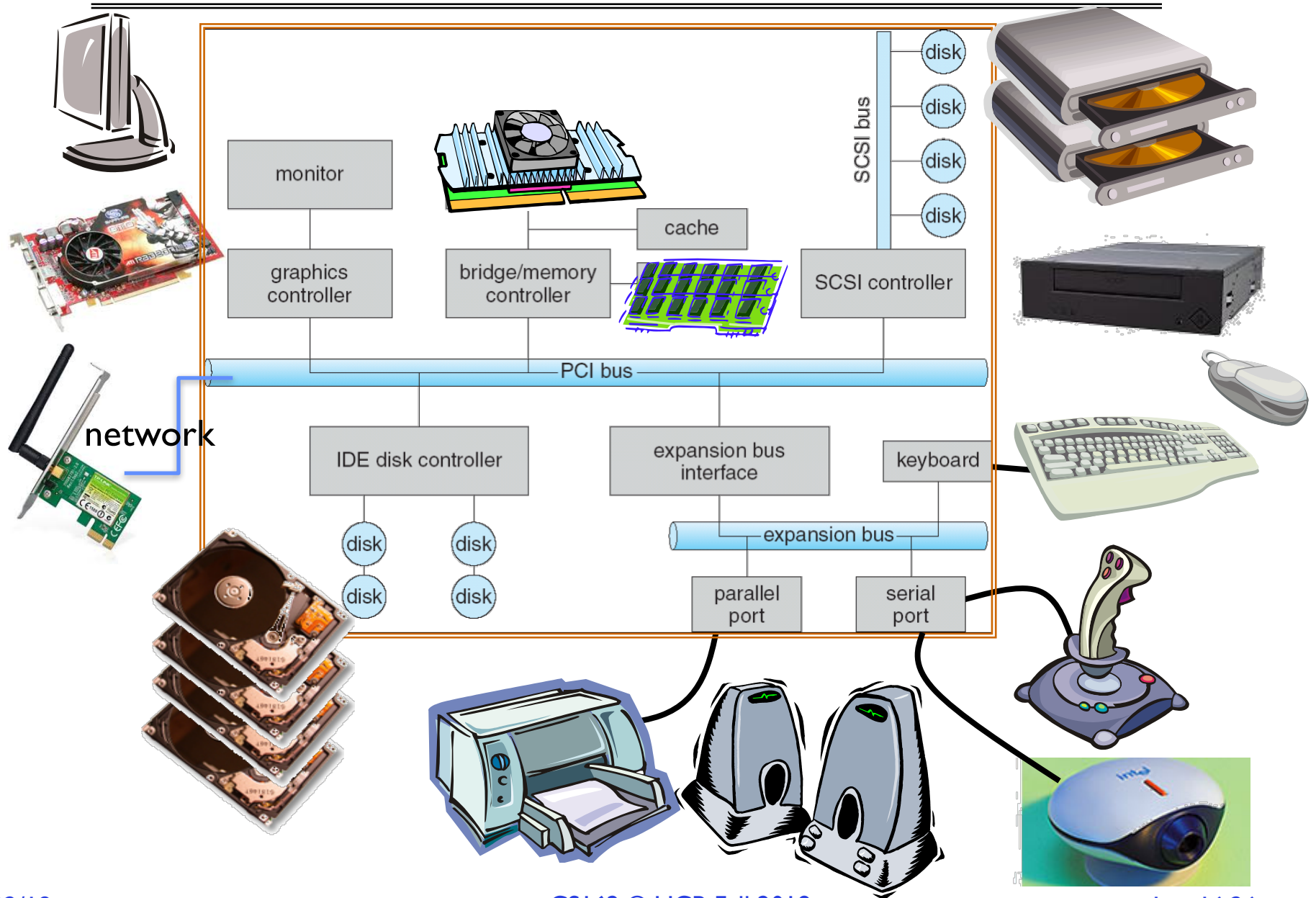
# Sky Lake I/O: PCH



**Sky Lake System Configuration**

- Platform Controller Hub
  - Connected to processor with proprietary bus
    » Direct Media Interface
- Types of I/O on PCH:
  - USB, Ethernet
  - Thunderbolt 3
  - Audio, BIOS support
  - More PCI Express (lower speed than on Processor)
  - SATA (for Disks)

# Modern I/O Systems

monitor

cache

bridge/memory
controller

SCSI controller

graphics
controller

SCSI bus

disk

disk

disk

disk

network

PCI bus

IDE disk controller

expansion bus
interface

keyboard

disk

disk

disk

disk

expansion bus

parallel
port

serial
port

# Example: PCI Architecture

# Example Device-Transfer Rates in Mb/s (Sun Enterprise 6000)



- Device Rates vary over 12 orders of magnitude !!!
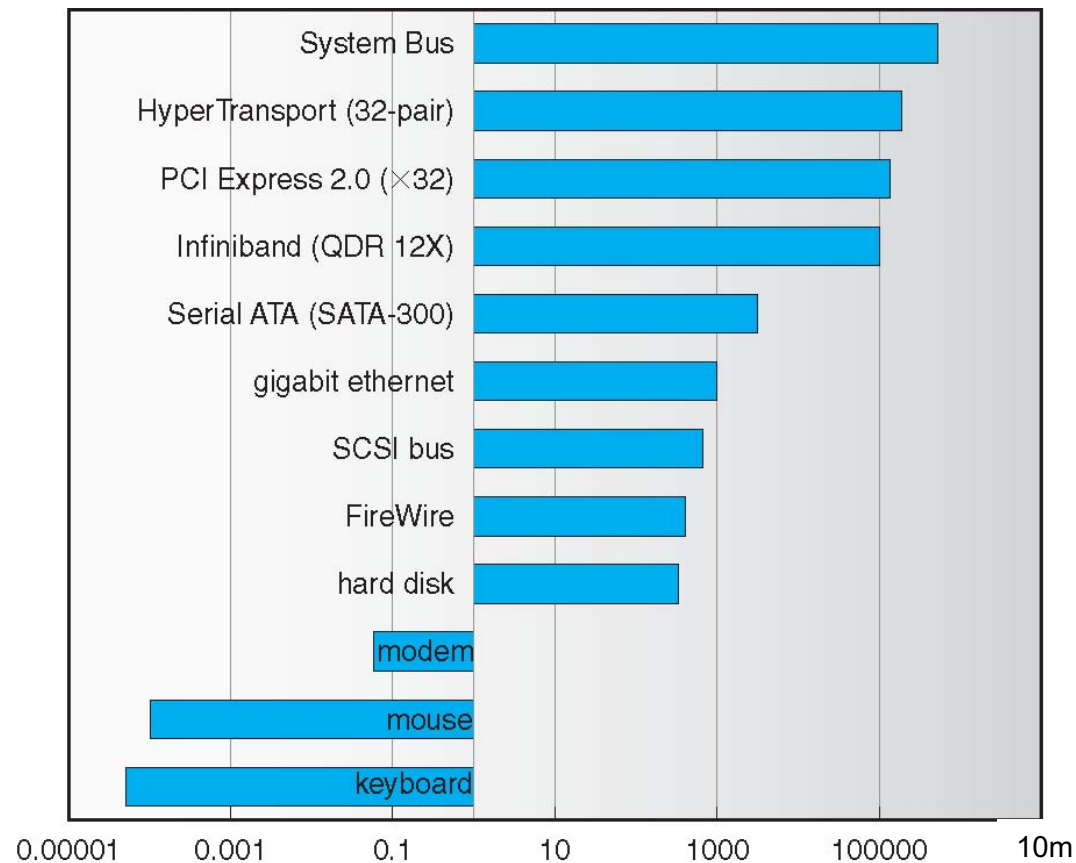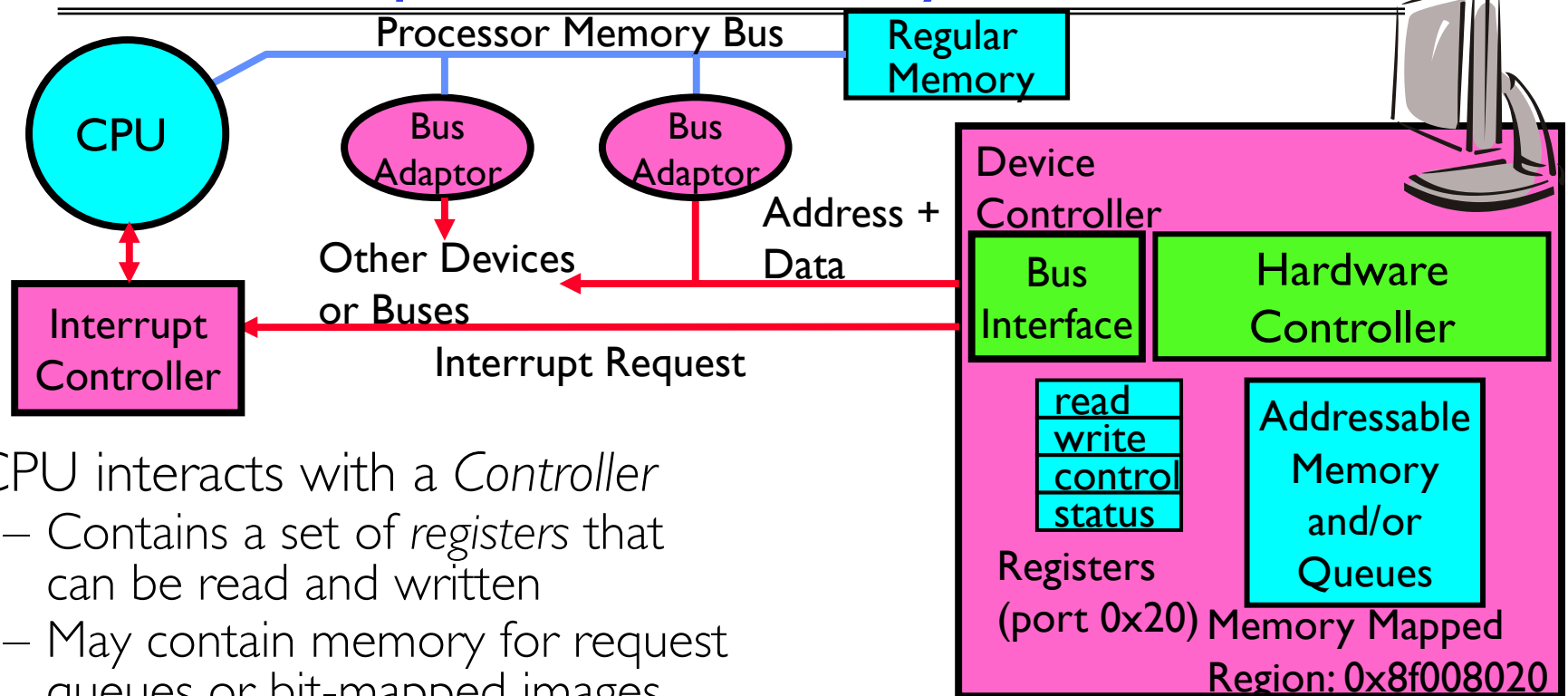  - System better be able to handle this wide range
  - Better not have high overhead/byte for fast devices!
  - Better not waste time waiting for slow devices

# How does the processor actually talk to the device?

Processor Memory Bus

**Regular Memory**

**CPU**

**Bus Adaptor**

**Bus Adaptor**

Other Devices or Buses

Address + Data

Interrupt Request

**Interrupt Controller**

**Device Controller**

**Bus Interface**

**Hardware Controller**

read
write
control
status

**Addressable Memory and/or Queues**

Registers (port 0x20)

Memory Mapped Region: 0x8f008020

- CPU interacts with a *Controller*
  - Contains a set of *registers* that can be read and written
  - May contain memory for request queues or bit-mapped images
- Regardless of the complexity of the connections and buses, processor accesses registers in two ways:
  - I/O instructions: in/out instructions
    - » Example from the Intel architecture: `out 0x21,AL`
  - Memory mapped I/O: load/store instructions
    - » Registers/memory appear in physical address space
    - » I/O accomplished with load and store instructions

# Example: Memory-Mapped Display Controller

- Memory-Mapped:
  - Hardware maps control registers and display memory into physical address space
    - » Addresses set by HW jumpers or at boot time
  - Simply writing to display memory (also called the "frame buffer") changes image on screen
    - » Addr: `0x8000F000 – 0x8000FFFF`
  - Writing graphics description to cmd queue
    - » Say enter a set of triangles describing some scene
    - » Addr: `0x80010000 – 0x8001FFFF`
  - Writing to the command register may cause on-board graphics hardware to do something
    - » Say render the above scene
    - » Addr: `0x0007F004`
- Can protect with address translation

| Address | |
|---|---|
| 0x80020000 | **Graphics Command Queue** |
| 0x80010000 | **Display Memory** |
| 0x8000F000 | |
| 0x0007F004 | **Command** |
| 0x0007F000 | **Status** |

**Physical Address Space**

# Transferring Data To/From Controller

- **Programmed I/O:**
  - Each byte transferred via processor in/out or load/store
  - Pro: Simple hardware, easy to program
  - Con: Consumes processor cycles proportional to data size

- **Direct Memory Access:**
  - Give controller access to memory bus
  - Ask it to transfer data blocks to/from memory directly

- Sample interaction with DMA controller (from OSC book):



1. device driver is told to transfer disk data to buffer at address X
2. device driver tells disk controller to transfer C bytes from disk to buffer at address X
3. disk controller initiates DMA transfer
4. disk controller sends each byte to DMA controller
5. DMA controller transfers bytes to buffer X, increasing memory address and decreasing C until C = 0
6. when C = 0, DMA interrupts CPU to signal transfer completion

CPU
cache
DMA/bus/ interrupt controller
memory bus
memory    x buffer
PCI bus
IDE disk controller
disk  disk
disk  disk

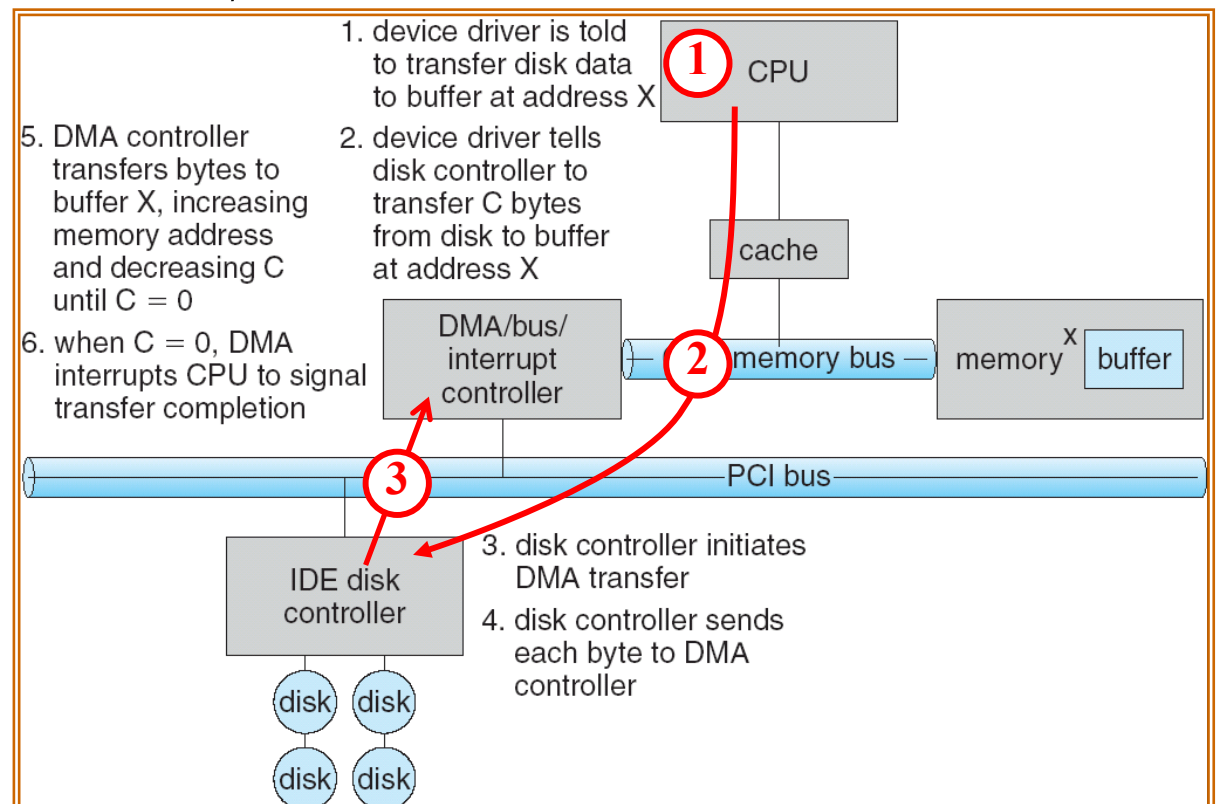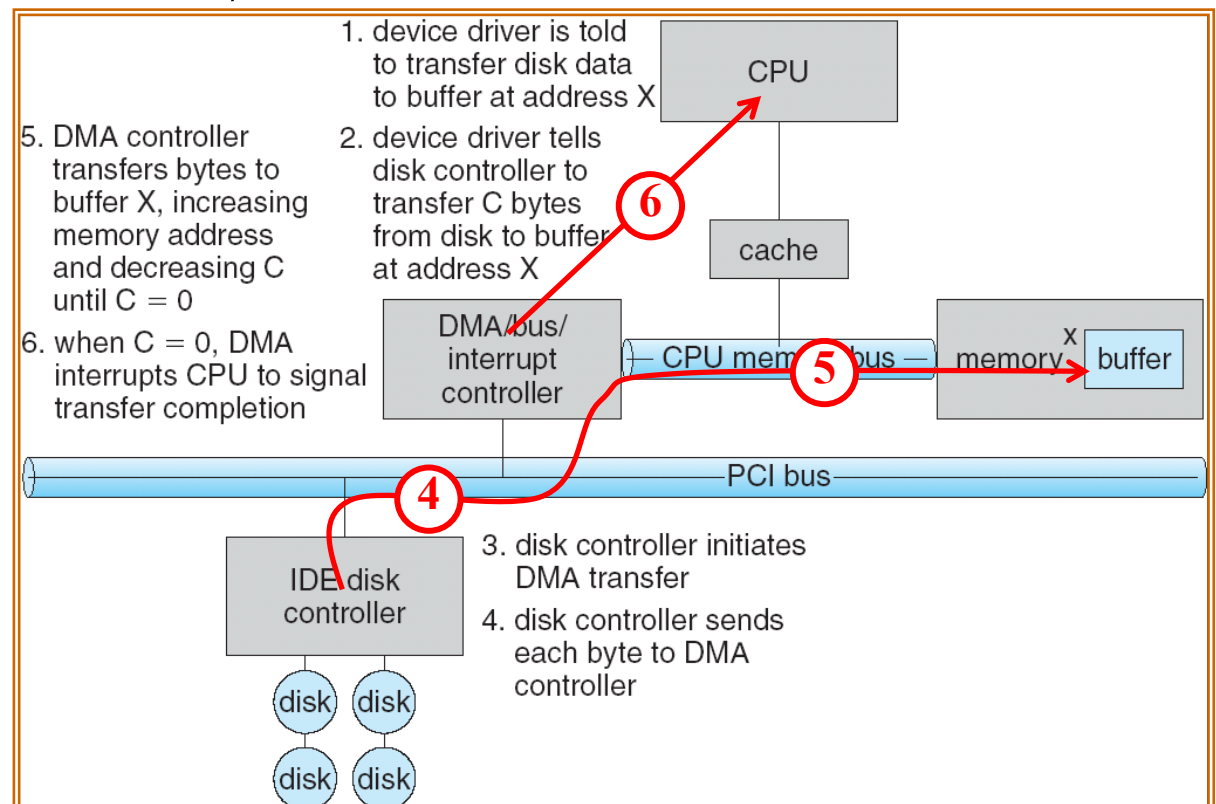# Transferring Data To/From Controller

- **Programmed I/O:**
  - Each byte transferred via processor in/out or load/store
  - Pro: Simple hardware, easy to program
  - Con: Consumes processor cycles proportional to data size

- **Direct Memory Access:**
  - Give controller access to memory bus
  - Ask it to transfer data blocks to/from memory directly

- Sample interaction with DMA controller (from OSC book):

1. device driver is told to transfer disk data to buffer at address X

2. device driver tells disk controller to transfer C bytes from disk to buffer at address X

3. disk controller initiates DMA transfer

4. disk controller sends each byte to DMA controller

5. DMA controller transfers bytes to buffer X, increasing memory address and decreasing C until C = 0

6. when C = 0, DMA interrupts CPU to signal transfer completion

CPU

cache

DMA/bus/interrupt controller

CPU memory bus

memory — X buffer

PCI bus

IDE disk controller
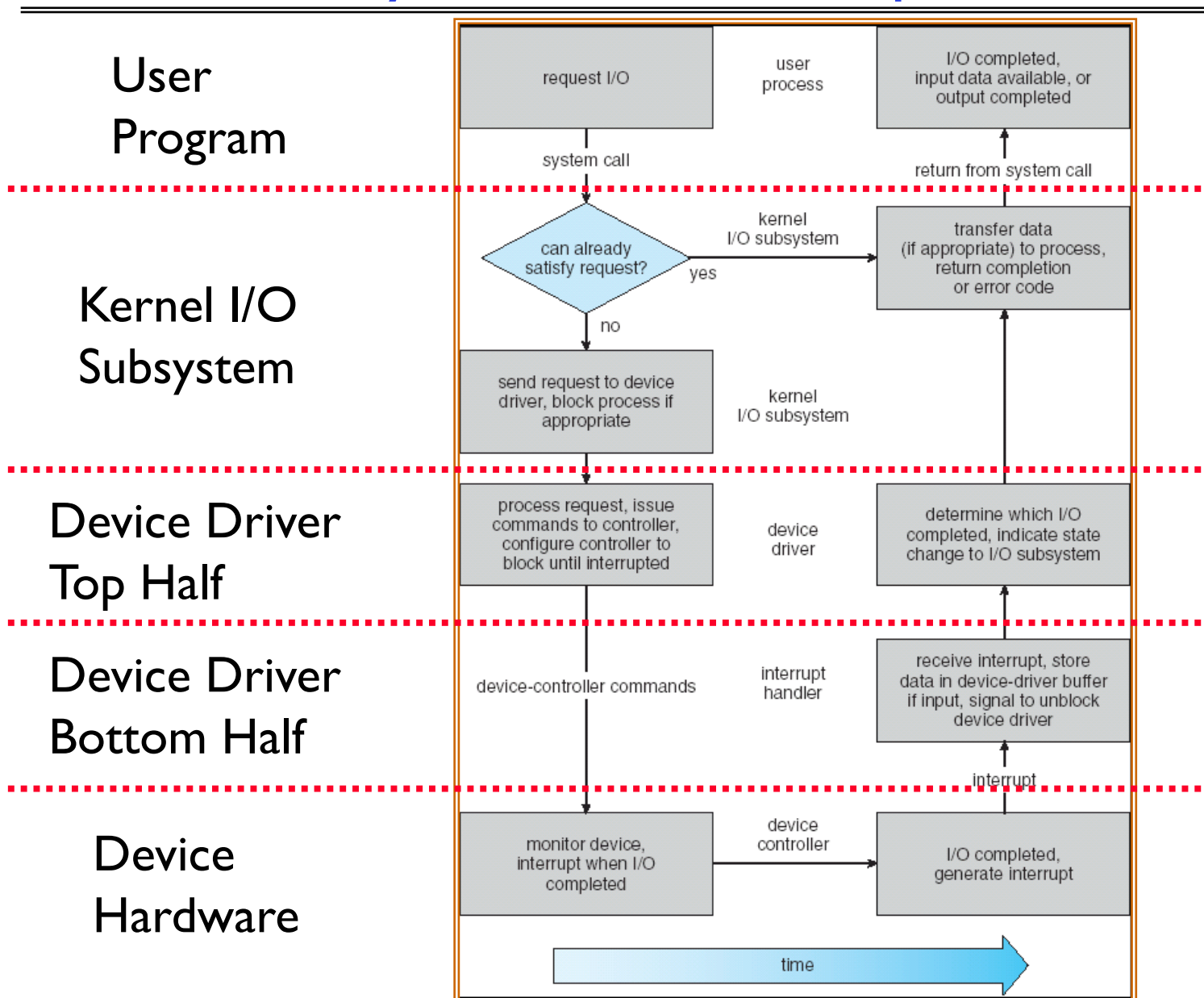
disk   disk
disk   disk

# I/O Device Notifying the OS

- The OS needs to know when:
  - The I/O device has completed an operation
  - The I/O operation has encountered an error
- I/O Interrupt:
  - Device generates an interrupt whenever it needs service
  - Pro: handles unpredictable events well
  - Con: interrupts relatively high overhead
- Polling:
  - OS periodically checks a device-specific status register
    - » I/O device puts completion information in status register
  - Pro: low overhead
  - Con: may waste many cycles on polling if infrequent or unpredictable I/O operations
- Actual devices combine both polling and interrupts
  - For instance – High-bandwidth network adapter:
    - » Interrupt for first incoming packet
    - » Poll for following packets until hardware queues are empty

# Device Drivers

- Device Driver: Device-specific code in the kernel that interacts directly with the device hardware
  - Supports a standard, internal interface
  - Same kernel I/O system can interact easily with different device drivers
  - Special device-specific configuration supported with the `ioctl()` system call

- Device Drivers typically divided into two pieces:
  - Top half: accessed in call path from system calls
    » implements a set of standard, cross-device calls like `open()`, `close()`, `read()`, `write()`, `ioctl()`, `strategy()`
    » This is the kernel's interface to the device driver
    » Top half will *start* I/O to device, may put thread to sleep until finished
  - Bottom half: run as interrupt routine
    » Gets input or transfers next block of output
    » May wake sleeping threads if I/O now complete

# Life Cycle of An I/O Request

**User Program**

**Kernel I/O Subsystem**

**Device Driver Top Half**

**Device Driver Bottom Half**

**Device Hardware**

# Summary

- I/O Devices Types:
  - Many different speeds (0.1 bytes/sec to GBytes/sec)
  - Different Access Patterns:
    - » Block Devices, Character Devices, Network Devices
  - Different Access Timing:
    - » Blocking, Non-blocking, Asynchronous
- I/O Controllers: Hardware that controls actual device
  - Processor Accesses through I/O instructions, load/store to special physical memory
- Notification mechanisms
  - Interrupts
  - Polling: Report results through status register that processor looks at periodically
- Device drivers interface to I/O devices
  - Provide clean Read/Write interface to OS above
  - Manipulate devices through PIO, DMA & interrupt handling
  - Three types: block, character, and network