## Slide 1

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

## General I/O

March 28th, 2016

Prof. Anthony D. Joseph

http://cs162.eecs.Berkeley.edu

## Slide 2
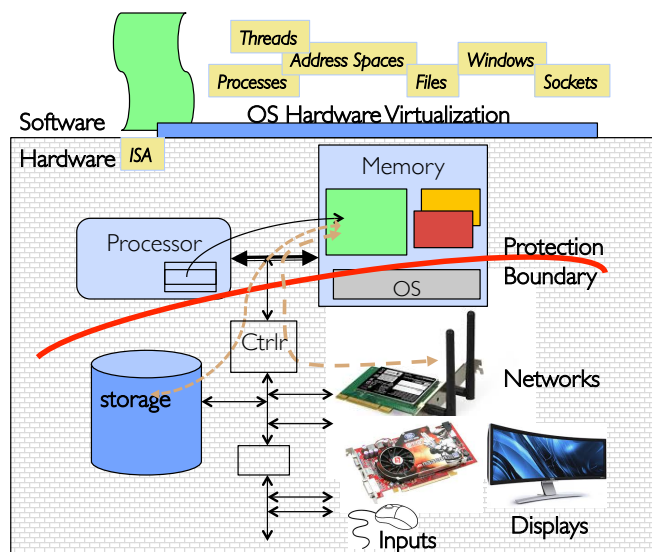
# The Requirements of I/O

- So far in this course:
  - We have learned how to manage CPU, 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?

## Slide 3

# OS Basics: I/O

## Slide 4

# 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
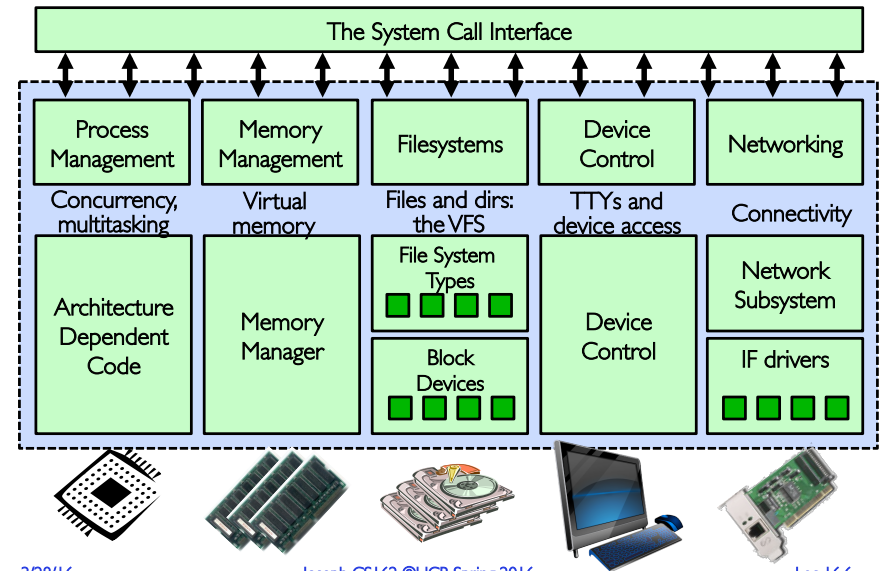
## Operational Parameters for I/O

- 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 sequential transfer (more later)

- Transfer Notification: Polling vs. Interrupts
  - Some devices require continual monitoring
  - Others generate interrupts when they need service

- Transfer Mechanism: Programmed IO and DMA

## Kernel Device Structure

## 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

- Project 2 initial design doc due today

## BREAK

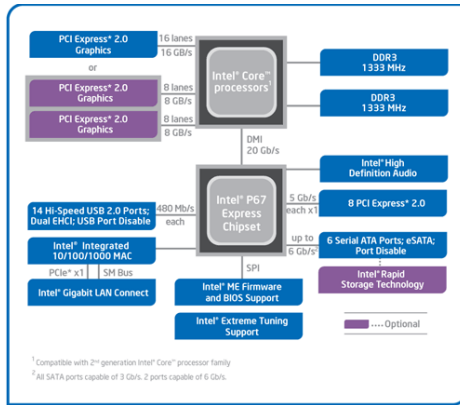## Chip-scale Features of Recent x86 (SandyBridge)



- Significant pieces:
  - Four OOO cores
    - » New Advanced Vector eXtensions (256-bit FP)
    - » Special purpose instructions: AES, Galois-Field mult
    - » 4 µ-ops/cycle
  - Integrated GPU, System Agent (Mem, Fast I/O)
  - Shared L3 cache divided in 4 banks
  - On-chip Ring bus network
    - » High-BW access to L3 Cache
- Integrated I/O
  - Integrated memory controller (IMC)
    - » Two independent channels of DDR3 DRAM
  - High-speed PCI-Express (for Graphics cards)
  - DMI Connection to SouthBridge (PCH)
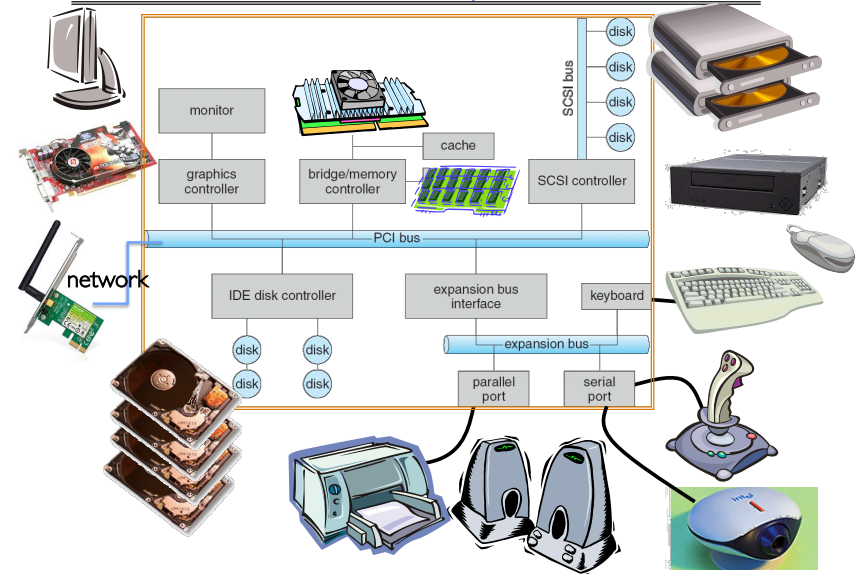
## SandyBridge I/O: PCH



SandyBridge
System Configuration

- Platform Controller Hub
  - Used to be "SouthBridge," but no "NorthBridge" now
  - Connected to processor with proprietary bus
    » Direct Media Interface
- Types of I/O on PCH:
  - USB, Ethernet
  - Audio, BIOS support
  - More PCI Express (lower speed than on Processor)
  - SATA (for Disks)

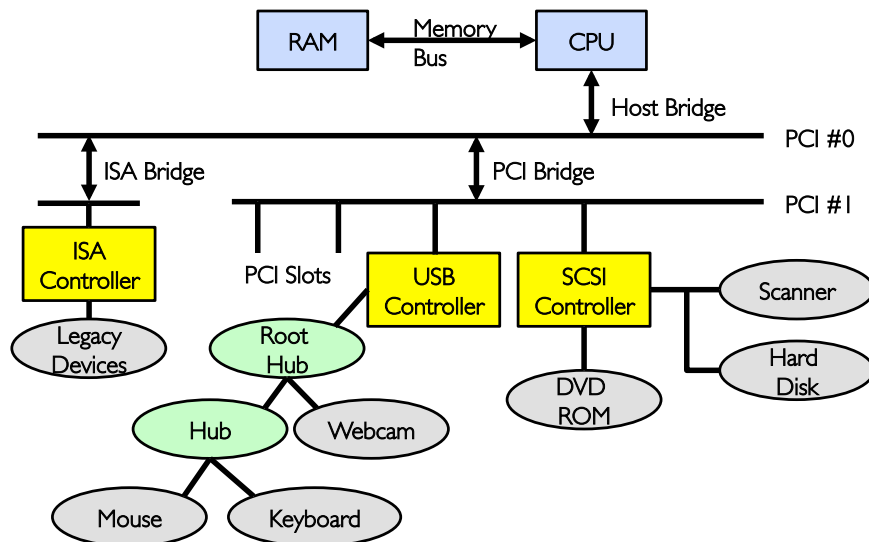## Modern I/O Systems

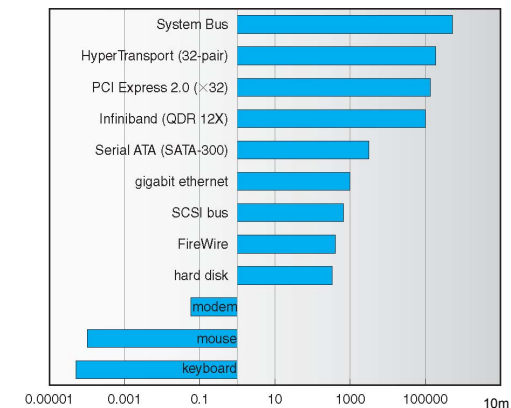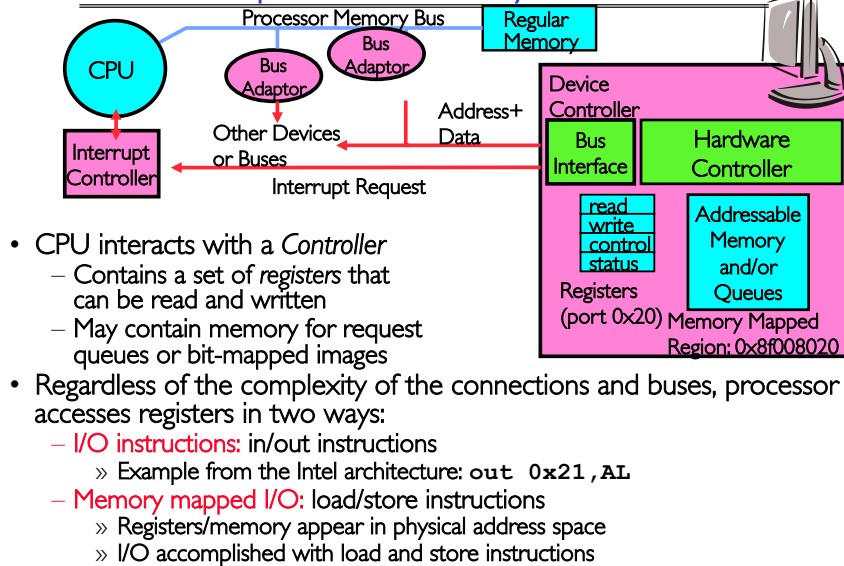## Example: PCI Architecture

## 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
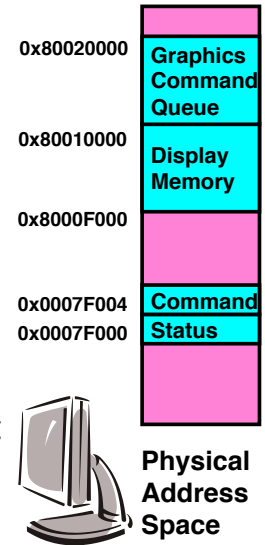
## How does the processor actually talk to the device?



- 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
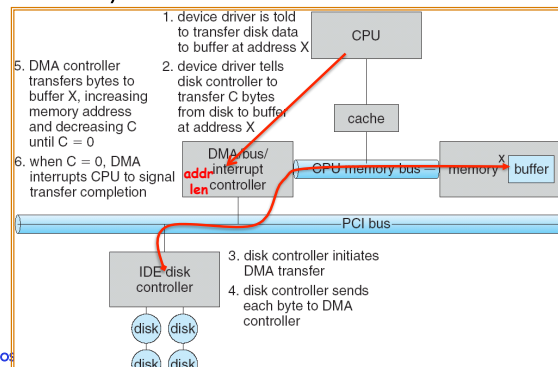- Can protect with address translation

## 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):

## 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
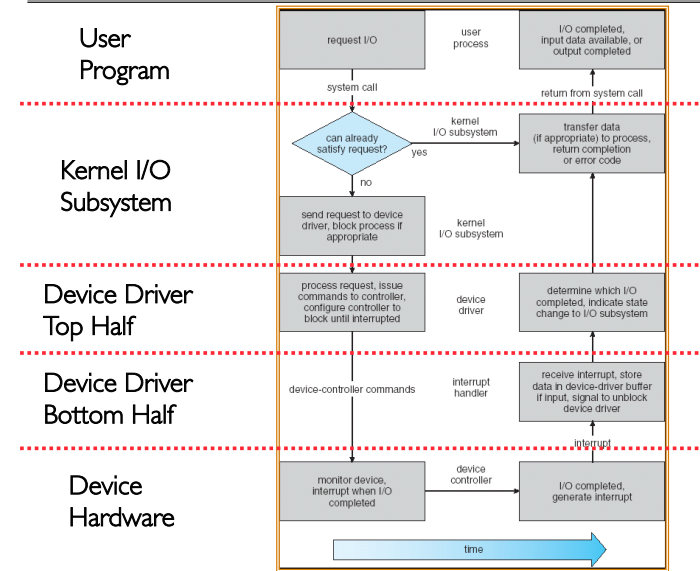
## Device Drivers

- **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

## Basic Performance Concepts

- *Response Time* or *Latency*: Time to perform an operation(s)

- *Bandwidth* or *Throughput*: Rate at which operations are performed (op/s)
  - Files: mB/s, Networks: mb/s, Arithmetic: GFLOP/s

- *Start up* or "Overhead": time to initiate an operation

- Most I/O operations are roughly linear
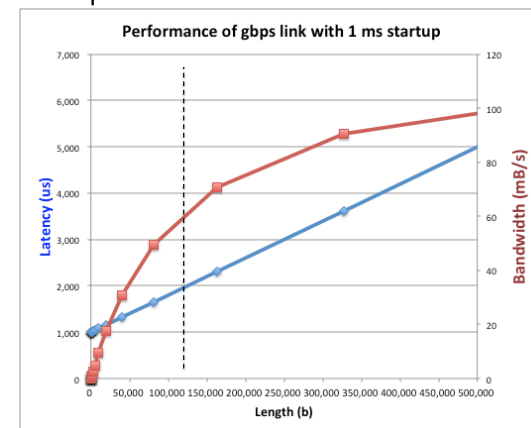  - Latency(n) = Overhead + n/Bandwidth

## Example (fast network)

- Consider a gpbs link (125 MB/s)
  - With a startup cost S = 1 ms



- Theorem: half-power point occurs at n=S*B:
  - When transfer time = startup $T(S*B) = S + S*B/B$

## Example: at 10 ms startup (like Disk)

**Performance of gbps link with 10 ms startup**



Chart axes: Latency (us) [left, 0–18,000], Bandwidth (mB/s) [right, 0–50], Length (b) [x-axis, 0–500,000]

## What Determines Peak BW for I/O ?

- Bus Speed
  - PCI-X: 1064 MB/s = 133 MHz × 64 bit (per lane)
  - ULTRA WIDE SCSI: 40 MB/s
  - Serial Attached SCSI & Serial ATA & IEEE 1394 (firewire) : 1.6 Gbps full duplex (200 MB/s)
  - USB 3.0 – 5 gb/s

- Device Transfer Bandwidth
  - Rotational speed of disk
  - Write / Read rate of NAND flash
  - Signaling rate of network link

- Whatever is the bottleneck in the path

## BREAK

## Storage Devices

- Magnetic disks
  - Storage that rarely becomes corrupted
  - Large capacity at low cost
  - Block level random access (except for SMR – later!)
  - Slow performance for random access
  - Better performance for streaming access

- Flash memory
  - Storage that rarely becomes corrupted
  - Capacity at intermediate cost (50x disk ???)
  - Block level random access
  - Good performance for reads; worse for random writes
  - Erasure requirement in large blocks
  - Wear patterns issue

## Are We in an Inflection Point?



An Accelerating Trend towards PC SSD

Notebook PC: SSD vs HDD Price

SSD vs HDD

Usually 10 000 or 15 000 rpm SAS drives

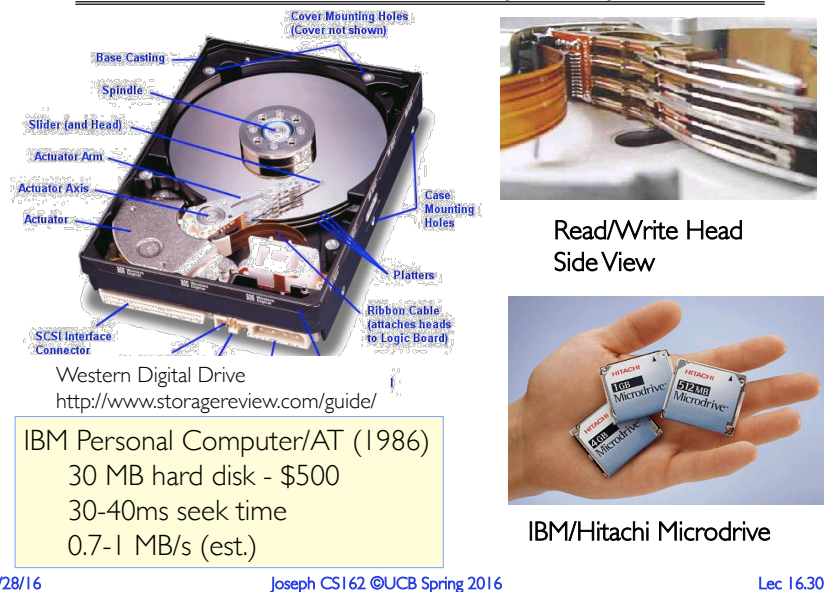| | | |
|---|---|---|
| **0.1** ms | **Access times** SSDs exhibit virtually no access time | **5.5 ~ 8.0** ms |
| SSDs deliver at least **6000** io/s SSDs are at least 15 times faster than HDDs | **Random I/O Performance** | HDDs reach up to **400** io/s |
| SSDs have a failure rate of less than **0.5** % | **Reliability** This means SSDs 4 - 10 times more reliable | HDD's failure rate fluctuates between **2 & 5** % |
| SSDs consume between **2 & 5** watts | **Energy savings** This means that on a large server like ours, approximately 100 watts are saved | HDDs consume between **6 & 15** watts |
| SSDs have an average I/O wait of **1** % You will have an extra 6% of CPU power for other operations | **CPU Power** | HDDs' average I/O wait is about **7** % |
| the average service time for an I/O request while running a database remains below **20** ms | **Input/Output request times** SSDs allow for much faster data access | the I/O request time with HDDs during backup rises up to **400~500** ms |
| SSD backups take about **6** hours SSDs allows for 3 - 5 times faster backups for your data | **Backup Rates** | HDD backups take up to **20~24** hours |

---

## Hard Disk Drives (HDDs)



Western Digital Drive
http://www.storagereview.com/guide/

Read/Write Head Side View

IBM/Hitachi Microdrive

IBM Personal Computer/AT (1986)
  30 MB hard disk - $500
  30-40ms seek time
  0.7-1 MB/s (est.)

---

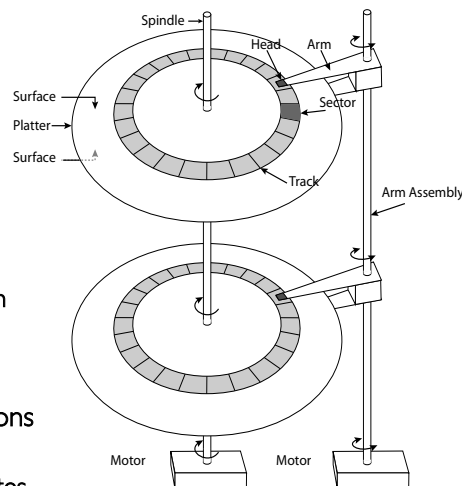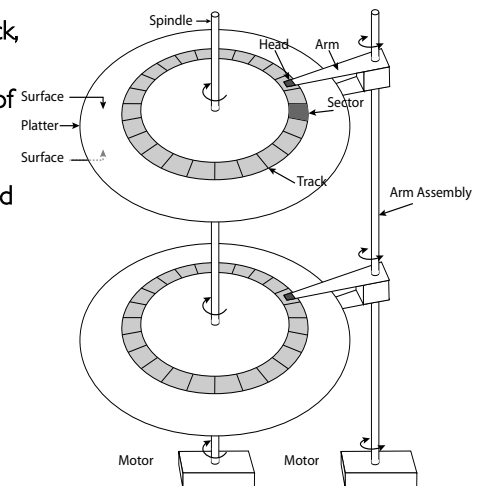## The Amazing Magnetic Disk

- Unit of Transfer: Sector
  - Ring of sectors form a track
  - Stack of tracks form a cylinder
  - Heads position on cylinders

- Disk Tracks ~ 1μm (micron) wide
  - Wavelength of light is ~ 0.5μm
  - Resolution of human eye: 50μm
  - 100K 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)
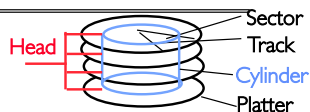
---

## 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

- New: Shingled Magnetic Recording (SMR)
  - Overlapping tracks ⇒ greater density, restrictions on writing
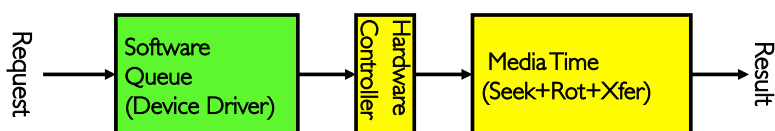  - Seagate (8TB), Hitachi (10TB)

## Magnetic Disk Characteristic

- Cylinder: all the tracks under the head at a given point on all surfaces
- Read/write: three-stage process:
  - Seek time: position head/arm over proper track (into proper cylinder)
  - Rotational latency: wait for desired sector to rotate under R/W head
  - Transfer time: transfer a block of bits (sector) under the R/W head
- Disk Latency = Queuing Time + Controller time + Seek Time + Rotation Time + Xfer Time

Request → Software Queue (Device Driver) → Hardware Controller → Media Time (Seek+Rot+Xfer) → Result

- Highest Bandwidth:
  - Transfer large group of blocks sequentially from one track

## Typical Numbers for Magnetic Disk

| Parameter | Info / Range |
|---|---|
| Space/Density | Space: 8TB (Seagate), 10TB (Hitachi) in 3½ inch form factor! Areal Density: ≥ 1 Terabit/square inch! (SMR, Helium, …) |
| Average seek time | Typically 5-10 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 100 MB/s. Depends on: <br> • Transfer size (usually a sector): 512B – 1KB per sector <br> • Rotation speed: 3600 RPM to 15000 RPM <br> • Recording density: bits per inch on a track <br> • Diameter: ranges from 1 in to 5.25 in |
| Cost | Drops by a factor of two every 1.5 years (or even faster). $0.03-0.07/GB in 2013 |

## Disk Performance Example

- Assumptions:
  - Ignoring queuing and controller times for now
  - Avg seek time of 5ms,
  - 7200RPM ⇒ Time for rotation: 60000(ms/M)/7200(rev/M) ~= 8ms
  - Transfer rate of 4MByte/s, sector size of 1 Kbyte ⇒ 1024 bytes/$4 \times 10^6$ (bytes/s) = $256 \times 10^{-6}$ sec ≈ .26 ms
- Read sector from random place on disk:
  - Seek (5ms) + Rot. Delay (4ms) + Transfer (0.26ms)
  - *Approx* 10ms to fetch/put data: 100 KByte/sec
- Read sector from random place in same cylinder:
  - Rot. Delay (4ms) + Transfer (0.26ms)
  - *Approx* 5ms to fetch/put data: 200 KByte/sec
- Read next sector on same track:
  - Transfer (0.26ms): 4 MByte/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

- …

# 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
  - Report their results through either interrupts or a status register that processor looks at occasionally (polling)
- Notification mechanisms
  - Interrupts
  - Polling: Report results through status register that processor looks at periodically
- Drivers interface to I/O devices

  - Provide clean Read/Write interface to OS above

  - Manipulate devices through PIO, DMA & interrupt handling

  - 2 types: block, character, and network