

CS162

Operating Systems and Systems Programming

Lecture 16

Introduction to I/O subsystems

October 24, 2019

Prof. David Culler

<http://cs162.eecs.Berkeley.edu>

Read: A&D Ch 11.2-4,12

So what are we supposed to understand about Virtual Machines ???

Today: turn system / user inside out



- Audio in capture failed at 6 mins
- Will re-record Lec 15

10/22/19

UCB CS162 Fa19 L15

14

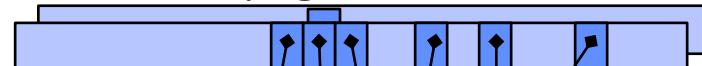
Virtual Machine Concepts

- Despite its seeming impossibility (running a full Guest OS and its Appl Processes at User privilege) it is all built on concepts you understand of classic operating systems
 - Processes, Threads, Page Tables, Interrupts, Syscall, RPC, Drivers
- Direct execution of the Guest OS Application Process is made possible by VMM constructing Shadow Page Table that is the composition of two translations

Recall: VMM Shadow Page Tables

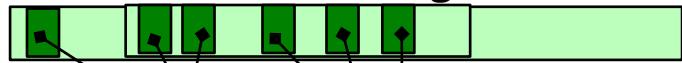
Guest OS Page Tables

Guest VAS pages



Guest Physical Frames

Host VM Process Pages

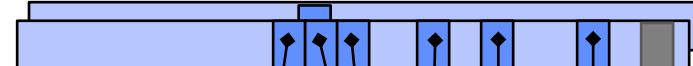


Host Physical Frames

Host OS Page Table for VM Process

VMM Page Tables for VM

Shadow VAS pages

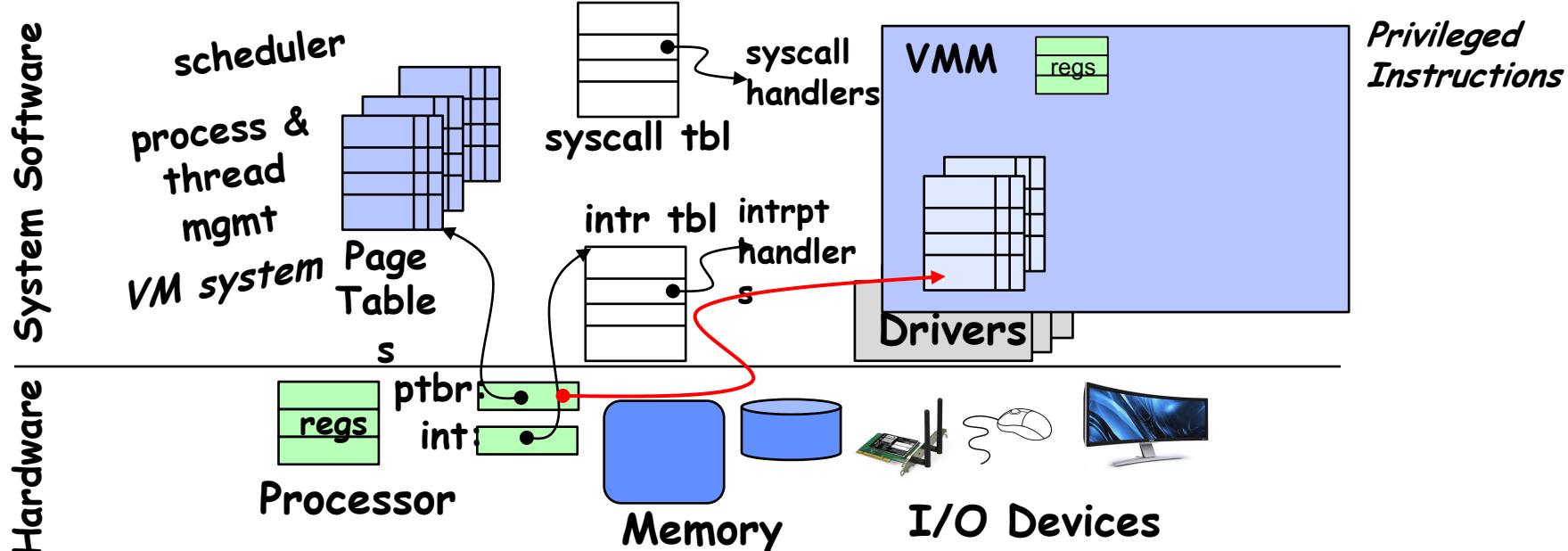
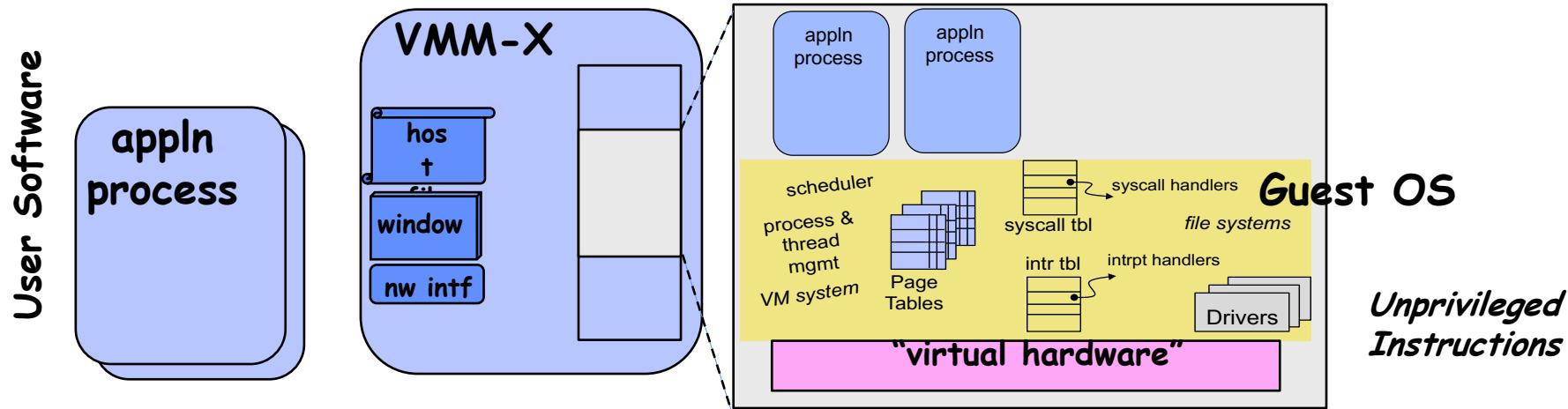


Host Physical Frames

Virtual Machine Concepts

- Despite its seeming impossibility (running a full Guest OS and its Appl Processes at User privilege) it is all built on concepts you understand of classic operating systems
 - Processes, Threads, Page Tables, Interrupts, Syscall, RPC, Drivers
- Direct execution of the Guest OS Application Process is made possible by VMM constructing Shadow Page Table that is the composition of two translations
- In order for the hardware to use a Shad-PT to translate Guest-Proc VAS, the VMM needs to switch out the Host OS PT

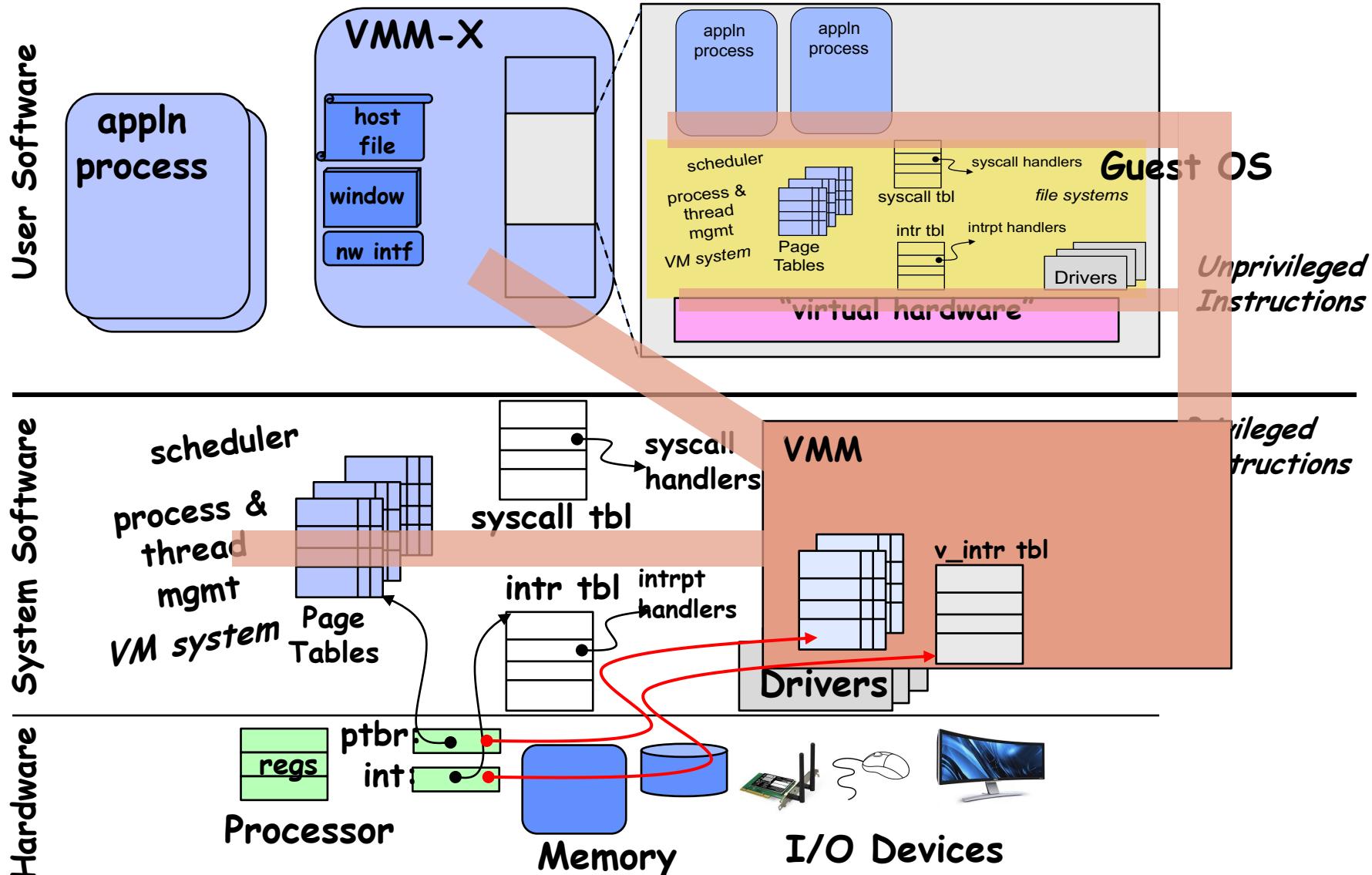
Recall: VMM – VM – VMM extension process



Virtual Machine Concepts

- Despite its seeming impossibility (running a full Guest OS and its Appl Processes at User privilege) it is all built on concepts you understand of classic operating systems
 - Processes, Threads, Page Tables, Interrupts, Syscall, RPC, Drivers
- Direct execution of the Guest OS Application Process is made possible by VMM constructing Shadow Page Table that is the composition of two translations
- In order for the hardware to use a Shad-PT to translate Guest-Proc VAS, the VMM needs to switch out the Host OS PT
- VMM interposes on all Guest Proc \leftrightarrow Guest OS interactions
 - Every Guest U \rightarrow K and K \rightarrow U goes through the VMM

Recall: VMM – VM – VMM extension process



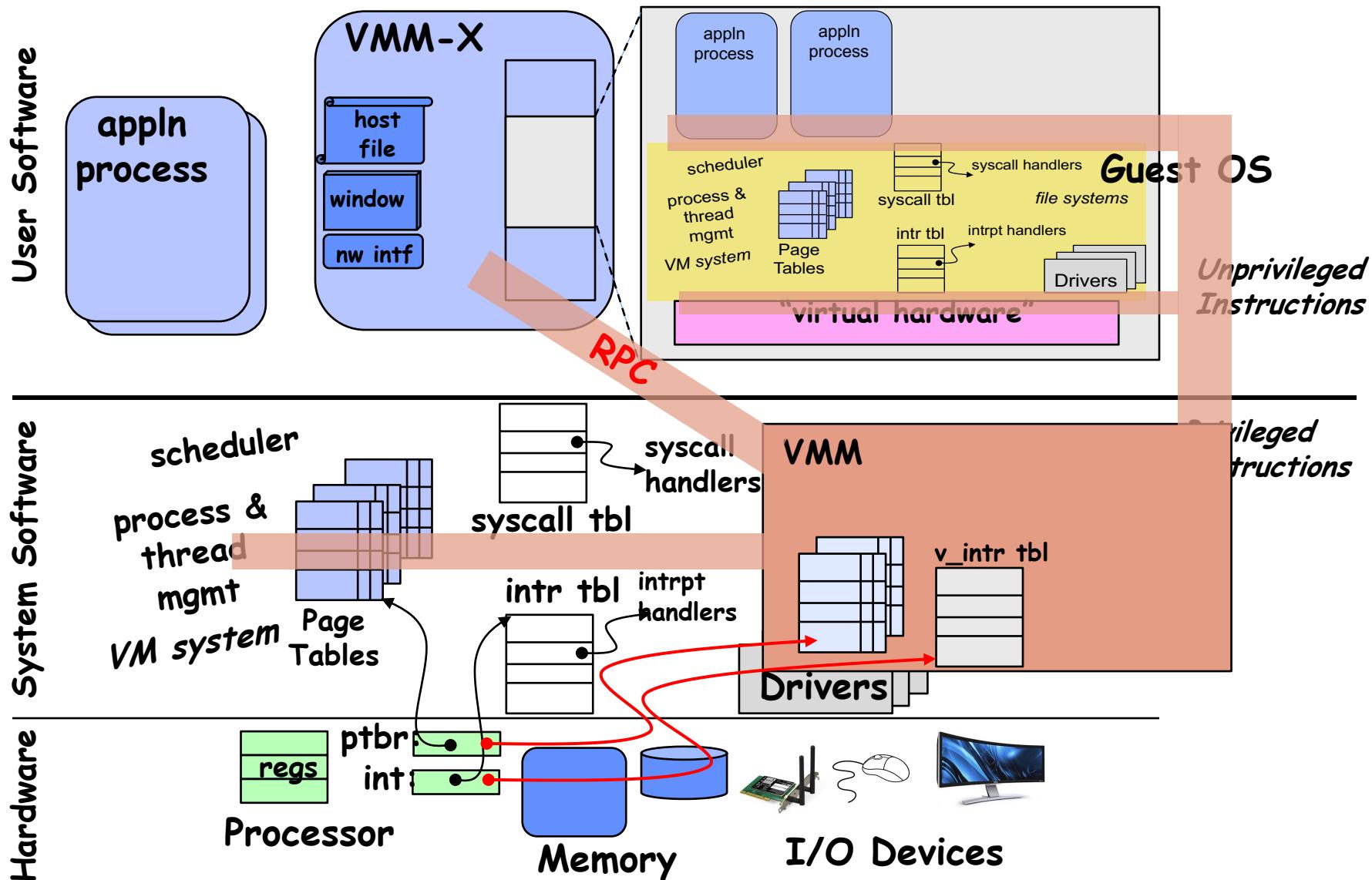
Virtual Machine Concepts

- Despite its seeming impossibility (running a full Guest OS and its Appl Processes at User privilege) it is all built on concepts you understand of classic operating systems
 - Processes, Threads, Page Tables, Interrupts, Syscall, RPC, Drivers
- Direct execution of the Guest OS Application Process is made possible by VMM constructing *Shadow Page Table* that is the composition of two translations
 - In order for the hardware to use a Shad-PT to translate Guest-Proc VAS, the VMM needs to switch out the Host OS PT
- VMM *interposes* on all Guest Proc $\leftarrow\rightarrow$ Guest OS interactions
 - Every Guest U \rightarrow K and K \rightarrow U goes through the VMM
- Trap and emulate: everything the Guest OS does it does by reading/writing its kernel address region. These trap and VMM emulates
 - Especially device drivers

Virtual Machine Concepts

- Despite its seeming impossibility (running a full Guest OS and its Application Processes at User privilege) it is all built on concepts you understand of classic operating systems
 - Processes, Threads, Page Tables, Interrupts, Syscall, RPC, Drivers
- Direct execution of the Guest OS Application Process is made possible by VMM constructing *Shadow Page Table* that is the composition of two translations
 - In order for the hardware to use a Shad-PT to translate Guest-Proc VAS, the VMM needs to switch out the Host OS PT
- VMM *interposes* on all Guest Proc $\leftarrow\rightarrow$ Guest OS interactions
 - Every Guest U \rightarrow K and K \rightarrow U goes through the VMM
- Trap and emulate: everything the Guest OS does it does by reading/writing its kernel address region. These trap and VMM emulates
 - Especially device drivers
- All the Virtual Devices are provided by Host Application Process resources
 - Virtual memory, files, windows, sockets, ...

Recall: VMM – VM – VMM extension process



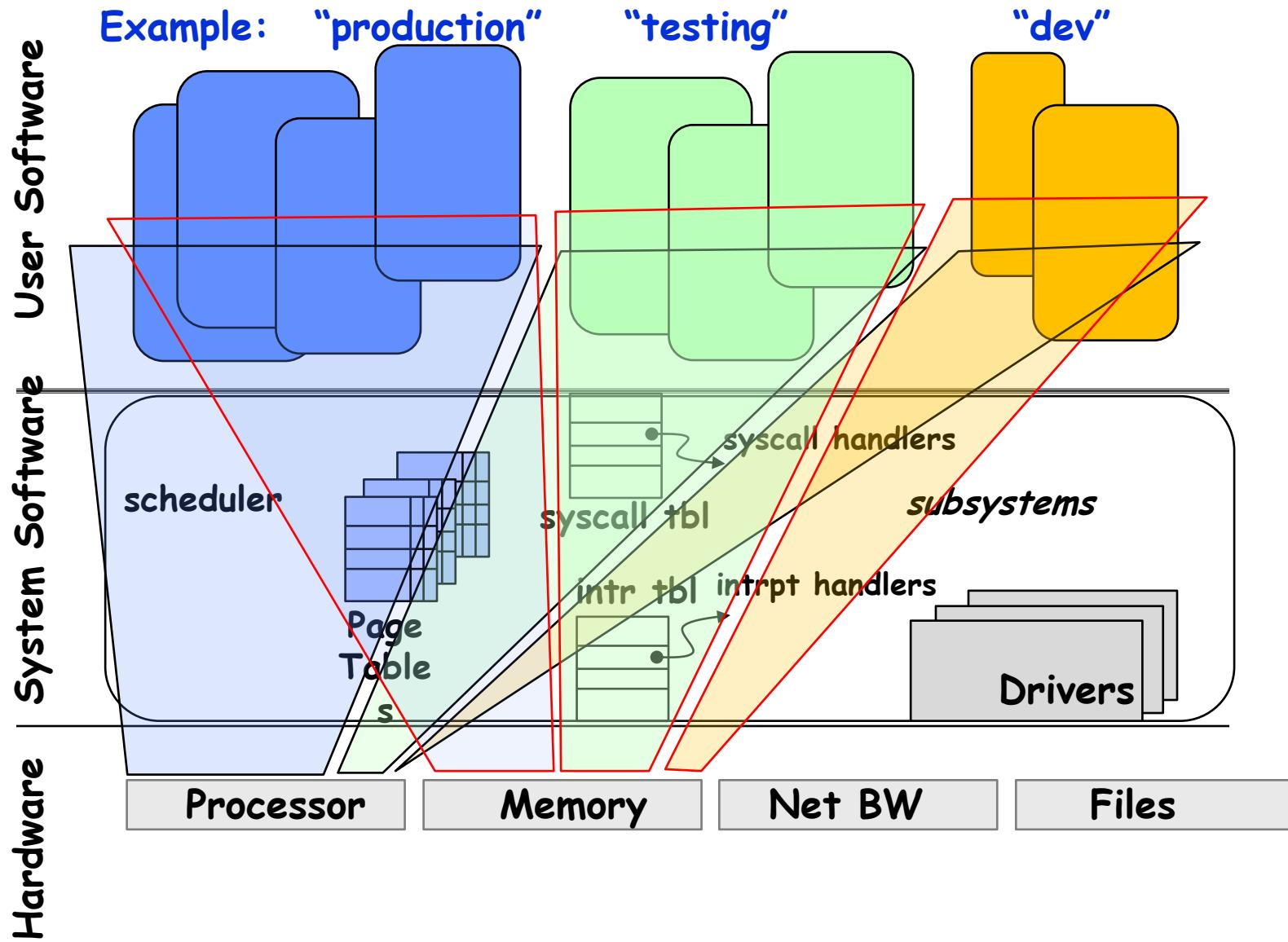
Virtual Machine Concepts

- all built on concepts you understand of classic operating systems
 - Processes, Threads, Page Tables, Interrupts, Syscall, RPC, Drivers
- *Shadow Page Table*
- VMM *interposes* on all Guest Proc \leftrightarrow Guest OS interactions
- Trap and emulate
- Virtual Devices by Host Application Process resources
- *Predictable allocation of shared resources*

Stepping Back

- In creating a virtual machine we configured a set of resources within which all of its activities – its OS and all its processes – would operate
 - Total amount of physical memory (the MMAP)
 - Total size of all its disk storage (the file backing its disk)
 - Total network bandwidth
- These constraints are valuable even without a distinct (possibly heterogenous) Guest OS
- Modern operating systems provide *performance isolation*, in addition to traditional protection isolation
 - Without the additional machinery and capability of VMs

Performance Isolation: CGroups



CGroups

- Identify collections of processes that will be treated as a *group* for resource allocation
 - Groups can have hierarchical structure
 - i.e., a Group can be comprised of subgroups
 - Process parent-child relationship defines a hierarchy
 - » Generalize this beyond `fork()`
- Set of key resource dimensions
 - Processor share (CPU), CPU set (bound to particular cores)
 - Physical memory share,
 - block IO, net priority, net class
 - Namespace (ns), i.e., containers
- Resource limiting, prioritization, accounting and control
- Containers define a collection of libraries and executables that should be a group

How is all the cgroups info recorded?

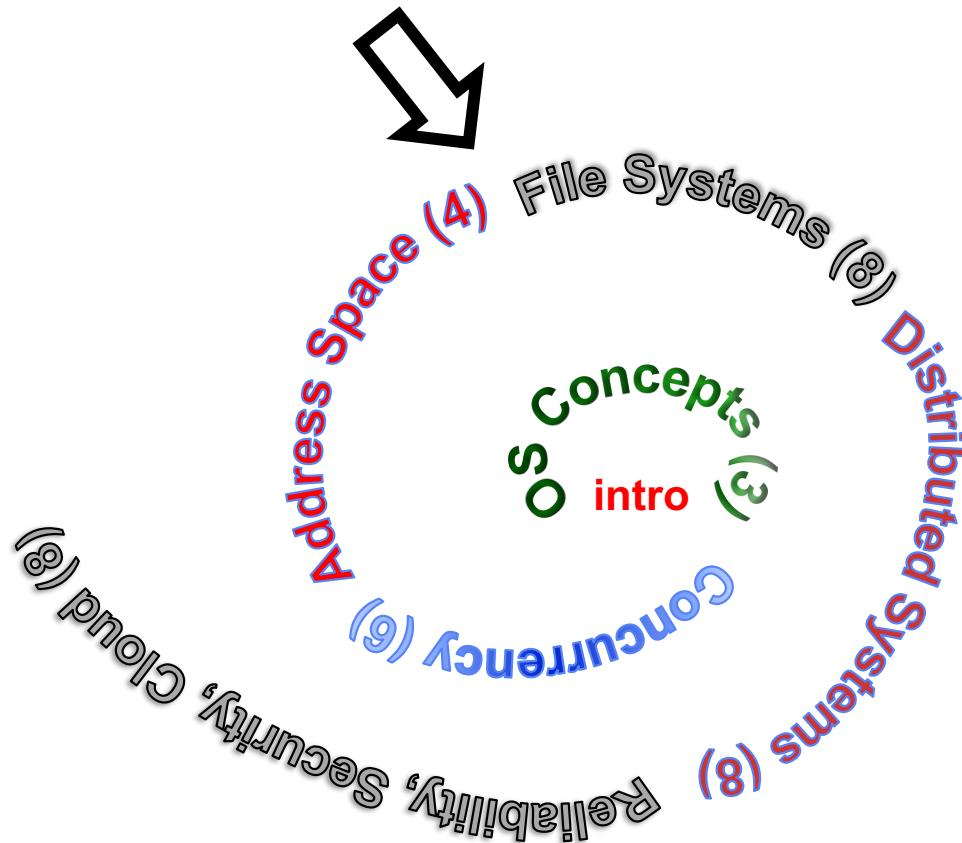
- Unix-based systems use `/proc` to represent information about processes
 - `/proc/<pid>` describes process `<pid>`
- `/proc/cgroups/*`
 - Describes what controllers are implemented
- Each cgroup controller has a directory under `/sys/fs/cgroup/<controller>`
 - `/sys/fs/cgroup/cpu/production`
 - `/sys/fs/cgroup/memory/foo/memory.limit_in_bytes`
- Kernel monitors and controls processes/threads in accordance with cgroup controllers

The Requirements of I/O

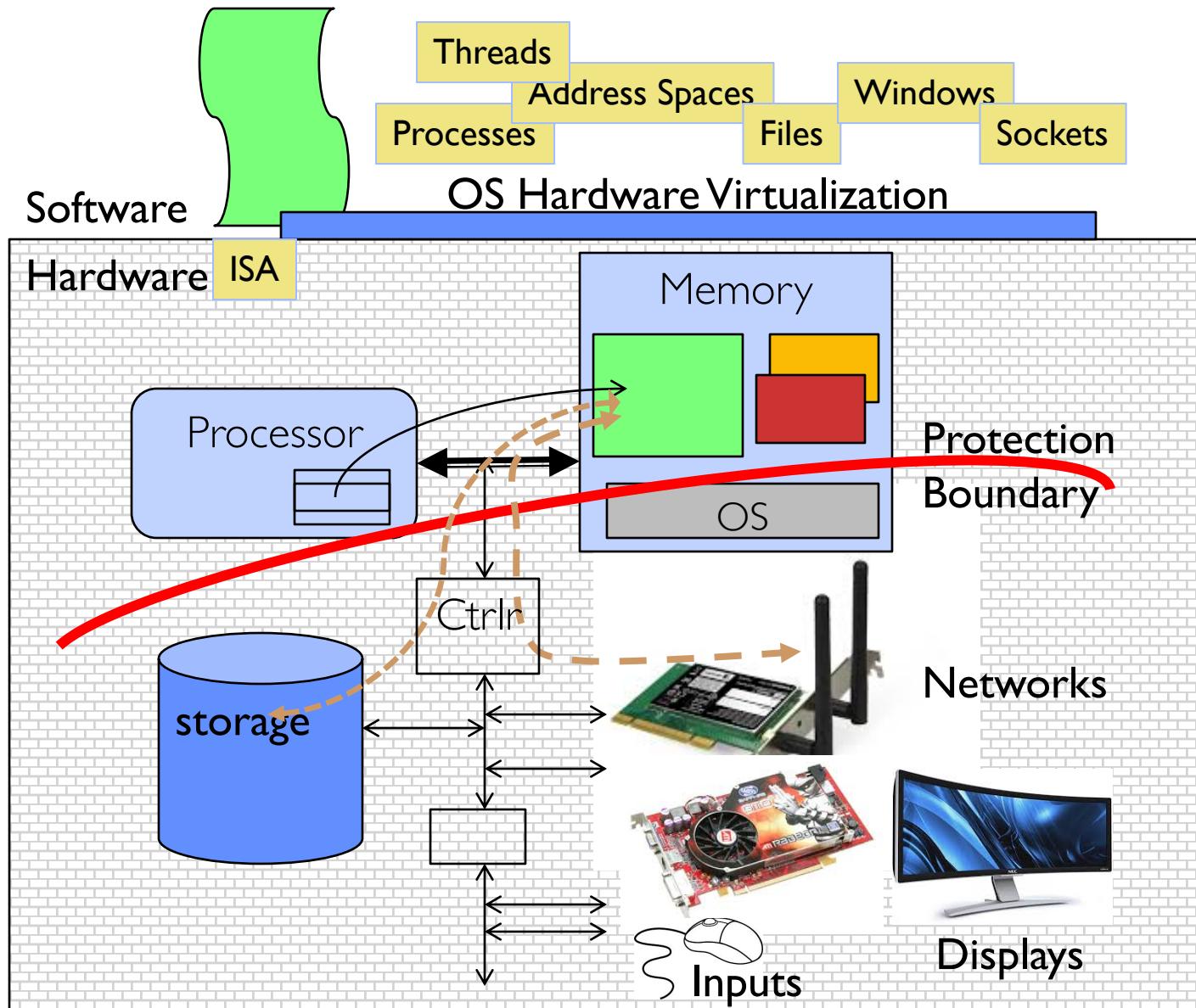
- So far in this course:
 - We have learned how to manage CPU and memory
- What about I/O?
 - Without I/O, computers are useless (disembodied brains?)
 - But... thousands of devices, each slightly different
 - » How can we standardize the interfaces to these devices?
 - Devices unreliable: media failures 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?



Course Structure: Spiral



Recall: OS Basics: I/O



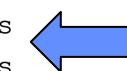
Recall: I/O is at every different timescale

And Range of Timescales



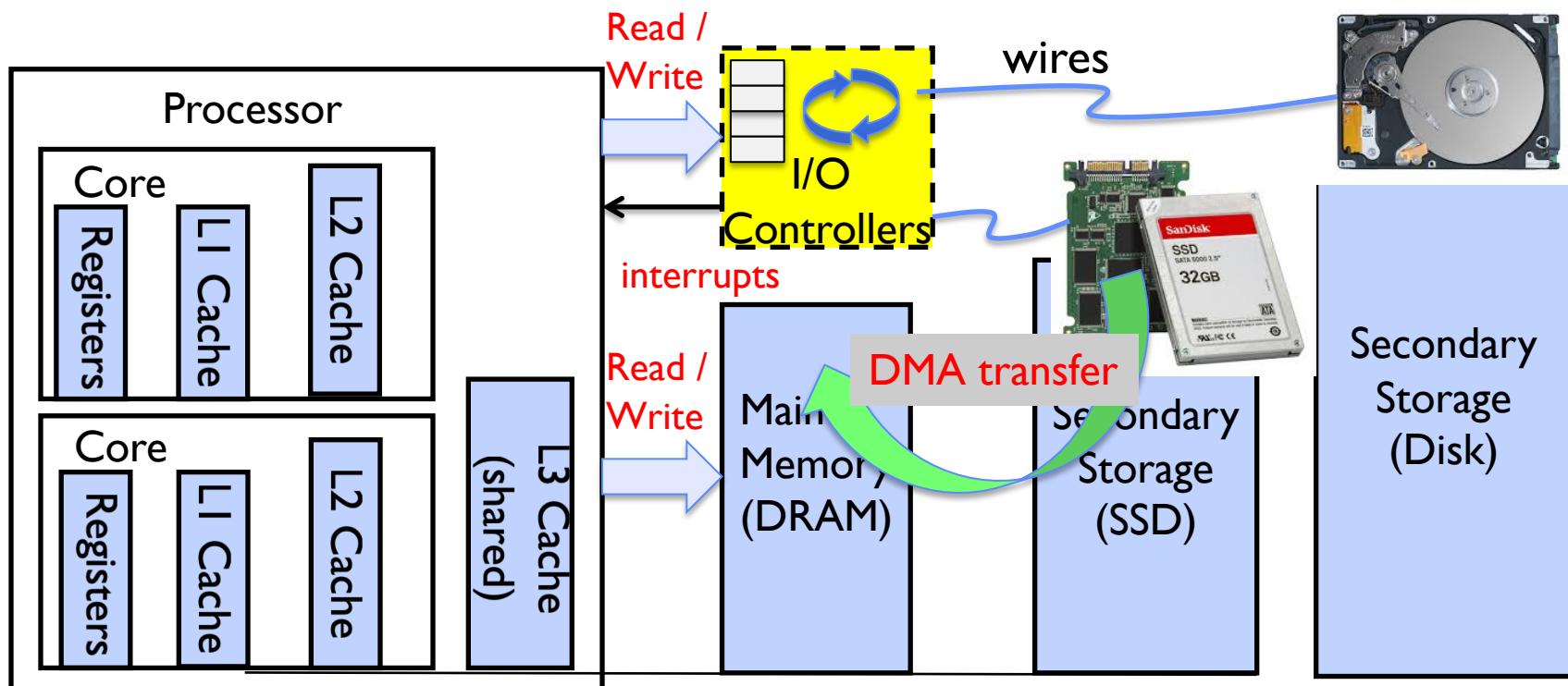
Jeff Dean: "Numbers Everyone Should Know"

L1 cache reference	0.5 ns
Branch mispredict	5 ns
L2 cache reference	7 ns
Mutex lock/unlock	25 ns
Main memory reference	100 ns
Compress 1K bytes with Zippy	3,000 ns
Send 2K bytes over 1 Gbps network	20,000 ns
Read 1 MB sequentially from memory	250,000 ns
Round trip within same datacenter	500,000 ns
Disk seek	10,000,000 ns
Read 1 MB sequentially from disk	20,000,000 ns
Send packet CA->Netherlands->CA	150,000,000 ns



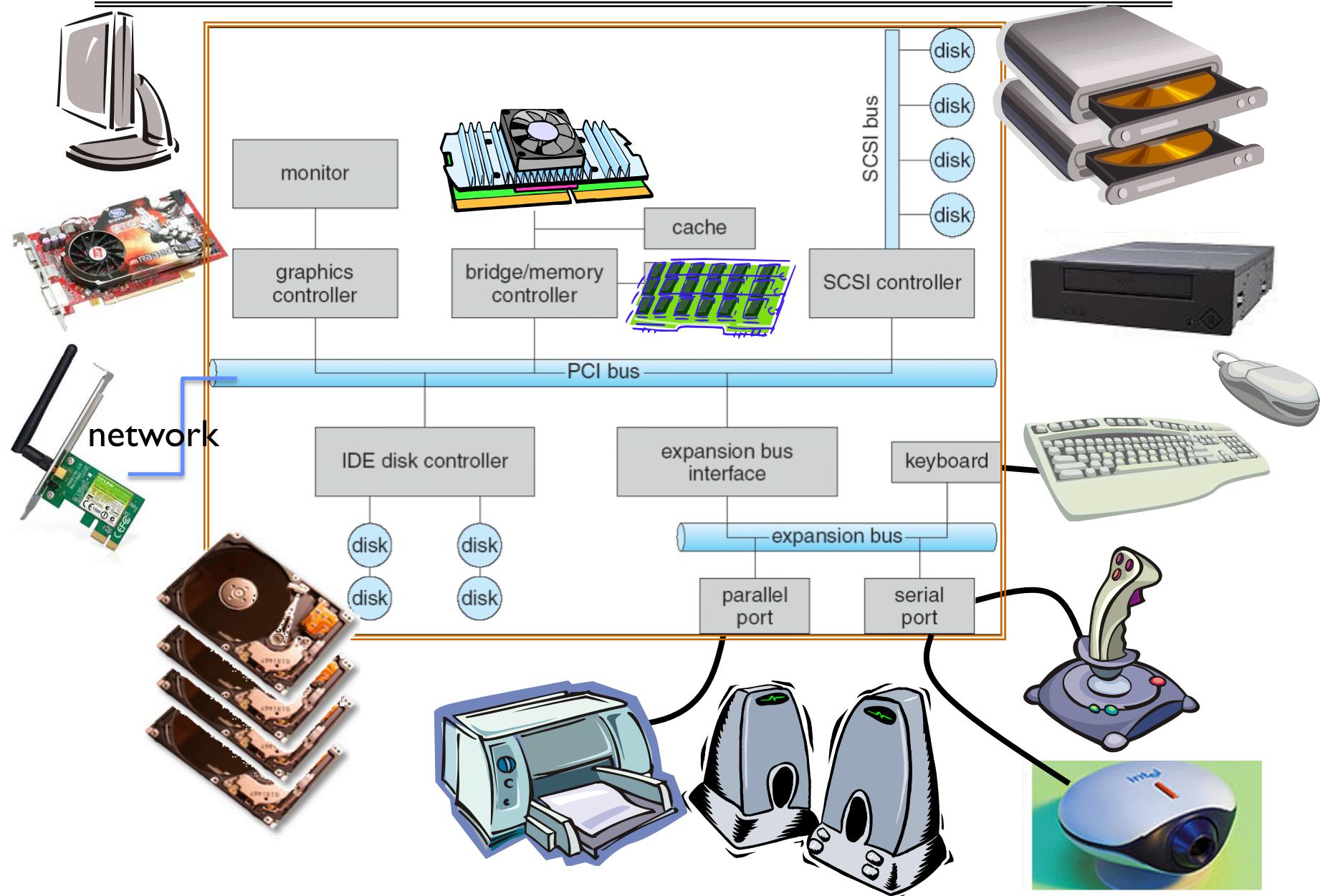
Key Stroke / Click
100 ms

In a Picture

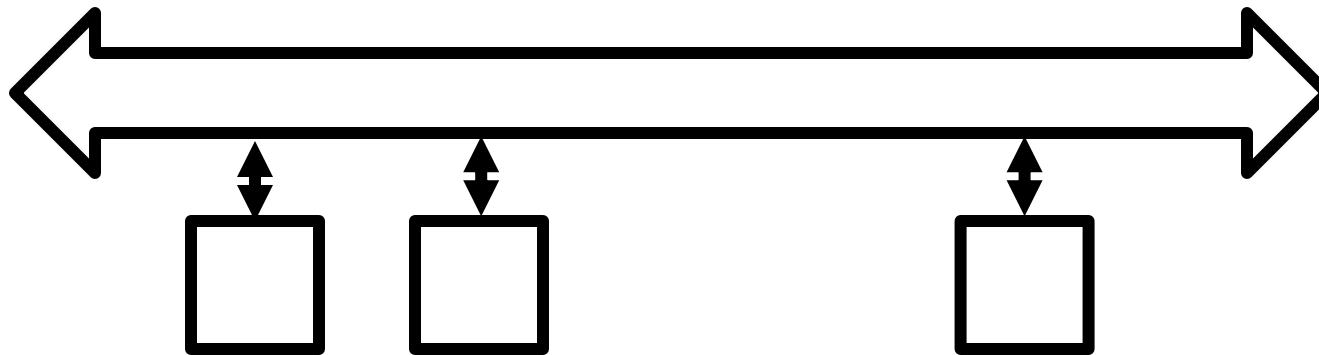


- I/O devices you recognize are supported by I/O Controllers
- Processors access them by reading and writing IO registers as if they were memory
 - Write commands and arguments, read status and results

Modern I/O Systems

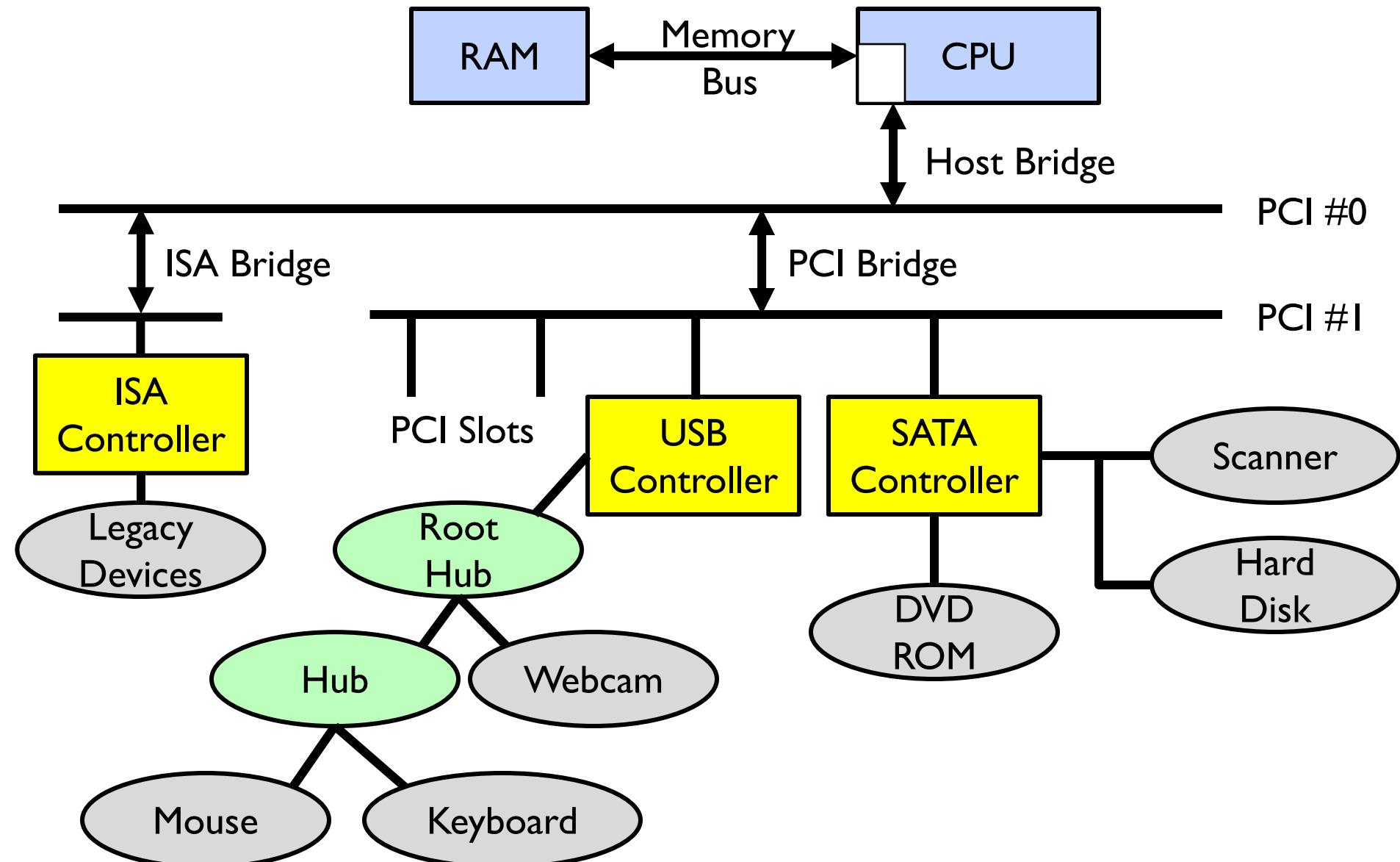


What's a bus?



- Common set of wires for communication among hardware devices plus protocols for carrying out data transfer transactions
 - Operations: e.g., Read, Write
 - Control lines, Address lines, Data lines
 - Typically multiple devices
- Protocol: initiator requests access, arbitration to grant, identification of recipient, handshake to convey address, length, data
- Very high BW close to processor (wide, fast, and inflexible), low BW with high flexibility out in I/O subsystem

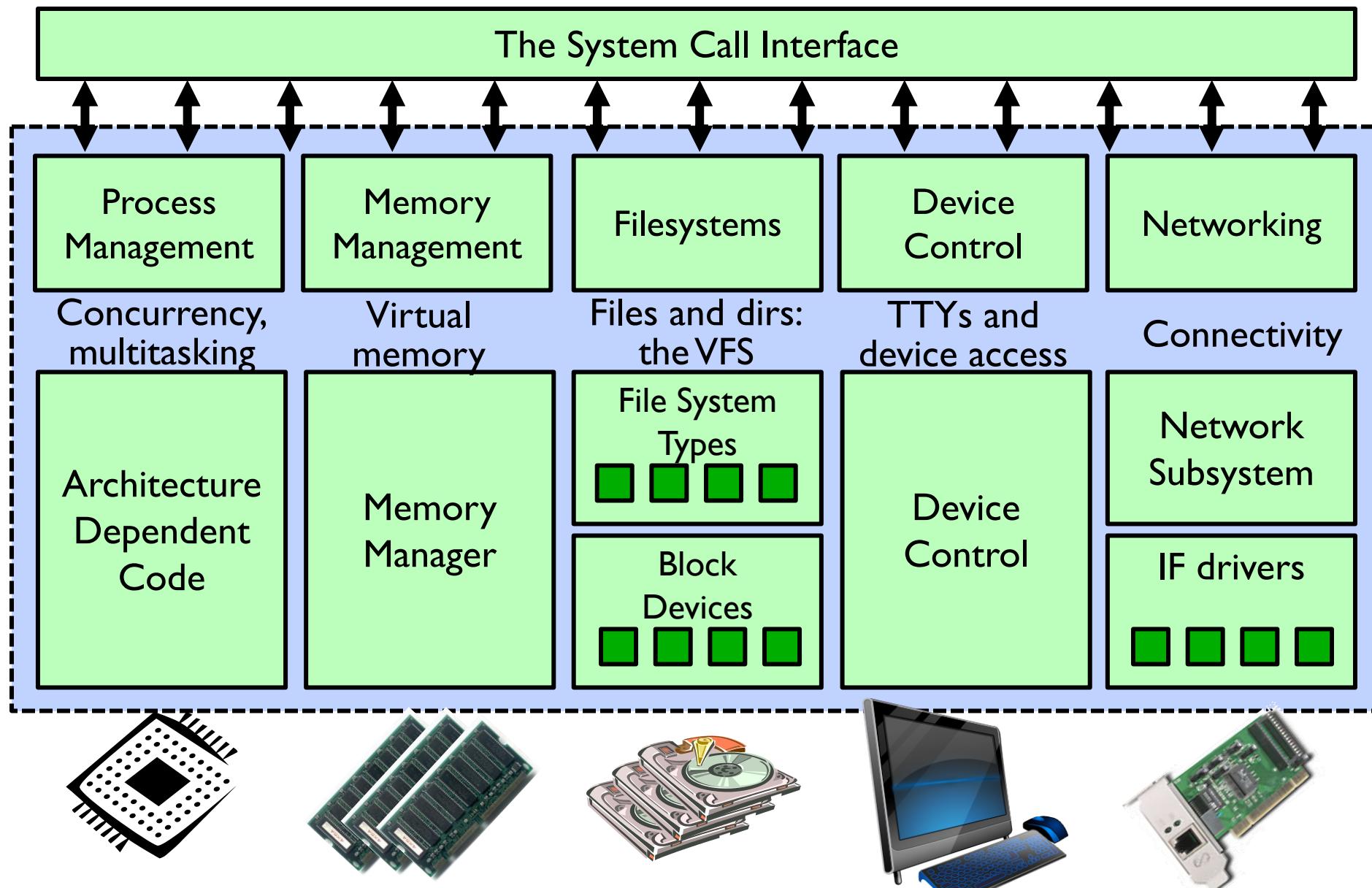
Example: Typical PCI Architecture



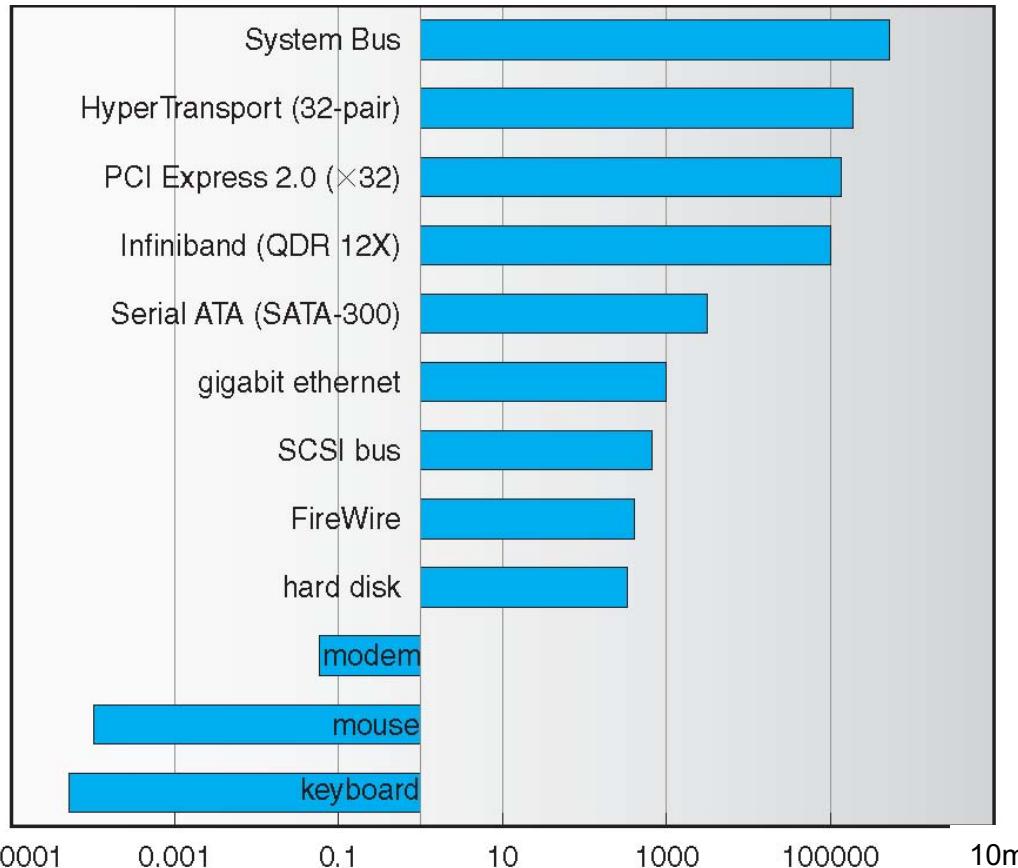
Operational Parameters for I/O

- 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 (1st peek)

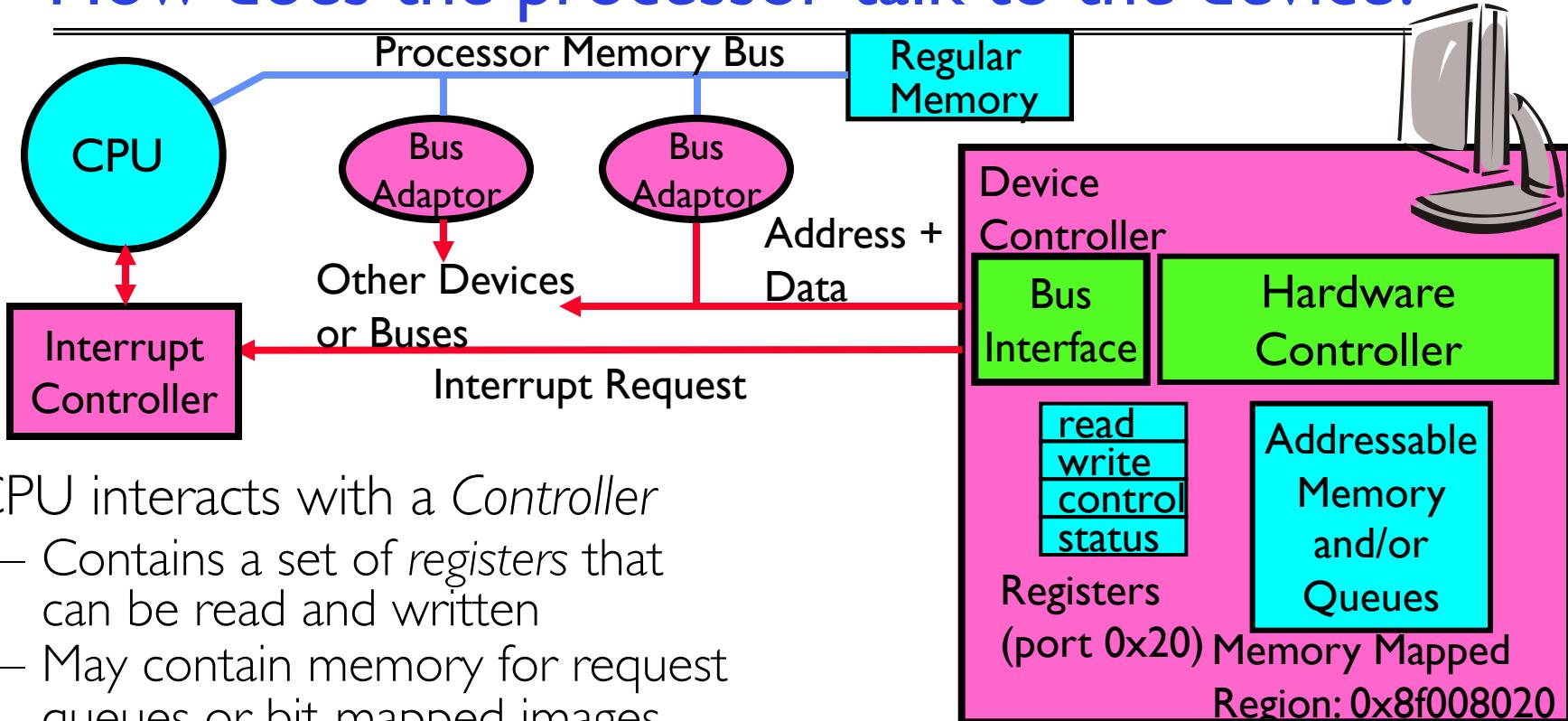


Example Device-Transfer Rates in Mb/s (Sun Enterprise 6000)



- 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

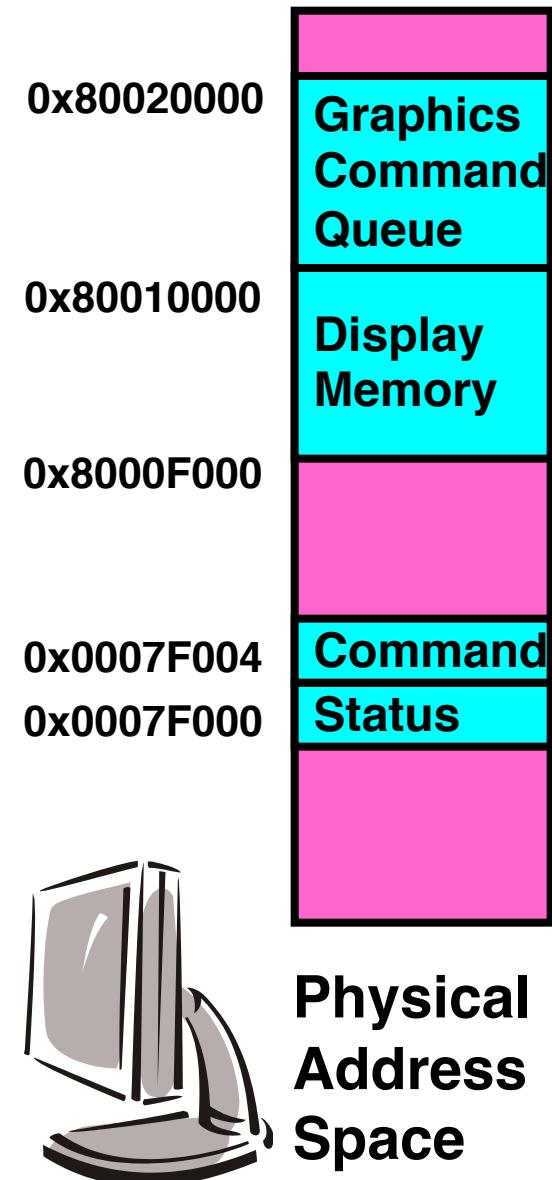
How does the processor talk to the device?



- 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

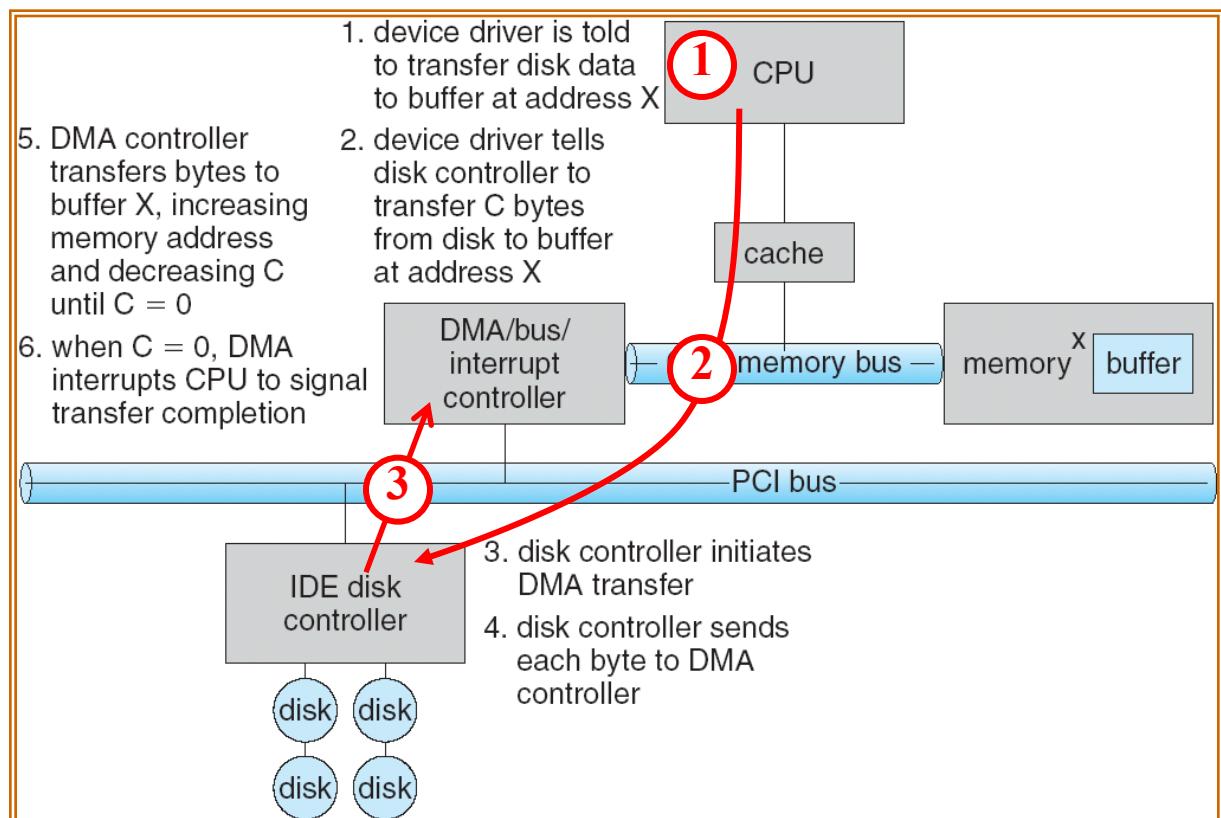
Example: Memory-Mapped Display Controller

- 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



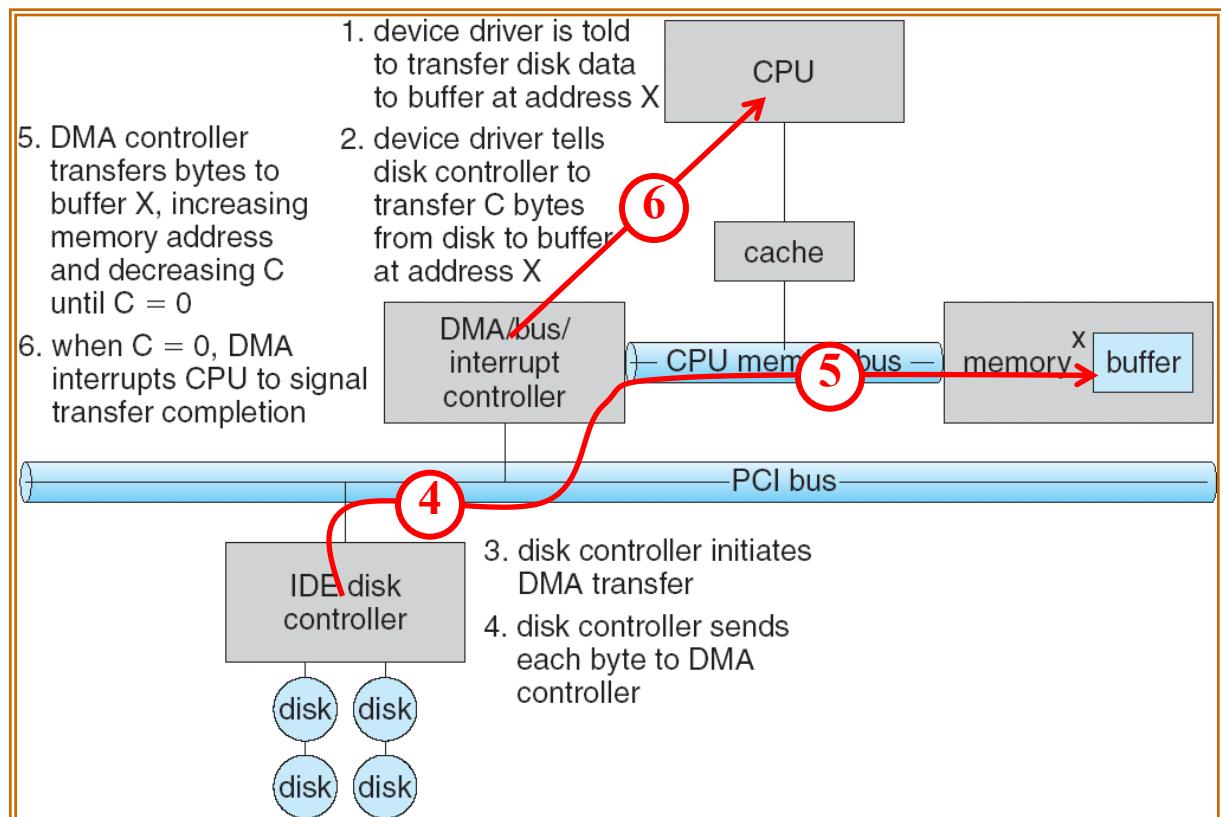
Transferring Data To/From Controller

- 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):



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):



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

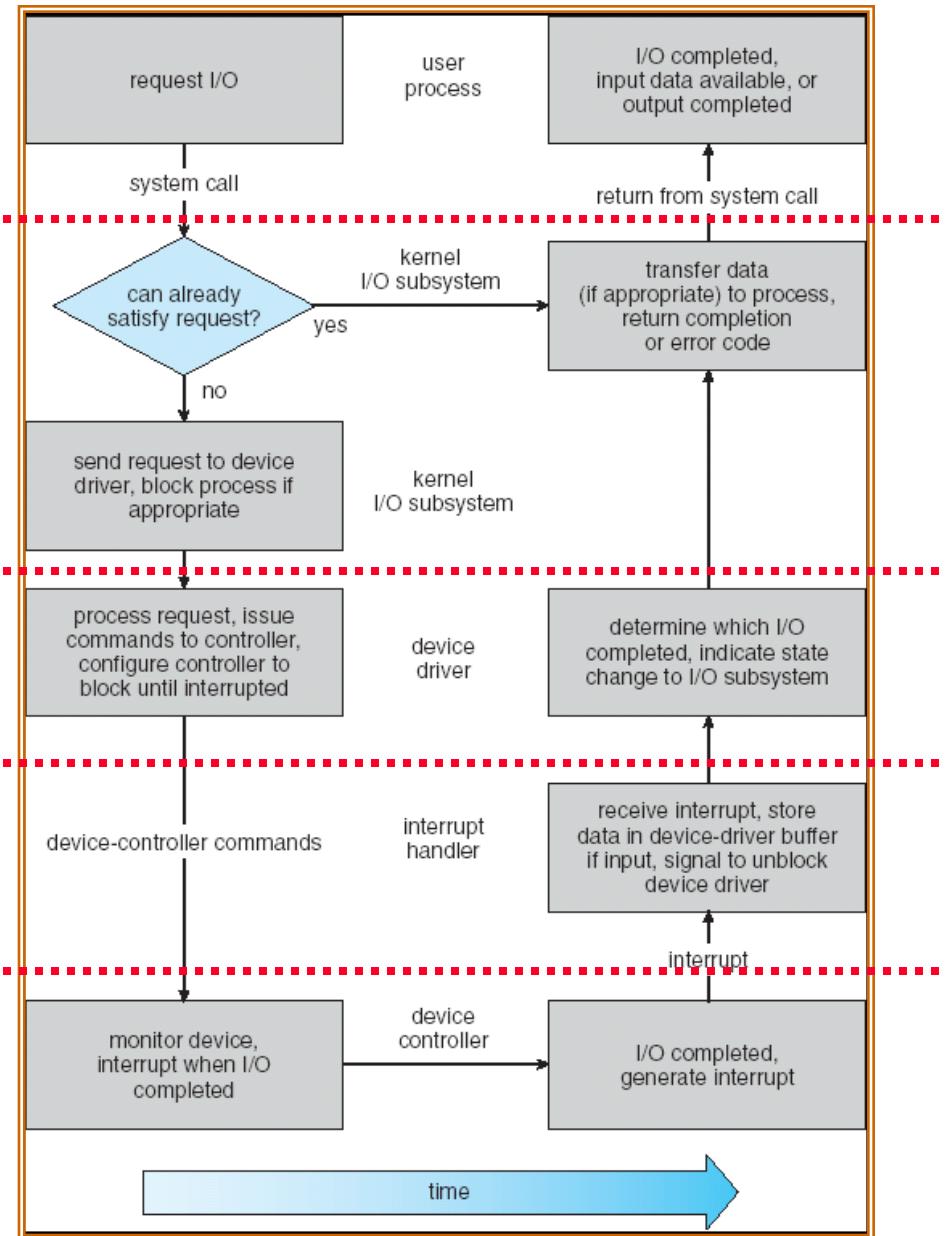
User
Program

Kernel I/O
Subsystem

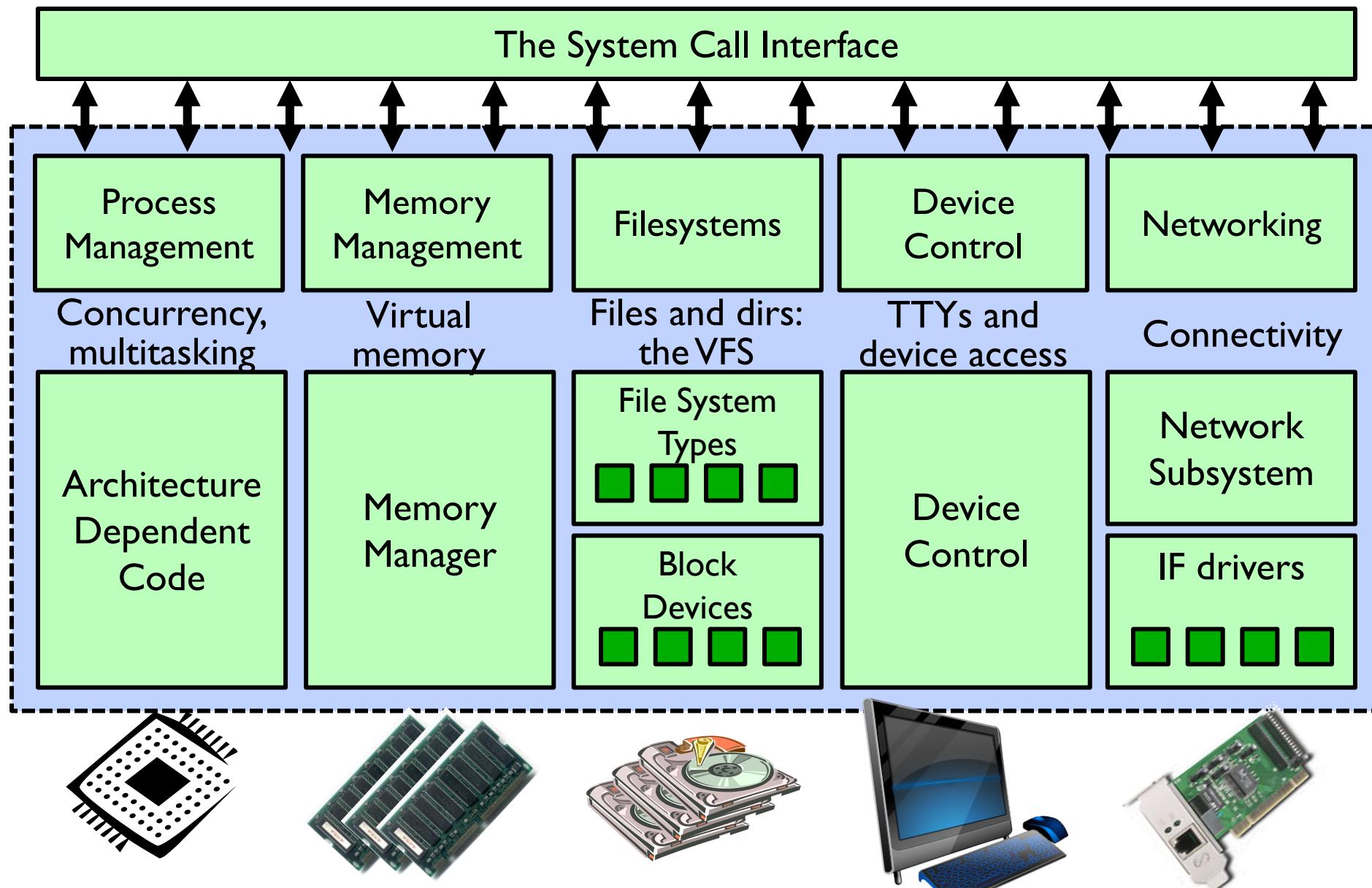
Device Driver
Top Half

Device Driver
Bottom Half

Device
Hardware



Kernel Device Structure (1st peek)



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!

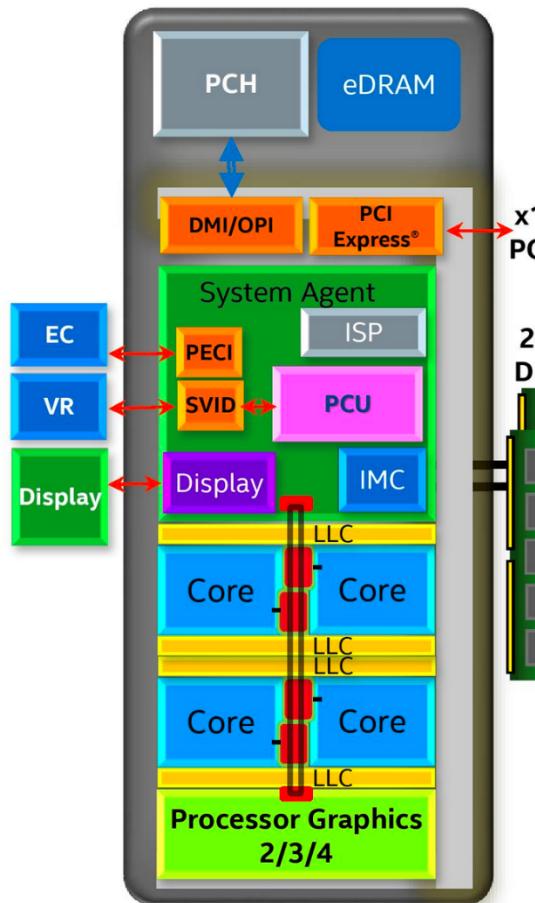
Want Standard Interfaces to Devices

- **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?

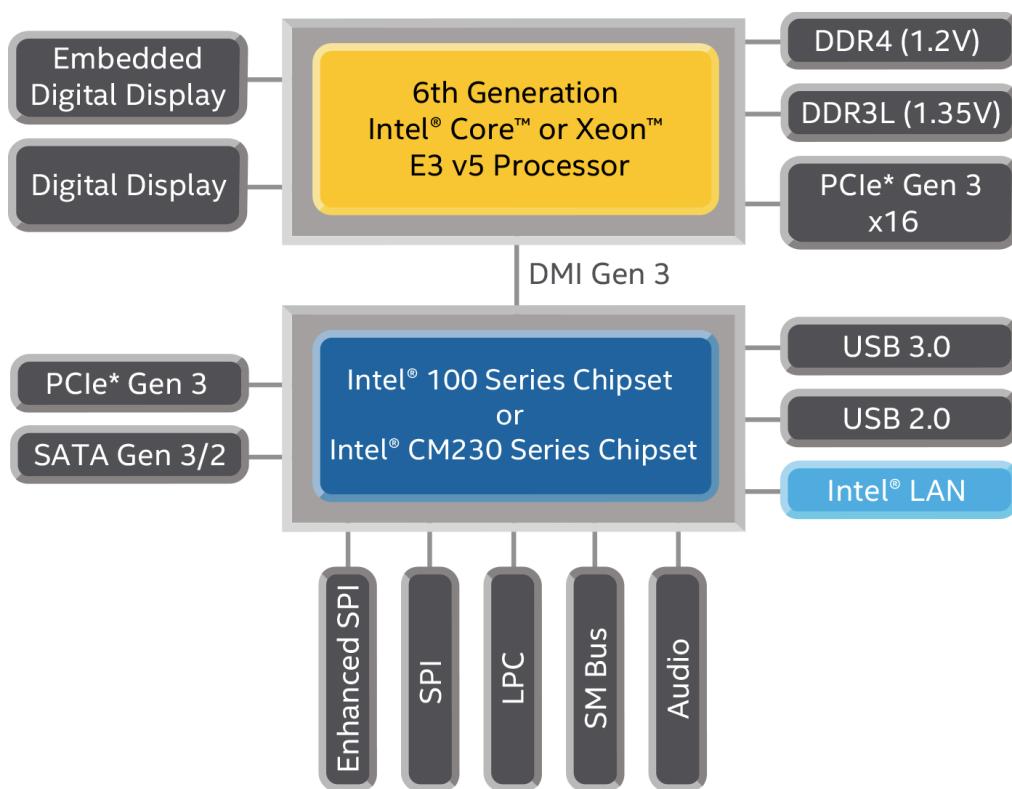
- **Blocking Interface:** “Wait”
 - When request data (e.g. `read()` system call), put process to sleep until data is ready
 - When write data (e.g. `write()` system call), put process to sleep until device is ready for data
- **Non-blocking Interface:** “Don’t Wait”
 - Returns quickly from read or write 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

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 DRAM
 - High-speed PCI-Express (for Graphics cards)
 - Direct Media Interface (DMI) Connection to PCH (Platform Control Hub)

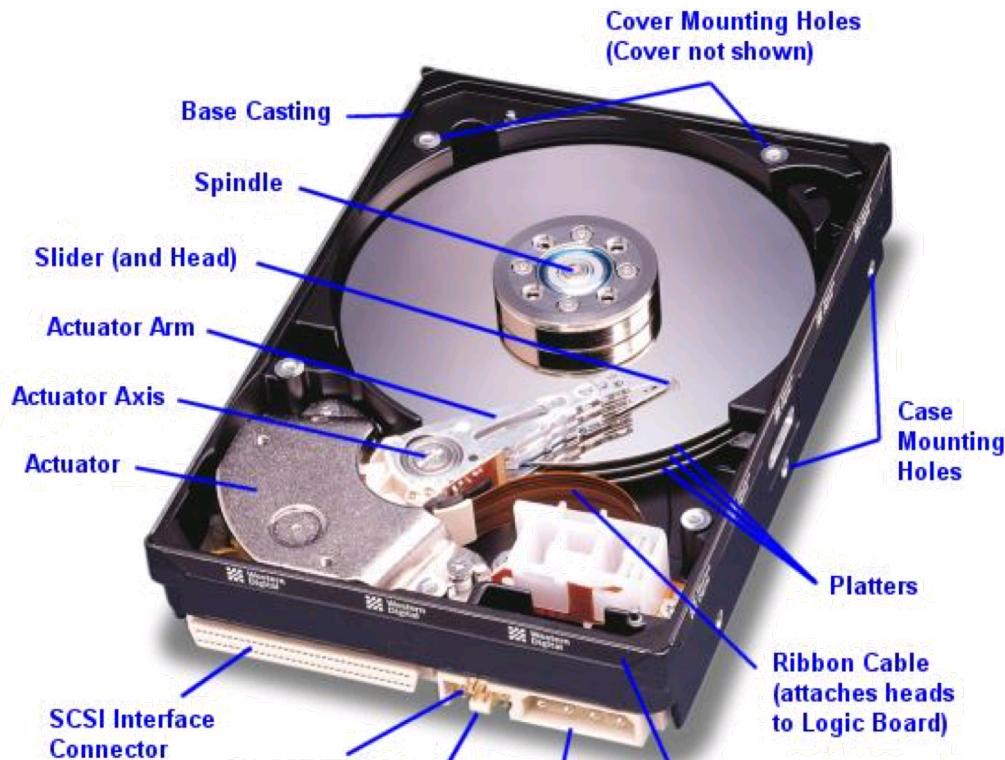
Sky Lake I/O: "platform controller hub" (PCH)



- 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)

Sky Lake
System Configuration

Hard Disk Drives (HDDs)



Western Digital Drive

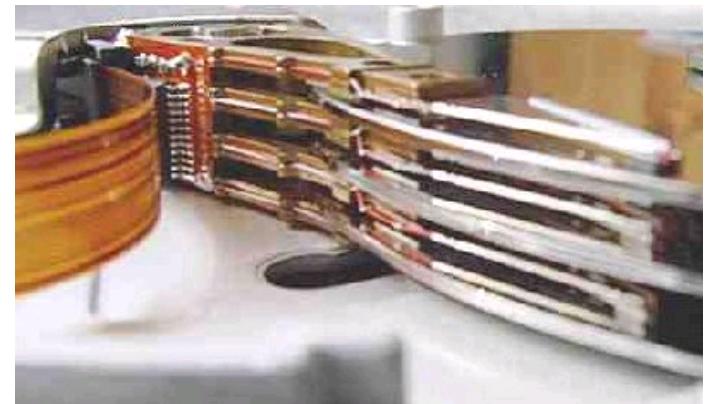
<http://www.storagereview.com/guide/>

IBM Personal Computer/AT (1986)

30 MB hard disk - \$500

30-40ms seek time

0.7-1 MB/s (est.)



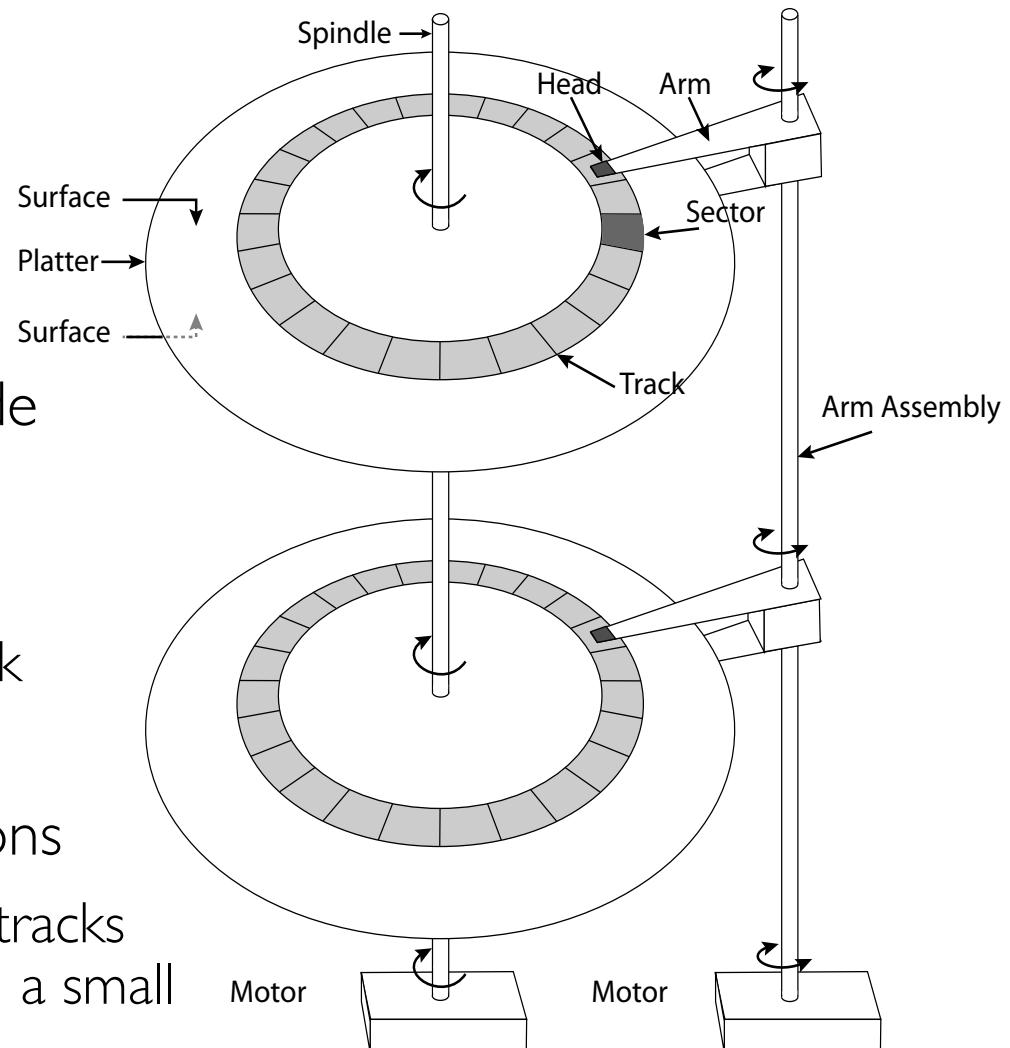
**Read/Write Head
Side View**



IBM/Hitachi Microdrive

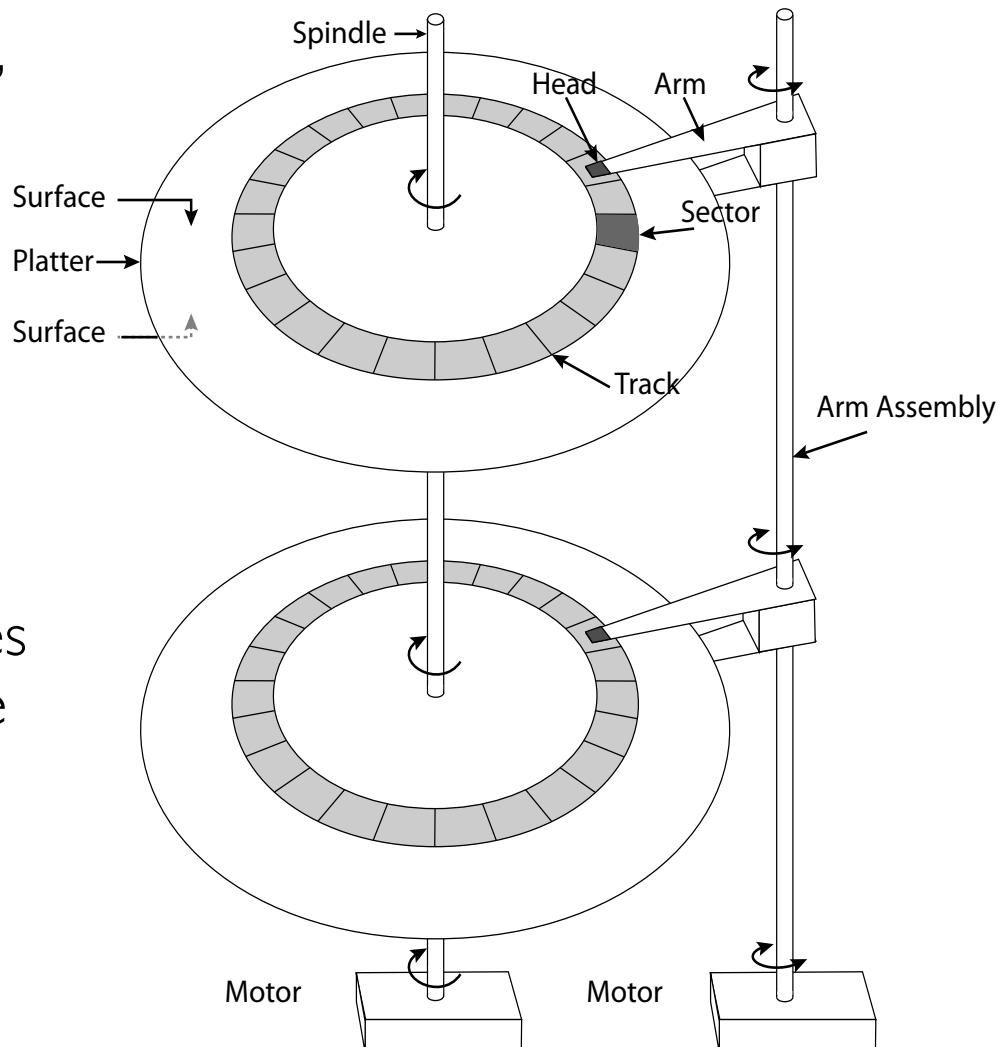
The Amazing Magnetic Disk

- Unit of Transfer: Sector
 - Ring of sectors form a track
 - Stack of tracks form a cylinder
 - Heads position on cylinders
- Disk Tracks $\sim 1 \mu\text{m}$ (micron) wide
 - Wavelength of light is $\sim 0.5 \mu\text{m}$
 - Resolution of human eye: $50 \mu\text{m}$
 - 100K tracks on a typical 2.5" disk
- Separated by unused guard regions
 - Reduces likelihood neighboring tracks are corrupted during writes (still a small non-zero chance)



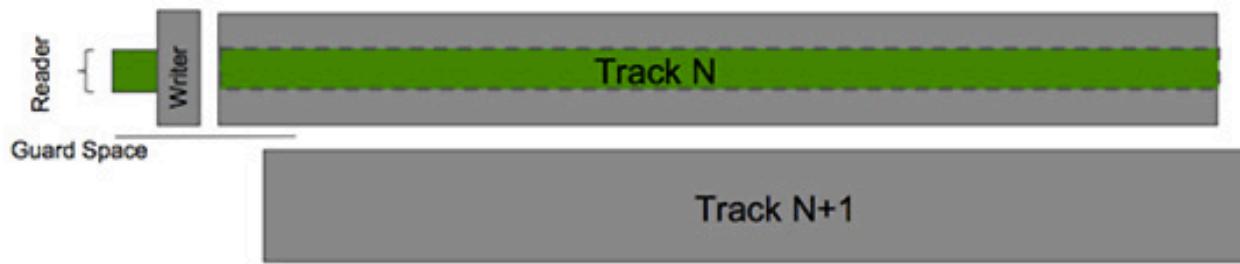
The Amazing Magnetic Disk

- Track length varies across disk
 - Outside: More sectors per track, higher bandwidth
 - Disk is organized into regions of tracks with same # of sectors/track
 - Only outer half of radius is used
 - » Most of the disk area in the outer regions of the disk
- Disks so big that some companies (like Google) reportedly only use part of disk for active data
 - Rest is archival data

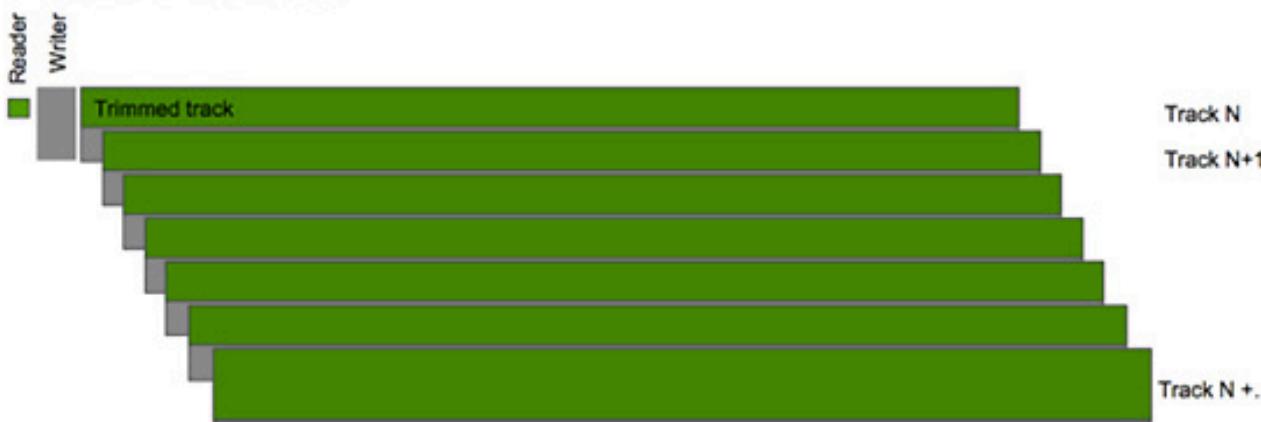


Shingled Magnetic Recording (SMR)

Conventional Writes



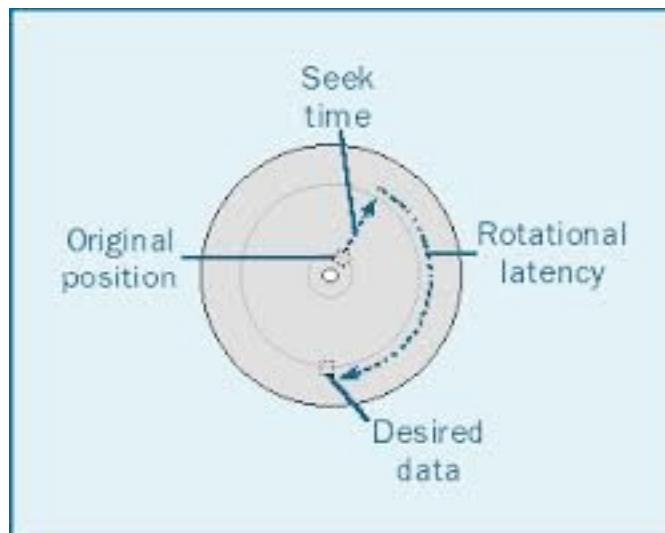
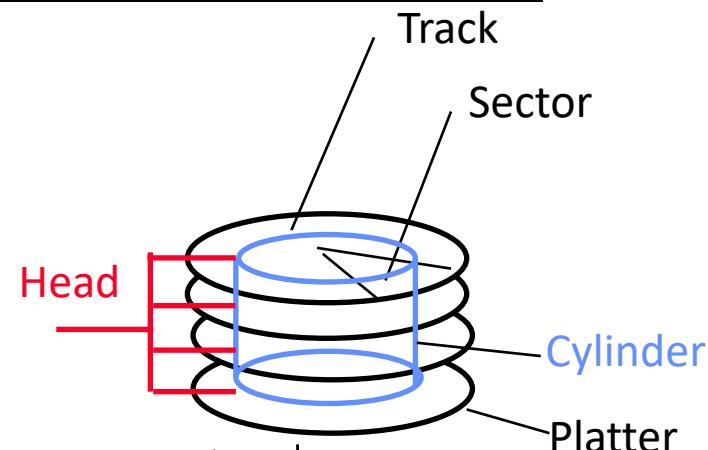
SMR Writes



- Overlapping tracks yields greater density, capacity
- Restrictions on writing, complex DSP for reading
- Examples: Seagate (8TB), Hitachi (10TB)

Review: Magnetic Disks

- **Cylinders:** all the tracks under the head at a given point on all surface
- Read/write data is a three-stage process:
 - **Seek time:** position the head/arm over the proper track
 - **Rotational latency:** wait for desired sector to rotate under r/w head
 - **Transfer time:** transfer a block of bits (sector) under r/w head

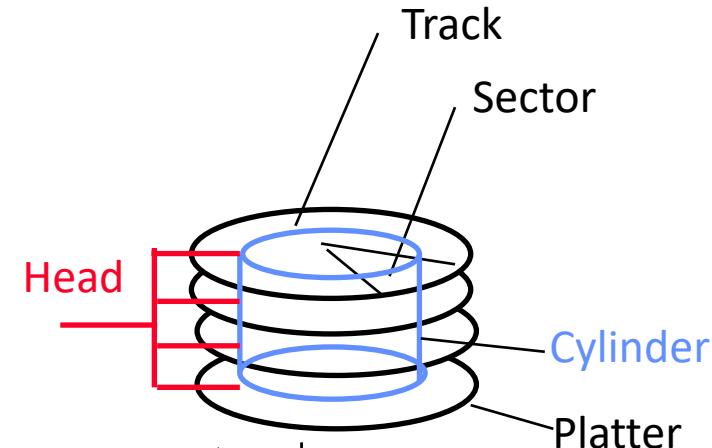


Seek time = 4-8ms
One rotation = 1-2ms
(3600-7200 RPM)

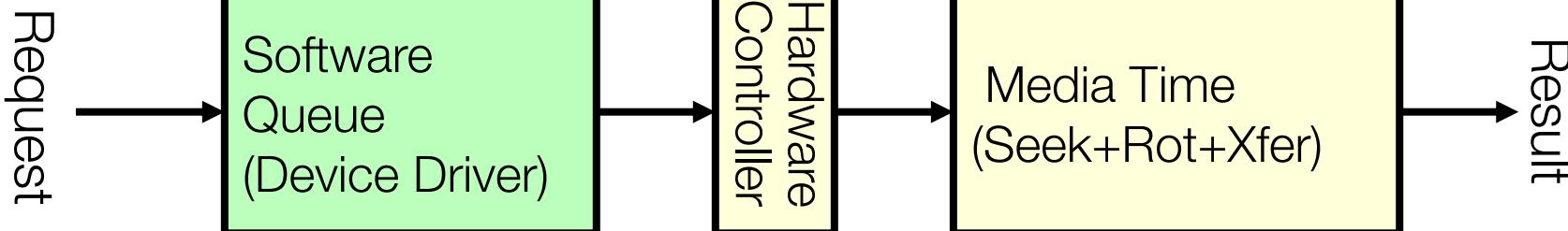
Review: Magnetic Disks

- **Cylinders:** all the tracks under the head at a given point on all surface

- Read/write data is a three-stage process:
 - **Seek time:** position the head/arm over the proper track
 - **Rotational latency:** wait for desired sector to rotate under r/w head
 - **Transfer time:** transfer a block of bits (sector) under r/w head



$$\text{Disk Latency} = \text{Queueing Time} + \text{Controller time} + \\ \text{Seek Time} + \text{Rotation Time} + \text{Xfer Time}$$



Typical Numbers for Magnetic Disk

Parameter	Info / Range
Space/Density	Space: 14TB (Seagate), 8 platters, in 3½ inch form factor! Areal Density: ≥ 1 Terabit/square inch! (PMR, Helium, ...)
Average seek time	Typically 4-6 milliseconds. Depending on reference locality, actual cost may be 25-33% of this number.
Average rotational latency	Most laptop/desktop disks rotate at 3600-7200 RPM (16-8 ms/rotation). Server disks up to 15,000 RPM. Average latency is halfway around disk so 8-4 milliseconds
Controller time	Depends on controller hardware
Transfer time	Typically 50 to 250 MB/s. Depends on: <ul style="list-style-type: none">Transfer size (usually a sector): 512B – 1KB per sectorRotation speed: 3600 RPM to 15000 RPMRecording density: bits per inch on a trackDiameter: ranges from 1 in to 5.25 in
Cost	Used to drop by a factor of two every 1.5 years (or even faster); now slowing down

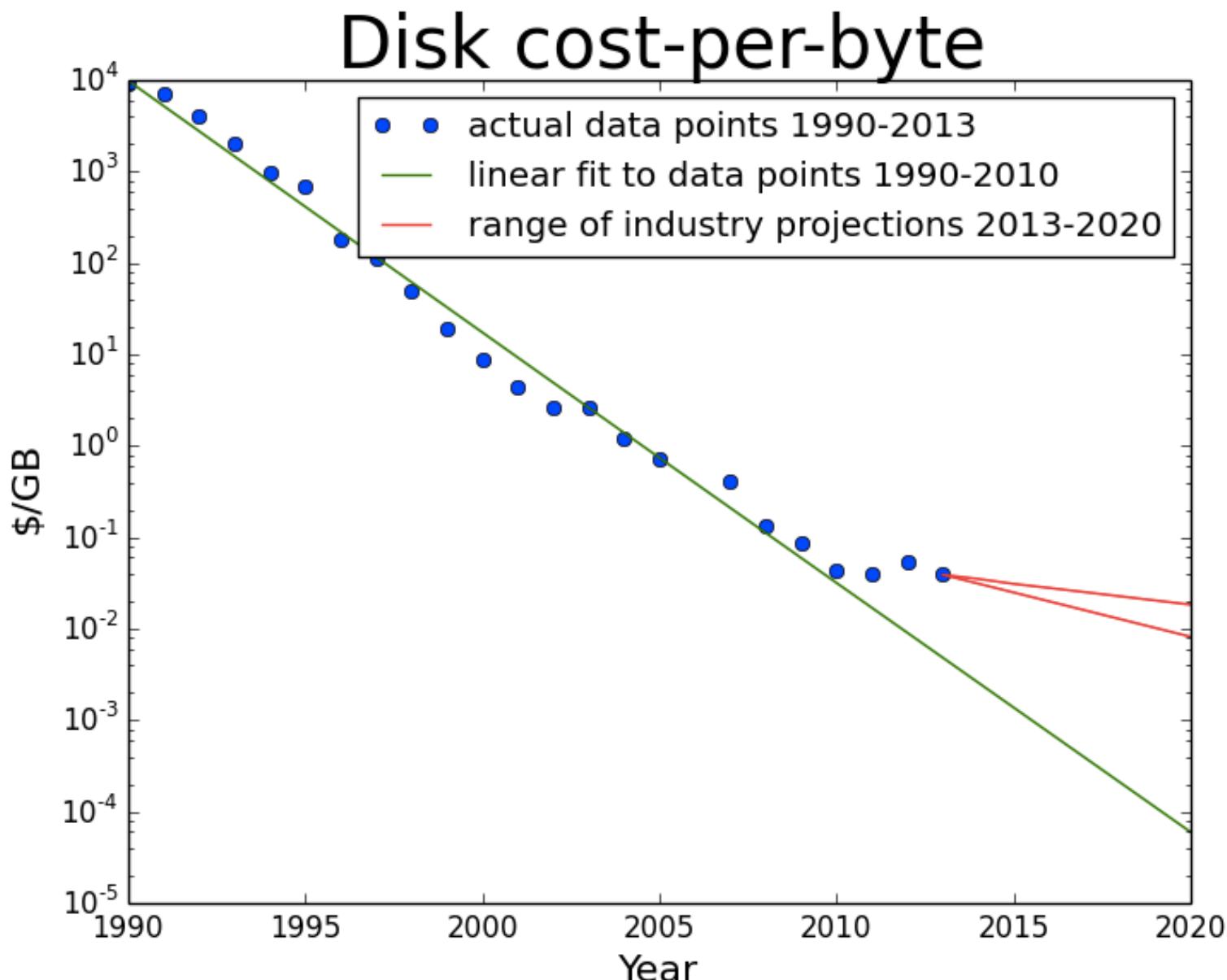
Disk Performance Example

- Assumptions:
 - Ignoring queuing and controller times for now
 - Avg seek time of 5ms,
 - 7200RPM \Rightarrow Time for rotation: $60000 \text{ (ms/min)} / 7200(\text{rev/min}) \sim= 8\text{ms}$
 - Transfer rate of 50MByte/s, block size of 4Kbyte \Rightarrow
 $4096 \text{ bytes} / 50 \times 10^6 \text{ (bytes/s)} = 81.92 \times 10^{-6} \text{ sec} \cong 0.082 \text{ ms for 1 sector}$
- Read block from random place on disk:
 - Seek (5ms) + Rot. Delay (4ms) + Transfer (0.082ms) = 9.082ms
 - Approx 9ms to fetch/put data: $4096 \text{ bytes} / 9.082 \times 10^{-3} \text{ s} \cong 451 \text{ KB/s}$
- Read block from random place in same cylinder:
 - Rot. Delay (4ms) + Transfer (0.082ms) = 4.082ms
 - Approx 4ms to fetch/put data: $4096 \text{ bytes} / 4.082 \times 10^{-3} \text{ s} \cong 1.03 \text{ MB/s}$
- Read next block on same track:
 - Transfer (0.082ms): $4096 \text{ bytes} / 0.082 \times 10^{-3} \text{ s} \cong 50 \text{ MB/sec}$
- Key to using disk effectively (especially for file systems) is to minimize seek and rotational delays

(Lots of) Intelligence in the Controller

- Sectors contain sophisticated error correcting codes
 - Disk head magnet has a field wider than track
 - Hide corruptions due to neighboring track writes
- Sector sparing
 - Remap bad sectors transparently to spare sectors on the same surface
- Slip sparing
 - Remap all sectors (when there is a bad sector) to preserve sequential behavior
- Track skewing
 - Sector numbers offset from one track to the next, to allow for disk head movement for sequential ops
- ...

Hard Drive Prices over Time



Example of Current HDDs

- Seagate Exos X14 (2018)
 - 14 TB hard disk
 - » 8 platters, 16 heads
 - » Helium filled: reduce friction and power
 - 4.16ms average seek time
 - 4096 byte physical sectors
 - 7200 RPMs
 - 6 Gbps SATA /12Gbps SAS interface
 - » 261MB/s MAX transfer rate
 - » Cache size: 256MB
 - Price: \$615 (< \$0.05/GB)



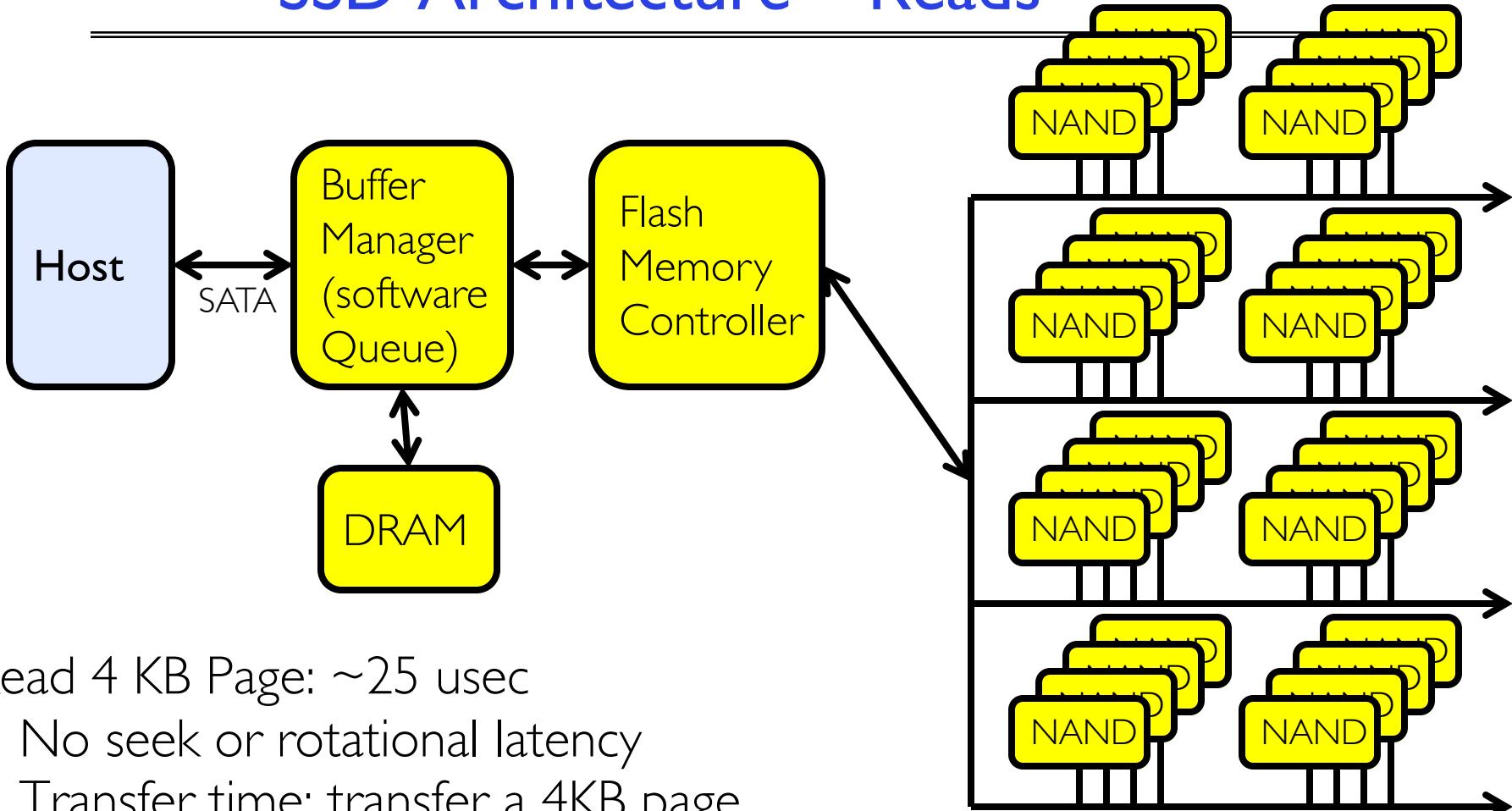
- IBM Personal Computer/AT (1986)
 - 30 MB hard disk
 - 30-40ms seek time
 - 0.7-1 MB/s (est.)
 - Price: \$500 (\$17K/GB, 340,000x more expensive !!)

Solid State Disks (SSDs)



- 1995 – Replace rotating magnetic media with non-volatile memory (battery backed DRAM)
- 2009 – Use NAND Multi-Level Cell (2 or 3-bit/cell) flash memory
 - Sector (4 KB page) addressable, but stores 4-64 “pages” per memory block
 - Trapped electrons distinguish between 1 and 0
- No moving parts (no rotate/seek motors)
 - Eliminates seek and rotational delay (0.1-0.2ms access time)
 - Very low power and lightweight
 - Limited “write cycles”
- Rapid advances in capacity and cost ever since!

SSD Architecture – Reads

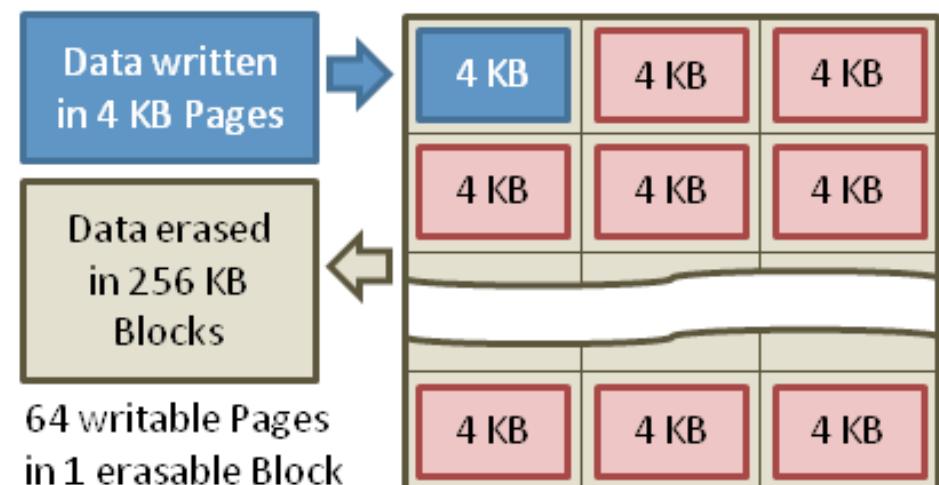


Read 4 KB Page: ~25 usec

- No seek or rotational latency
- Transfer time: transfer a 4KB page
 - » SATA: $300\text{-}600\text{MB/s} \Rightarrow \sim 4 \times 10^3 \text{ b} / 400 \times 10^6 \text{ bps} \Rightarrow 10 \text{ us}$
- Latency = Queuing Time + Controller time + Xfer Time
- Highest Bandwidth: Sequential OR Random reads

SSD Architecture – Writes

- Writing data is complex! ($\sim 200\mu\text{s} - 1.7\text{ms}$)
 - Can only write empty pages in a block
 - Erasing a block takes $\sim 1.5\text{ms}$
 - Controller maintains pool of empty blocks by coalescing used pages (read, erase, write), also reserves some % of capacity
- Rule of thumb: writes 10x reads, erasure 10x writes



Typical NAND Flash Pages and Blocks

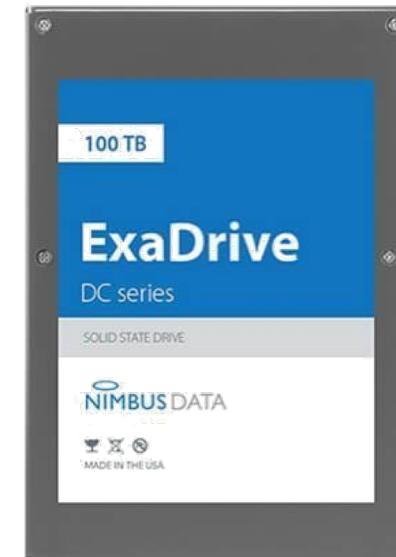
https://en.wikipedia.org/wiki/Solid-state_drive

Some “Current” 3.5in SSDs

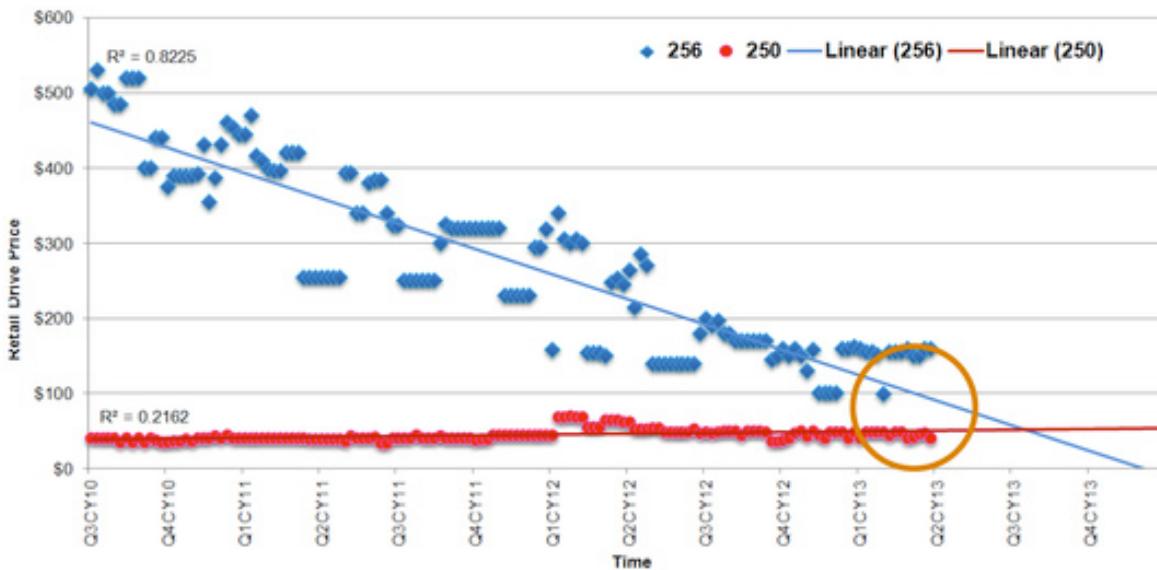
- Seagate Nytro SSD: 15TB (2017)
 - Dual 12Gb/s interface
 - Seq reads 860MB/s
 - Seq writes 920MB/s
 - Random Reads (IOPS): 102K
 - Random Writes (IOPS): 15K
 - Price (Amazon): \$6325 (\$0.41/GB)



- Nimbus SSD: 100TB (2019)
 - Dual port: 12Gb/s interface
 - Seq reads/writes: 500MB/s
 - Random Read Ops (IOPS): 100K
 - *Unlimited writes for 5 years!*
 - Price: ~ \$50K? (\$0.50/GB)



HDD vs SSD Comparison



SSD vs HDD

Usually 10 000 or 15 000 rpm SAS drives

0.1 ms

Access times

SSDs exhibit virtually no access time

5.5 ~ 8.0 ms

6000 io/s

Random I/O Performance

HDDs reach up to

400 io/s

SSDs have a failure rate of less than
0.5 %

SSDs are at least 15 times faster than HDDs

HDD's failure rate fluctuates between
2 ~ 5 %

SSDs have a failure rate of less than
0.5 %

This makes SSDs 4 - 10 times more reliable

HDD's failure rate fluctuates between
2 ~ 5 %

SSDs consume between
2 & 5 watts

SSDs consume between
6 & 15 watts

SSDs have an average I/O wait of
1 %

SSDs have an average I/O wait of
7 %

the average service time for an I/O request while running a backup remains below
20 ms

the I/O request time with HDDs during backup rises up to
400~500 ms

SSD backups take about
6 hours

HDD backups take up to
20~24 hours

Price Crossover Point for HDD and SSD

	2012	2013	2014	2015E	2016F	2017F
HDD	0.09	0.08	0.07	0.06	0.06	0.06
2.5" SSD	0.99	0.68	0.55	0.39	0.24	0.17

SSD prices drop much faster than HDD

Amusing calculation: Is a full Kindle heavier than an empty one?

- Actually, “Yes”, but not by much
- Flash works by trapping electrons:
 - So, erased state lower energy than written state
- Assuming that:
 - Kindle has 4GB flash
 - $\frac{1}{2}$ of all bits in full Kindle are in high-energy state
 - High-energy state about 10^{-15} joules higher
 - Then: Full Kindle is 1 attogram (10^{-18} gram) heavier
(Using $E = mc^2$)
- Of course, this is less than most sensitive scale can measure (it can measure 10^{-9} grams)
- Of course, this weight difference overwhelmed by battery discharge, weight from getting warm,
- Source: John Kubiatowicz (New York Times, Oct 24, 2011)

SSD Summary

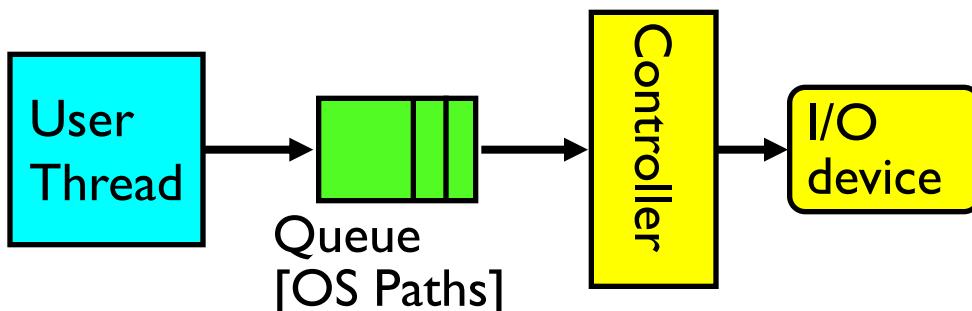
- Pros (vs. hard disk drives):
 - Low latency, high throughput (eliminate seek/rotational delay)
 - No moving parts:
 - » Very light weight, low power, silent, very shock insensitive
 - Read at memory speeds (limited by controller and I/O bus)
- Cons
 - Small storage (0.1-0.5x disk), expensive (3-20x disk)
 - » Hybrid alternative: combine small SSD with large HDD

SSD Summary

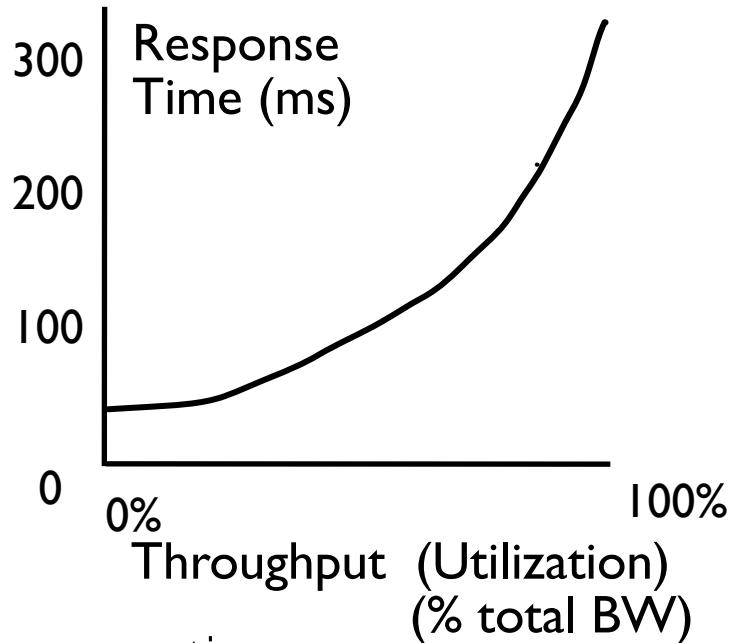
- Pros (vs. hard disk drives):
 - Low latency, high throughput (eliminate seek/rotational delay)
 - No moving parts:
 - » Very light weight, low power, silent, very shock insensitive
 - Read at memory speeds (limited by controller and I/O bus)
- Cons
 - ~~Small storage (0.1-0.5x disk), expensive (3-20x disk)~~

No longer true!
 - » Hybrid alternative: combine small SSD with large HDD
 - Asymmetric block write performance: read pg/erase/write pg
 - » Controller garbage collection (GC) algorithms have major effect on performance
 - Limited drive lifetime
 - » 1-10K writes/page for MLC NAND
 - » Avg failure rate is 6 years, life expectancy is 9–11 years
- These are changing rapidly!

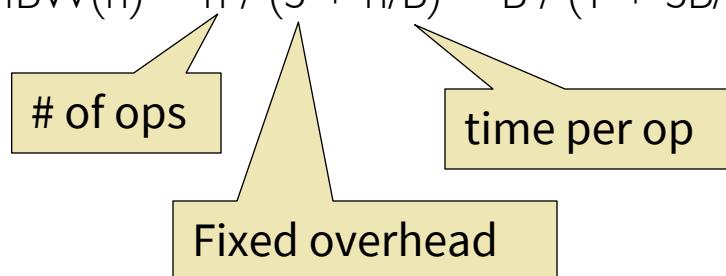
I/O Performance



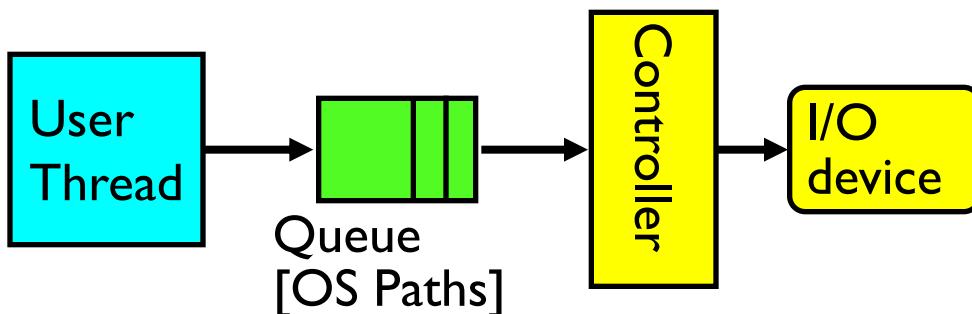
Response Time = Queue + I/O device service time



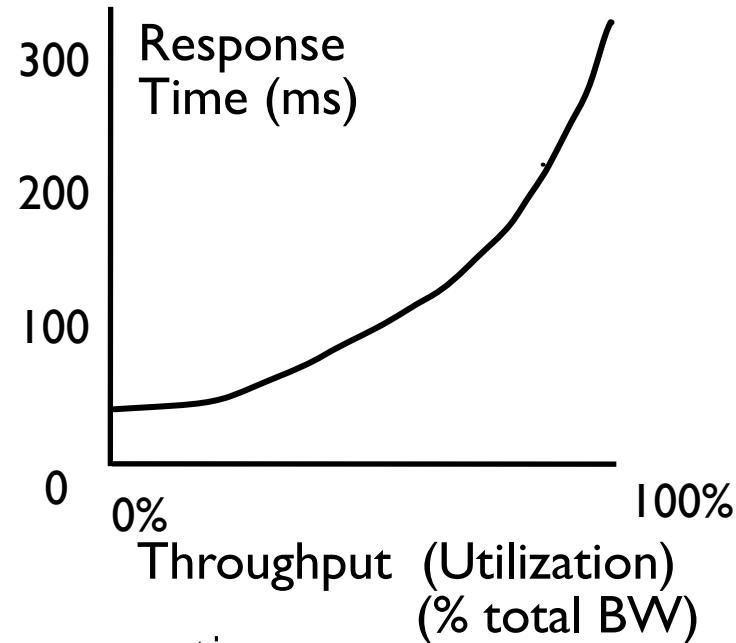
- Performance of I/O subsystem
 - Metrics: Response Time, Throughput
 - Effective BW per op = transfer size / response time
 - » $\text{EffBW}(n) = n / (S + n/B) = B / (I + SB/n)$



I/O Performance



Response Time = Queue + I/O device service time



- Performance of I/O subsystem
 - Metrics: Response Time, Throughput
 - Effective BW per op = transfer size / response time
 - » $\text{EffBW}(n) = n / (S + n/B) = B / (I + SB/n)$
 - Contributing factors to latency:
 - » Software paths (can be loosely modeled by a queue)
 - » Hardware controller
 - » I/O device service time
- Queuing behavior:
 - Can lead to big increases of latency as utilization increases
 - Solutions?

Summary (1/2)

- 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

Summary (2/2)

- Disk Performance:
 - Queuing time + Controller + Seek + Rotational + Transfer
 - Rotational latency: on average $\frac{1}{2}$ rotation
 - Transfer time: spec of disk depends on rotation speed and bit storage density
- Devices have complex interaction and performance characteristics
 - Response time (Latency) = Queue + Overhead + Transfer
 - » Effective BW = $BW * T/(S+T)$
 - HDD: Queuing time + controller + seek + rotation + transfer
 - SSD: Queuing time + controller + transfer (erasure & wear)
- Systems (e.g., file system) designed to optimize performance and reliability
 - Relative to performance characteristics of underlying device
- Bursts & High Utilization introduce queuing delays
- Queuing Latency:
 - M/M/I and M/G/I queues: simplest to analyze
 - As utilization approaches 100%, latency $\rightarrow \infty$